# **EIDER: Evidence-enhanced Document-level Relation Extraction**

**Anonymous ACL submission** 

#### Abstract

Document-level relation extraction (DocRE) 002 aims to extract the semantic relations among entity pairs in a document. Typical DocRE methods blindly take the full document as input, while a subset of the sentences in the document, noted as the evidence, are often sufficient for humans to predict the relation of an entity pair. In this paper, we propose an evidence-enhanced framework EIDER that automatically extracts and utilizes evidence for DocRE. We first train an evidence extraction model together with relation extraction via multi-task learning, which allows the two tasks to benefit from shared representations and improve each other. We show that 016 even in the absence of human annotated evidence, using silver evidence labels extracted by heuristic rules still leads to better RE performance. We further design a simple yet effective evidence-enhanced inference process that makes RE predictions on both extracted evidence and the full document, then fuses the predictions through a blending layer. This allows EIDER to focus on important sentences while still having access to all the information in the document. Extensive experiments show that EIDER outperforms state-of-the-art methods on three benchmark datasets, e.g., by 1.37/1.26 Ign F1/F1 on DocRED.

#### 1 Introduction

011

017

022

024

034

040

041

Relation extraction (RE) is the task of extracting semantic relations among entities within a given text, which has abundant applications such as knowledge graph construction, question answering, and biomedical text analysis (Yu et al., 2017; Shi et al., 2019; Trisedya et al., 2019). Prior studies mostly focus on predicting the relation between two entity mentions in a single sentence. However, in reality, an entity may have multiple mentions throughout a document. It is also common that a relation can only be inferred given multiple sentences as the

Head: Hero of the Day Tail: the United States Relation: [country of origin]
Ground truth evidence sentences: [1,10] Extracted evidence: [1,10]
Original document as input: [1] Load is the sixth studio album by the
American heavy metal band Metallica, released on June 4, 1996 by
Elektra Records in the United States and by Vertigo Records
internationally [9] It was certified 5×platinum by the Recording
Industry Association of America ( RIAA ) for shipping five million
copies in the United States. [10] Four singles-"Until It Sleeps",
"Hero of the Day", "Mama Said", and "King Nothing" - were released
as part of the marketing campaign for the album.
Prediction result (logits): NA: 17.63 country of origin: 14.79
Extracted evidence as input: [1] Load is the sixth studio album in
the United States and by Vertigo Records internationally. [10] Four
singles —"Until It Sleeps", "Hero of the Day", for the album.
Prediction result (logits): country of origin: 18.31 NA: 13.45
Final prediction of our model: country of origin

Figure 1: A test sample in the DocRED dataset (Yao et al., 2019), where the  $i^{th}$  sentence in the document is marked with [i] at the start. Our model correctly predicts [1,10] as evidence, and if we only use the extracted evidence as input, the model can predict the relation "country of origin" correctly.

context. As a result, recent studies have been moving towards the more realistic setting of documentlevel relation extraction (DocRE) (Peng et al., 2017; Yao et al., 2019; Zeng et al., 2020).

042

043

044

046

052

054

056

058

060

061

062

063

064

Unlike typical DocRE models that blindly take the whole document as input, a human may only need a few sentences to infer the relation of an entity pair. For each entity pair, we define the minimal set of sentences required by human annotators to infer their relation as their *evidence sentences*. As shown in Figure 1, to predict the relation between "Hero of the Day" and "the United States", it is sufficient to know that *Load (the album)* was released in the United States from the  $1^{st}$  sentence, and "Hero of the Day" is a single of Load from the  $10^{th}$  sentence. In other words, the  $1^{st}$  and  $10^{th}$ sentences serve as the evidence to infer this relation. Although the  $9^{th}$  sentence also mentions "the United States", it is irrelevant to this specific relation. Including such irrelevant sentences in input might sometimes introduce noise to the model and be more detrimental than beneficial. Despite the usefulness of evidence, previous studies regarding

065 066 067

069

072

073

074

075

077

081

087

095

100

101

102

104

105

106

107

109

110

111

112

113

114

115

evidence either require expensive human-annotated evidence sentences (Huang et al., 2021a) or fail to show improvements when paired up with pretrained language models (Huang et al., 2021b).

In this paper, we propose an evidence-enhanced DocRE framework EIDER, which automatically extracts evidence and effectively leverages the extracted evidence to improve DocRE without extensive human annotation. During training, we enhance DocRE by jointly extracting relations and evidence using multi-task learning. Intuitively, both relation extraction and evidence extraction should focus on the information relevant to the current entity pair, such as the underlined "Load" and "the album" in Figure 1. This suggests that the two tasks have certain commonalities and can provide additional training signals for each other. Experimental results show that these two tasks can mutually enhance each other. One remaining issue is that human-annotated evidence sentences are not always available and heavily relying on them may limit model applicability. To reduce the need for evidence annotation, we design several heuristic rules to construct silver labels if evidence annotation is unavailable. We observe that EIDER still improves RE performance even trained with our silver labels, and sometimes even performs on par with using gold labels.

With the evidence extracted, either by rules or our evidence extraction model, we further enhance DocRE by leveraging the evidence in inference. In the extreme case, if there is only one sentence related to the relation, one can make predictions solely based on this sentence and reduce the problem to sentence-level relation extraction. One naive approach is thus to directly replace the original document with the extracted evidence. However, since no system can extract evidence perfectly, solely relying on extracted sentences may miss important information and harm model performance in certain cases (see Table 5). To avoid information loss, we fuse the prediction results of the original document and extracted evidence through a blending layer (Wolpert, 1992). In this way, EIDER pays more attention to the extracted important sentences, while still having access to all the information in the document. Empirical analysis demonstrates that removing either source would lead to degenerate performance.

We conduct extensive experiments on three widely-adopted DocRE benchmarks: DocRED

(Yao et al., 2019), CDR (Li et al., 2016) and GDA (Wu et al., 2019). Experiment results show that EIDER achieves state-of-the-art performance on all the datasets. Performance analysis further shows that the improvement of EIDER is most significant on inter-sentence entity pairs, suggesting that leveraging evidence is especially effective in reasoning over multiple sentences. In particular, EIDER significantly improves the performance on entity pairs that require co-reference/multi-hop reasoning by 1.98/2.08 F1 on DocRED, respectively.

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

**Contributions**. (1) We propose a joint relation and evidence extraction model that allows the two tasks to mutually enhance each other without heavily relying on evidence annotation. (2) We design a simple and effective DocRE inference process enhanced by the extracted evidence, enabling more focus on the important sentences with no information loss. (3) We demonstrate that our evidence-enhanced framework outperforms state-of-the-art methods on three DocRE datasets.

# 2 Problem Formulation

Given a document d comprised of N sentences  $\{s_n\}_{n=1}^N$ , L tokens  $\{h_l\}_{l=1}^{\overline{L}}$ , E named entities  $\{e_i\}_{i=1}^{E}$  and all the proper-noun mentions of each entity,  $\{m_i^i\}$ , the task of document-level relation extraction (DocRE) is to predict the set of all possible relations between all entity pairs  $(e_h, e_t)$  from a pre-defined relation set  $\mathcal{R} \mid J\{NA\}$ . We refer to  $e_h$ and  $e_t$  as the head entity and tail entity, respectively. A relation r belongs to the positive class  $\mathcal{P}_{h,t}^T$  if it exists between  $(e_h, e_t)$  and otherwise the negative class  $\mathcal{N}_{h,t}^T$ . For each entity pair  $(e_h, e_t)$  that possesses a non-NA relation, we define its  $evidence^1$  $V_{h,t} = \{s_{v_k}\}_{k=1}^K$  as the subset of sentences in the document that are sufficient for human annotators to infer the relation. Human annotation of evidence may or may not be given in training, depending on the datasets, but is not available in inference.

# 3 Methodology

An illustration of the framework of EIDER is shown in Figure 2. In training, we jointly extract relation and evidence using multi-task learning, where the two tasks have their own classifier and share the base encoder (Sec. 3.1). In inference, we fuse the predictions on the original document and the extracted evidence using a blending layer (Sec. 3.2).

<sup>&</sup>lt;sup>1</sup>We use "*evidence sentence*" and "*evidence*" interchangeably throughout the paper.



Figure 2: The overall architecture of EIDER. The left part illustrates the first stage (training) and the right shows the second and third stages (inference) of EIDER. We highlight head entities, tail entities and extracted evidences.

In case the evidence annotation is not available, we also provide several heuristic rules to construct silver evidence labels as an alternative (Sec. 3.3).

163

164

165

166

167

168

169

170

171

172

175

176

177

178

179

183

184

187

190

191

192

193

194

#### Joint Relation and Evidence Extraction 3.1

In our framework, we jointly train the relation extraction model with an evidence extraction model using multi-task learning. As shown in Figure 2, the two tasks have their own classifier but share the base encoder. Intuitively, tokens relevant to predicting the relation are essential in both models. By sharing the base encoder, the two tasks can provide additional training signals for each other and hence mutually enhance each other (Ruder, 2017).

**Base Encoder**. We leverage pre-trained language models (Devlin et al., 2019) to encode the semantic meanings of each token in the document. Specifically, given a document  $d = [h_l]_{l=1}^L$ , we insert a special token "\*" before and after each entity mention  $\{m_i^i\}$  and leverage the encoder to obtain the sdim token embeddings  $\boldsymbol{H} = [\mathbf{h}_1, ..., \mathbf{h}_L], \mathbf{h}_l \in \mathbb{R}^s$ and the cross token attention  $A \in \mathbb{R}^{L \times L}$ :

$$\boldsymbol{H}, \boldsymbol{A} = \text{Encoder}([h_1, ..., h_L]), \quad (1)$$

where A is the average of the attention heads in the last transformer layer (Vaswani et al., 2017). For each mention of an entity  $e_i$ , we use the embedding of the start symbol "\*" as its mention embedding  $\mathbf{m}_{i}^{i}$ . Then, we obtain the embedding of entity  $e_{i}$ by adopting LogSumExp pooling (Jia et al., 2019; Zhou et al., 2021) over the embeddings of all its mentions:  $\mathbf{e}_i = \log \sum_j \exp(\mathbf{m}_j^i)$ .

To predict the relation of different entity pairs, a model may need to focus on different parts of the context. To capture the context relevant to each entity pair  $(e_h, e_t)$ , we compute its context embedding  $\mathbf{c}_{h,t} \in \mathbb{R}^s$  based on the attention matrix Afrom the pre-trained encoder (Zhou et al., 2021):

$$c_{h,t} = \boldsymbol{H}^T \frac{\boldsymbol{A}_h \circ \boldsymbol{A}_t}{\boldsymbol{A}_h^T \boldsymbol{A}_t},$$
 (2)

195

197

199

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

where  $\circ$  is the Hadamard product and  $A_h \in \mathbb{R}^L$ is  $e_h$ 's attention to all the tokens in the document, obtained by averaging  $e_h$ 's mention-level attention. Similarly for  $A_t$ . The intuition is that tokens with high attention towards both  $e_h$  and  $e_t$  are important to both entities. Hence, these tokens are likely to be essential to the relation and should contribute more to the context embedding.

Relation Classifier. To predict the relation between an entity pair  $(e_h, e_t)$ , we first compute their context-aware representations  $(\mathbf{z}_{\mathbf{h}}, \mathbf{z}_{\mathbf{t}})$  by combining their entity embeddings  $(\mathbf{e}_h, \mathbf{e}_t)$  with their context embedding  $\mathbf{c}_{h,t}$  and then utilize a bilinear function to calculate the logit of how likely a relation  $r \in \mathcal{R}$  exists between  $e_h$  and  $e_t$ :

$$z_{h} = \tanh \left( W_{h} \mathbf{e}_{h} + W_{c_{h}} \mathbf{c}_{h,t} \right),$$

$$z_{t} = \tanh \left( W_{t} \mathbf{e}_{t} + W_{c_{t}} \mathbf{c}_{h,t} \right), \quad (3)$$

$$\mathbf{y}_{r} = z_{h} W_{r} z_{t} + \boldsymbol{b}_{r},$$

$$\boldsymbol{x} = \boldsymbol{z}_h \boldsymbol{W}_r \boldsymbol{z}_t + \boldsymbol{b}_r,$$

where  $W_h, W_t, W_{c_h}, W_{c_t}, W_r$  and  $b_r$  are learnable parameters. As the model may have different confidence for different entity pairs or classes, we apply the adaptive-thresholding loss (Zhou et al., 2021), which learns a dummy relation class TH that serves as the dynamic threshold for each entity pair:

$$\mathbf{y}_{\mathrm{TH}} = \boldsymbol{z}_h \boldsymbol{W}_{\mathrm{TH}} \boldsymbol{z}_t + \boldsymbol{b}_r. \tag{4}$$

224 During inference, for each tuple  $(e_h, e_t, r), r \in \mathcal{R}$ , 225 we obtain the prediction score:  $S_{h,t,r}^{(O)} = \mathbf{y}_r - \mathbf{y}_{TH}$ . 226 Finally, we define our training objective for relation 227 extraction as follows:

228

229

232

239

240

241

242

244

245

246

247

252

256

257

260

261

262

264

$$\mathcal{L}_{RE} = -\sum_{h \neq t} \sum_{r \in \mathcal{P}_{h,t}^{T}} \log \left( \frac{\exp\left(\mathbf{y}_{r}\right)}{\sum_{r' \in \mathcal{P}_{h,t}^{T} \cup \{\text{TH}\}} \exp\left(\mathbf{y}_{r'}\right)} \right) - \log \left( \frac{\exp\left(\mathbf{y}_{\text{TH}}\right)}{\sum_{r' \in \mathcal{N}_{h,t}^{T} \cup \{\text{TH}\}} \exp\left(\mathbf{y}_{r'}\right)} \right).$$
(5)

**Evidence Classifier**. In addition to the relation, we also predict whether each sentence  $s_n$  is an evidence sentence of entity pair  $(e_h, e_t)$ . Similar to entity embeddings, to obtain sentence embedding  $\mathbf{s}_n$ , we apply a LogSumExp pooling over all the tokens in  $s_n$ :  $\mathbf{s}_n = \log \sum_{h_l \in s_n} \exp(\mathbf{h}_l)$ . Intuitively, if  $s_n$  is an evidence sentence of  $(e_h, e_t)$ , the tokens in  $s_n$  would be relevant to the relation prediction, and should contribute more to  $\mathbf{c}_{h,t}$ . Hence, we use a bilinear function between context embedding  $\mathbf{c}_{h,t}$ and sentence embedding  $\mathbf{s}_n$  to measure the importance of sentence  $s_n$  to entity pair  $(e_h, e_t)$ :

$$P(s_n|e_h, e_t) = \sigma\left(\mathbf{s}_n \boldsymbol{W}_v \mathbf{c}_{h,t} + \boldsymbol{b}_v\right), \quad (6)$$

where  $W_v$  and  $b_v$  are learnable parameters.

As an entity pair may have more than one evidence sentence, we use the binary cross entropy as the objective to train the evidence extraction model.

$$\mathcal{L}_{Evi} = -\sum_{h \neq t, \mathbf{NA} \notin \mathcal{P}_{h,t}^T} \sum_{s_n \in \mathcal{D}} y_n \cdot \mathbf{P}\left(s_n | e_h, e_t\right) + (1 - y_n) \cdot \log(1 - \mathbf{P}\left(s_n | e_h, e_t\right)), \quad (7)$$

where the evidence label  $y_n$  is 1 when  $s_n \in V_{h,t}$ and otherwise 0. If golden labels are not provided, we use several heuristic rules to construct silver labels instead. Details are introduced in Sec 3.3.

Although it is possible that different relations possessed by  $(e_h, e_t)$  are inferred from different sentences, we observe that most entity pairs only have one set of evidence across relations. We thus only predict the evidence for each  $(e_h, e_t)$  pair instead of  $(e_h, e_t, r)$  tuple. This largely reduces the memory and run time of our method, especially when  $|\mathcal{R}|$  is large (e.g.,  $|\mathcal{R}| = 96$  in DocRED).

Note that only entity pairs with  $r \in \mathcal{R}$  have human-annotated evidence sentences. For entity pairs with r = NA, a naive way is to simply regard their evidence label as the empty set (Huang et al., 2021a). However, there may still exist some implicit relations out of the pre-defined set  $\mathcal{R}$ , which can be inferred from certain sentences. Regarding every sentence as non-evidence may not be reasonable. As a result, we only train the evidence extraction model on entity pairs with at least one non-NA relation  $r \in \mathcal{R}$ , which accounts for a small subset (e.g., 2.97% in DocRED) of the total possible entity pairs. This is another reason why our model is efficient in both memory and training time.

Finally, we optimize our model by the combination of the relation extraction loss  $\mathcal{L}_{RE}$  and evidence extraction loss  $\mathcal{L}_{Evi}$ :

$$\mathcal{L} = \mathcal{L}_{RE} + \mathcal{L}_{Evi}.$$
 (8)

265

266

267

269

270

271

272

273

274

275

276

277

278

279

281

282

283

284

286

287

289

290

291

292

294

297

299

300

301

302

303

304

305

306

307

309

#### 3.2 Evidence-enhanced Inference

Suppose our extracted evidence sentences already contain all the information relevant to the relation, then there is no need to use the whole document for relation extraction. However, no system can perfectly extract the evidence without missing any sentences. Solely relying on the extracted evidence may miss important information in the document and lead to sub-optimal performance. Therefore, we combine the prediction results on both the original document and the extracted evidence, which can either be learned by our evidence classifier (Sec. 3.1) or constructed by our heuristic rules (Sec. 3.3) if evidence annotation is unavailable.

Specifically, as shown in Figure 2, we first obtain a set of relation prediction scores  $S_{h,t,r}^{(O)}$  from the original documents. Then we construct a pseudo document  $d'_{h,t}$  for each entity pair by concatenating the extracted evidence sentences  $V'_{h,t}$  in the order they are presented in the original document and feed it into the relation classifier to obtain another set of prediction scores  $S_{h,t,r}^{(E)}$ . Finally, we fuse the results by aggregating the two sets of prediction sores through a blending layer (Wolpert, 1992):

$$P_{Fuse}(r|e_h, e_t) = \sigma(S_{h,t,r}^{(O)} + S_{h,t,r}^{(E)} - \tau).$$
(9)

We choose this design because it is simple and only includes one learnable parameter,  $\tau$ , alleviating over-fitting in the development set. We optimize the parameter  $\tau$  on the development set as follows:

$$\mathcal{L}_{Fuse} = -\sum_{d \in \mathcal{D}} \sum_{h \neq t} \sum_{r \in \mathcal{R}} y_r \cdot \mathcal{P}_{Fuse} \left( r | e_h, e_t \right) + (1 - y_r) \cdot \log(1 - \mathcal{P}_{Fuse} \left( r | e_h, e_t \right)), \quad (10)$$

where  $y_r = 1$  if relation r holds between  $(e_h, e_t)$ and  $y_r = 0$  otherwise. Empirically, using other loss functions does not affect the performance much.

# 310

#### 3.3 Heuristic Evidence Label Construction

In case human annotation of evidence is not available, we design a set of heuristic rules to automatically construct silver labels for evidence extraction. Then we train our joint model on the silver labels and directly use the silver labels as pseudo documents in inference. The percentage of test samples covered by each rule is shown in Table 6.

Co-occur. If the head and tail entities co-occur in the same sentence (e.g., "Load" and "the United States" co-occur in the  $1^{st}$  sentence in Figure 2), we use all the sentences they co-occur as evidence. 321 Coref. If the proper-noun mentions of the head and tail entity do not co-occur, but their coreferential mentions co-occur (e.g., "Hero of the Day" and "the album", the co-reference of "Load" co-occur 325 in the  $10^{th}$  sentence in Figure 2), we use all the sentences where their coreferential mentions co-327 occur as evidence. In practice, we directly apply a pre-trained coreference resolution model, HOI (Xu and Choi, 2020), without fine-tuning on our dataset. 330 Bridge. If the first two conditions are not met, but there exists a third bridge entity whose coreferential 332 mention co-occurs with both head and tail (e.g., "Load" or its coreferential mention "the album" co-334 occurs with both "the United States" and "Hero of the Day" in Figure 2), we take all the sentences 336 where the bridge co-occurs with head or tail as the evidence. If there is more than one bridge entity, we choose the one with the highest frequency. While 339 this rule can be easily extended to multiple bridges, we empirically observe that capturing one bridge 341 already leads to satisfying results.

# 4 Experiments

343

346

348

## 4.1 Experiment Setup

**Datasets**. We evaluate the effectiveness of EI-DER on three datasets: DocRED (Yao et al., 2019), CDR (Li et al., 2016) and GDA (Wu et al., 2019), where DocRED is the only dataset that provides evidence labels as part of the annotation. The details of the datasets are listed in Appendix A.1.

Evaluation Metrics. Following prior studies (Yao et al., 2019), we use F1 and Ign F1 as the main evaluation metrics for relation extraction, where Ign F1 measures the F1 score excluding the relations shared by the training and development/test set. We also report Intra F1 and Inter F1, where the former measures the performance on the co-occurred (intra-sentence) entity pairs and the latter evaluates

the inter-sentence entity pairs where none of their proper-noun mentions co-occurs. For evidence extraction, we compute the F1 score (denoted as **Evi F1**) and further introduce **PosEvi F1**, which measures the F1 score of evidence only on positive entity pairs (i.e., those with non-NA relations).

#### 4.2 Main Results

We compare our methods with both *Graph-based methods* and *transformer-based methods*. Graphbased methods explicitly perform inference on document-level graphs. Transformer-based methods, including EIDER, model cross-sentence relations by implicitly capturing the long-distance token dependencies via the transformer. We also compare to an ablation EIDER (Rule) that does not rely on human-annotated evidence and only uses silver evidence labels as described in Sec. 3.3. The implementation details are listed in Appendix A.2.

**Relation Extraction Results**. Table 1 and Table 2 present the relation extraction results, where we observe that EIDER outperforms the baseline methods in all datasets. Our improvement is especially large on Inter F1 (e.g., 1.21/2.01 Intra/Inter F1 under BERT<sub>base</sub> compared to ATLOP). We hypothesize that the bottleneck of inter-sentence pairs is to locate the relevant context, which often spreads through the whole document. EIDER learns to capture important sentences in training and focuses more on these important sentences in inference.

Among the baselines, the Inter F1 of GAIN is 0.70 higher while the Intra F1 of ATLOP is 0.16 higher, indicating that graph-based methods may capture the long-distance dependency between entities by directly connecting them on the graph. Although EIDER does not involve an explicit multihop reasoning module, it still notably outperforms the graph-based models in terms of Inter F1. This demonstrates that the evidence-enhanced inference also captures long-distance dependency by directly concatenating important sentences.

Finally, in both DocRED and the two biomedical datasets which do not have evidence annotation, EIDER (Rule) also outperforms all the baselines. This shows that EIDER does not heavily rely on evidence annotation. The improvement on DocRED and CDR is much larger than that on GDA. We hypothesize that it is because more than 85% relations in GDA are intra-sentence ones, where the evidence is normally the sentences where head and tail co-occur. After training on massive examples, 363 364

359

360

361

362

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

Model		D	ev		Te	st
hout	Ign F1	F1	Intra F1	Inter F1	Ign F1	F1
LSR-BERT <sub>base</sub> (Nan et al., 2020)	52.43	59.00	65.26	52.05	56.97	59.05
GLRE-BERT <sub>base</sub> (Wang et al., 2020)	-	-	-	-	55.40	57.40
Reconstruct-BERT <sub>base</sub> (Xu et al., 2021)	58.13	60.18	-	-	57.12	59.45
GAIN-BERT <sub>base</sub> (Zeng et al., 2020)	59.14	61.22	67.10	53.90	59.00	61.24
BERT <sub>base</sub> (Wang et al., 2019)	-	54.16	61.61	47.15	-	53.20
BERT-Two-Step (Wang et al., 2019)	-	54.42	61.80	47.28	-	53.92
HIN-BERT <sub>base</sub> (Tang et al., 2020)	54.29	56.31	-	-	53.70	55.60
E2GRE-BERT <sub>base</sub> (Huang et al., 2021a)	55.22	58.72	-	-	-	-
CorefBERT <sub>base</sub> (Ye et al., 2020)	55.32	57.51	-	-	54.54	56.96
ATLOP-BERT <sub>base</sub> (Zhou et al., 2021)	$59.11\pm0.14^{\dagger}$	$61.01\pm0.10^{\dagger}$	$67.26\pm0.15^{\dagger}$	$53.20\pm0.19^{\dagger}$	59.31	61.30
EIDER (Rule)-BERT <sub>base</sub>	$60.36\pm0.13$	$62.34\pm0.08$	$68.40\pm0.14$	$54.79\pm0.13$	60.23	62.21
EIDER-BERT <sub>base</sub>	$\textbf{60.51} \pm \textbf{0.11}$	$\textbf{62.48} \pm \textbf{0.13}$	$\textbf{68.47} \pm \textbf{0.08}$	$\textbf{55.21} \pm \textbf{0.21}$	60.42	62.47
BERT <sub>large</sub> (Ye et al., 2020)	56.67	58.83	-	-	56.47	58.69
CorefBERT <sub>large</sub> (Ye et al., 2020)	56.82	59.01	-	-	56.40	58.83
RoBERTalarge (Ye et al., 2020)	57.14	59.22	-	-	57.51	59.62
CorefRoBERTa <sub>large</sub> (Ye et al., 2020)	57.35	59.43	-	-	57.90	60.25
GAIN-BERT <sub>large</sub> (Zeng et al., 2020)	60.87	63.09	-	-	60.31	62.76
ATLOP-RoBERTa <sub>large</sub> (Zhou et al., 2021)	$61.30\pm0.22^{\dagger}$	$63.15\pm0.21^{\dagger}$	$69.61\pm0.25^{\dagger}$	$55.01\pm0.18^{\dagger}$	61.39	63.40
EIDER (Rule)-RoBERTa <sub>large</sub>	$61.73\pm0.07$	$63.91 \pm 0.07$	$69.99 \pm 0.09$	$56.27 \pm 0.11$	61.93	64.12
EIDER-RoBERTalarge	$\textbf{62.34} \pm \textbf{0.14}$	$\textbf{64.27} \pm \textbf{0.10}$	$\textbf{70.36} \pm \textbf{0.07}$	$\textbf{56.53} \pm \textbf{0.15}$	62.85	64.79

Table 1: Relation extraction results on DocRED. We report the mean and standard deviation on the development set by conducting 5 runs with different random seeds. We report the official test score of the best checkpoint on the development set. Results with † are based on our implementation. Others are reported in their original papers. We separate graph-based and transformer-based methods into two groups.

Model	CDR	GDA
LSR-BERT <sub>base</sub> (Nan et al., 2020)	64.8	82.2
SciBERT <sub>base</sub> (Zhou et al., 2021)	$65.1\pm0.6$	$82.5\pm0.3$
DHG-BERT <sub>base</sub> (Zhang et al., 2020b)	65.9	83.1
GLRE-SciBERT <sub>base</sub> (Wang et al., 2020)	68.5	-
ATLOP-SciBERT <sub>base</sub> (Zhou et al., 2021)	$69.4\pm1.1$	$83.9\pm0.2$
EIDER (Rule)-SciBERT <sub>base</sub>	$\textbf{70.63} \pm 0.49$	$\textbf{84.54} \pm 0.22$

Table 2: Relation extraction results on CDR and GDA.

it may be trivial even for the single RE model to focus on these sentences.

409

410

411

412

413

414

415

416

417

418

419

**Evidence Extraction Results**. To our knowledge, E2GRE is the only method that has reported their evidence extraction result. The results in Table 3 indicate that EIDER outperforms E2GRE significantly (e.g., by 3.57 Dev Evi F1 under BERT<sub>base</sub>). The results show that it may be sufficient to train the evidence classifier only on pairs with  $r \in \mathcal{R}$ and over each (entity, entity, sentence) tuple instead of (entity, entity, sentence, relation) as in E2GRE.

Our ablation studies in Table 4 show that our 420 three heuristic rules, denoted as Rules (ours), al-421 ready capture most of the evidence for positive en-422 tity pairs. The high quality of silver labels explains 423 why our model can perform well using silver la-424 bels only. Furthermore, training the RE model and 425 evidence extraction model separately (denoted as 426 NoJoint) results in a sharp performance drop. As 427

Model	Dev Evi F1	Test Evi F1
E2GRE-BERT <sub>base</sub>	47.14	48.35
EIDER-BERT <sub>base</sub>	<b>50.71</b>	<b>51.27</b>
E2GRE-RoBERTa <sub>large</sub>	51.11	50.50
EIDER-RoBERTa <sub>large</sub>	<b>52.54</b>	<b>53.01</b>

Table 3: Evidence extraction results. We compare EI-DER with E2GRE (Huang et al., 2021a).

	Rules (ours)	EIDER-BERT <sub>base</sub>	NoJoint
PosEvi F1	77.43	80.33	51.13

Table 4: Ablation study for evidence extraction.

the relation and evidence classifiers share the same base encoder, discarding the relation classifier will result in insufficient training of the base encoder and harm the performance. We use PosEvi F1 because Evi F1 depends on the relation extraction as well, which is not applicable for **Rules (ours)**.

428

429

430

431

432

433

434

435

436

437

438

439

440

441

### 4.3 Performance Analysis

**Ablation Study**. We conduct ablation studies to further analyze the utility of each module in EI-DER. The results are shown in Table 5. Compared to our full model, we observe that **NoJoint** leads to RE performance drop. Compared to ATLOP, EIDER (Rule)-Nojoint achieves significant "free gains" (0.90/1.08 Ign F1/F1) by simply conduct-

Ablation	Ign F1	F1	Intra F1	Inter F1
EIDER-BERT <sub>base</sub>	60.51	62.48	68.47	55.21
NoJoint	59.98	62.03	68.51	54.10
NoPseudo	59.70	61.53	67.55	54.01
NoOrigDoc	58.47	60.44	66.24	53.23
NoBlending	58.93	61.46	67.33	54.37
FinetuneOnEvi	60.11	62.29	68.13	54.84
EIDER (Rule)-BERT <sub>base</sub>	60.36	62.34	68.40	54.79
NoJoint	60.01	62.09	68.21	54.34

Table 5: Ablation study of EIDER on DocRED.

	Co-occur	Coref	Bridge	Total
Count	6711	984	3212	10,907
Percent	54.46%	7.99%	26.07%	88.52%

Table 6: Statistics of the 12,323 relations in the DocRED development set.

ing evidence-enhanced inference, which could be applied to general trained RE models in principle.

We also remove the pseudo document (constructed from the extracted evidence) and the original document separately, denoted as **NoPseudo** and **NoOrigDoc**, respectively. We observe that removing either source will lead to performance drops. Also, the drop of Inter F1 is much larger than Intra F1 for **NoPseudo**, indicating that our inference process is more effective for inter-sentence pairs where the evidence may not be consecutive.

As for **NoBlending**, we remove the blending layer and simply take the union of the two sets of results. The sharp drop of performance indicates the blending layer can successfully learn a dynamic threshold to combine the prediction results.

Finally, we further finetune the RE model on ground truth evidence before feeding it the extracted evidence (denoted as **FinetuneOnEvi**). We observe that the performance is not improved, probably because the encoded entity representations in evidence and original documents are already similar to each other. In fact, when performing relation extraction on the training set using the ground truth evidence alone, the F1 is already over 95%.

**Performance Breakdown**. To further analyze the performance of EIDER on different types of entity pairs, we categorize the relations into three categories based on our three heuristic rules in Sec. 3.3: *Co-occur, Coref* and *Bridge*. The number and percentage of relations covered by each rule are listed in Table 6. We can see that the three categories cover over 88% of the relations in the development set. The results on each category are shown in Figure 3. We can see that our full model has the



Figure 3: Performance gains in F1 by relation categories. The gains are relative to the second best base-line (ATLOP-RoBERTa<sub>large</sub>).

Model	Memory	Training time
ATLOP-BERT <sub>base</sub>	9,139 MB	5.19 it/s
E2GRE-BERT <sub>base</sub>	36,182 MB	0.53 it/s
EIDER-BERT <sub>base</sub>	10,933 MB	4.92 it/s

Table 7: Training time and memory usage on DocRED.

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

501

502

503

505

506

507

509

best performance in all three categories and our ablations also outperform ATLOP. For all our methods, the improvements over ATLOP is  $Bridge > Coref \gg Co-occur$ . This reveals that both modules mainly improve the model's reasoning ability from multiple sentences, either by coreference reasoning or by multi-hop reasoning over a third entity.

Efficiency Comparison. We benchmark the time and memory usage of EIDER on an RTX A6000 GPU. Table 7 show that our joint model incurs only ~5% training time and ~14% GPU memory overhead. Experiments also show that EIDER can be trained on a single consumer GPU (e.g., an 11GB GTX 1080 Ti) but E2GRE is not able to.

#### 4.4 Case Studies

Table 8 shows a few examples of EIDER. Detailed statistics and error analysis are provided in Appendix A.3. In the first example, the head entity is mentioned in the first sentence and the tail entity appears in the second. We can see that EIDER correctly extracts these sentences as evidence. Since the evidence sentences are consecutive, the predictions on both the original document and the evidence sentences are correct. In the second example, the prediction using only the original document is incorrect, possibly because the "King Louie" in the  $1^{st}$  and  $3^{rd}$  sentences are so far away from each other that the model fails to recognize them as coreference. Hence, it fails to distinguish "King Louie" as a bridge entity as wrongly predicts "NA". Instead, these two sentences are consecutive in the extracted evidence, making it easier for the model to find the bridge. In the last example, the  $6^{th}$  sen-

474

475

476

442 443

**Ground Truth Relation: Located in** Ground Truth Evidence Sentence(s): [1, 2] Extracted Evidence Sentence(s): [1, 2] Document: [1] The Portland Golf Club is a private golf club in the northwest United States , in suburban Portland, Oregon. [2] It is located in the unincorporated Raleigh Hills area of eastern Washington County, southwest of downtown Portland and east of Beaverton. [3] The club was established in the winter of 1914, when a group of nine businessmen assembled to form a new club after leaving their respective clubs ... Final Prediction: Located in Prediction on Orig. Doc: Located in Prediction on Extracted Evidences: Located in Ground Truth Relation: Characters Ground Truth Evidence Sentence(s): [1, 3] Extracted Evidence Sentence(s): [1, 3] Document: [1] King Louie is a fictional character introduced in Walt Disney's 1967 animated musical film, The Jungle Book. [2] Unlike the majority of the adapted characters in the film, Louie was not featured in Rudyard Kipling's original works. [3] King Louie was portrayed as an orangutan who was the leader of the other jungle primates, and who attempted to gain knowledge of fire from Mowgli, ... **Final Prediction: Characters** Prediction on Extracted Evidences: Characters Prediction on Orig. Doc: NA Ground Truth Evidence Sentence(s): [5, 6] Ground Truth Relation: Inception Extracted Evidence Sentence(s): [5] Document: [1] Oleg Tinkov (born 25 December 1967) is a Russian entrepreneur and cycling sponsor. .... [5] Tinkoff is the founder and chairman of the Tinkoff Bank board of directors (until 2015 it was called Tinkoff Credit Systems). [6] The bank was founded in 2007 and as of December 1, 2016, it is ranked 45 in terms of assets and 33 for equity among Russian banks. ... Final Prediction: Inception Prediction on Orig. Doc: Inception Prediction on Extracted Evidences<sup>•</sup> NA

Table 8: Case studies of our proposed framework EIDER. We use red, blue and green to color the **head entity**, **tail entity** and **relation**, respectively. The indices of extracted evidence sentences are highlighted with yellow.

tence is missing in the extracted evidence, so the
extracted evidence does not contain enough information to predict the relation. However, the prediction on the original document is correct, leading to
the correct final result.

# 5 Related Work

515

516

517

518

519

520

522

524

525

526

**Relation Extraction**. Previous research efforts on relation extraction mainly concentrate on predicting relations within a sentence (Cai et al., 2016; Zeng et al., 2015; Feng et al., 2018; Zheng et al., 2021; Zhang et al., 2018, 2019, 2020a). While these approaches tackle the sentence-level RE task effectively, in the real world, certain relations can only be inferred from multiple sentences. Consequently, recent studies (Quirk and Poon, 2017; Peng et al., 2017; Yao et al., 2019; Wang et al., 2019; Tang et al., 2020) have proposed to work on the document-level relation extraction (DocRE).

Graph-based DocRE. Graph-based DocRE meth-528 ods generally construct a graph with mentions, en-529 tities, sentences, or documents as the nodes, and 530 infer the relations by reasoning on this graph. Zeng 531 et al. (2020) performs multi-hop reasoning on both 532 a mention-level graph and an entity-level graph. 533 Xu et al. (2021) extracts a reasoning path between each entity pair holding at least one relation and en-535 courages the model to reconstruct the path during training. However, the extracted graph may omit 537 some important information in the text. Compli-538 cated operations on the graphs may also hinder the model from capturing the text structure. 540

541Transformer-based DocRE. Another line of stud-542ies model cross-sentence relations by implicitly543capturing the long-distance token dependencies via

the transformer (Vaswani et al., 2017). Zhou et al. (2021) uses attention in the transformers to extract useful context and adopts an adaptive threshold for each entity pair. Huang et al. (2021b) predicts on only a few sentences selected by rules, which may miss important information and does not show consistent improvements. Instead, EIDER puts more attention to the important sentences without information loss and shows significant improvements. Huang et al. (2021a) extracts relation and evidence together but highly relies on evidence annotations and suffers from massive time and memory overhead. In comparison, EIDER automatically constructs silver evidence labels and improves RE performance by training on these silver labels, achieving strong applicability.

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

# 6 Conclusion

In this work, we propose EIDER, an evidenceenhanced RE framework, which improves DocRE by joint relation and evidence extraction and evidence-enhanced inference. In training, the RE and evidence extraction model provide additional training signals for each other and mutually enhance each other. The joint model is efficient in time and memory and does not rely heavily on the human annotation of evidence. During inference, the prediction results on both the original document and the extracted evidence are combined, which encourages the model to focus on the important sentences while reducing information loss. Experiment results demonstrate that EIDER significantly outperforms existing methods on three public datasets (DocRED, CDR, and GDA), especially on inter-sentence relations.

### References

578

579

581

582

583

584

589

593

595

597

610

611

612

613

614

615

616

617

618

619

625

626

627

629

- Rui Cai, Xiaodong Zhang, and Houfeng Wang. 2016. Bidirectional recurrent convolutional neural network for relation classification. In *ACL*, pages 756–765.
  - 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, pages 4171–4186.
  - Jun Feng, Minlie Huang, Li Zhao, Yang Yang, and Xiaoyan Zhu. 2018. Reinforcement learning for relation classification from noisy data. In *AAAI*, pages 5779–5786.
  - Kevin Huang, Peng Qi, Guangtao Wang, Tengyu Ma, and Jing Huang. 2021a. Entity and evidence guided document-level relation extraction. In *Proceedings* of the 6th Workshop on Representation Learning for NLP (RepL4NLP-2021).
- Quzhe Huang, Shengqi Zhu, Yansong Feng, Yuan Ye, Yuxuan Lai, and Dongyan Zhao. 2021b. Three sentences are all you need: Local path enhanced document relation extraction. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers). Association for Computational Linguistics.
- Robin Jia, Cliff Wong, and Hoifung Poon. 2019. Document-level n-ary relation extraction with multiscale representation learning. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*, pages 3693–3704.
- Jiao Li, Yueping Sun, Robin J. Johnson, Daniela Sciaky, Chih-Hsuan Wei, Robert Leaman, Allan Peter Davis, Carolyn J. Mattingly, Thomas C. Wiegers, and Zhiyong Lu. 2016. Biocreative V CDR task corpus: a resource for chemical disease relation extraction. *Database*, 2016.
- Guoshun Nan, Zhijiang Guo, Ivan Sekulic, and Wei Lu. 2020. Reasoning with latent structure refinement for document-level relation extraction. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 1546–1557.
- Nanyun Peng, Hoifung Poon, Chris Quirk, Kristina Toutanova, and Wen-tau Yih. 2017. Cross-sentence n-ary relation extraction with graph LSTMs. *Transactions of the Association for Computational Linguistics*, 5:101–115.
- Chris Quirk and Hoifung Poon. 2017. Distant supervision for relation extraction beyond the sentence boundary. In Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers, pages 1171–1182.

Sebastian Ruder. 2017. An overview of multitask learning in deep neural networks. *ArXiv*, abs/1706.05098. 634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

688

- Y. Shi, Jiaming Shen, Yuchen Li, N. Zhang, Xinwei He, Zhengzhi Lou, Q. Zhu, M. Walker, Myung-Hwan Kim, and Jiawei Han. 2019. Discovering hypernymy in text-rich heterogeneous information network by exploiting context granularity. In *CIKM*.
- Hengzhu Tang, Yanan Cao, Zhenyu Zhang, Jiangxia Cao, Fang Fang, Shi Wang, and Pengfei Yin. 2020. HIN: hierarchical inference network for documentlevel relation extraction. In Advances in Knowledge Discovery and Data Mining - 24th Pacific-Asia Conference, PAKDD 2020, Singapore, May 11-14, 2020, Proceedings, Part I, volume 12084 of Lecture Notes in Computer Science, pages 197–209.
- Bayu Distiawan Trisedya, Gerhard Weikum, Jianzhong Qi, and Rui Zhang. 2019. Neural relation extraction for knowledge base enrichment. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 229–240, Florence, Italy. Association for Computational Linguistics.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pages 5998–6008.
- Difeng Wang, Wei Hu, Ermei Cao, and Weijian Sun. 2020. Global-to-local neural networks for document-level relation extraction. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3711–3721.
- Hong Wang, Christfried Focke, Rob Sylvester, Nilesh Mishra, and William Wang. 2019. Fine-tune bert for docred with two-step process. *Computing Research Repository*, arXiv:1909.11898.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface's transformers: Stateof-the-art natural language processing. *ArXiv*, pages arXiv–1910.
- David H. Wolpert. 1992. Stacked generalization. *Neural Networks*, 5:241–259.
- Ye Wu, Ruibang Luo, Henry C. M. Leung, Hing-Fung Ting, and Tak Wah Lam. 2019. RENET: A deep learning approach for extracting gene-disease associations from literature. In *Research in Computational Molecular Biology - 23rd Annual International Conference, RECOMB 2019, Washington, DC, USA, May 5-8, 2019, Proceedings*, volume 11467, pages 272–284.

- 693
- 700 701 702 703 704 706
- 710
- 711
- 712
- 713 714 715
- 716 717
- 718 719
- 721 724

725

- 726 727 728 729

- 734
- 737
- 738 739
- 740 741 742

- 743 744 745

- Liyan Xu and Jinho D. Choi. 2020. Revealing the myth of higher-order inference in coreference resolution. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 8527–8533. Association for Computational Linguistics.
- Wang Xu, Kehai Chen, and Tiejun Zhao. 2021. Document-level relation extraction with reconstruction. Proceedings of the AAAI Conference on Artificial Intelligence.
- Yuan Yao, Deming Ye, Peng Li, Xu Han, Yankai Lin, Zhenghao Liu, Zhiyuan Liu, Lixin Huang, Jie Zhou, and Maosong Sun. 2019. DocRED: A large-scale document-level relation extraction dataset. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 764–777.
- Deming Ye, Yankai Lin, Jiaju Du, Zhenghao Liu, Peng Li, Maosong Sun, and Zhiyuan Liu. 2020. Coreferential Reasoning Learning for Language Representation. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 7170-7186.
- Mo Yu, Wenpeng Yin, Kazi Saidul Hasan, Cicero dos Santos, Bing Xiang, and Bowen Zhou. 2017. Improved neural relation detection for knowledge base question answering. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 571-581, Vancouver, Canada. Association for Computational Linguistics.
- Daojian Zeng, Kang Liu, Yubo Chen, and Jun Zhao. 2015. Distant supervision for relation extraction via piecewise convolutional neural networks. In EMNLP, pages 1753-1762.
- Shuang Zeng, Runxin Xu, Baobao Chang, and Lei Li. 2020. Double graph based reasoning for documentlevel relation extraction. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 1630–1640.
- Ningyu Zhang, Shumin Deng, Zhanlin Sun, Jiaoyan Chen, Wei Zhang, and Huajun Chen. 2020a. Relation adversarial network for low resource knowledge graph completion. In Proceedings of The Web Conference 2020.
- Ningyu Zhang, Shumin Deng, Zhanlin Sun, Xi Chen, Wei Zhang, and Huajun Chen. 2018. Attentionbased capsule networks with dynamic routing for relation extraction. In EMNLP.
- Ningyu Zhang, Shumin Deng, Zhanlin Sun, Guanying Wang, Xi Chen, Wei Zhang, and Huajun Chen. 2019. Long-tail relation extraction via knowledge graph embeddings and graph convolution networks. In NAACL-HLT.
- Zhenyu Zhang, Bowen Yu, Xiaobo Shu, Tingwen Liu, Hengzhu Tang, Wang Yubin, and Li Guo. 2020b. Document-level relation extraction with dual-tier

heterogeneous graph. In Proceedings of the 28th International Conference on Computational Linguistics, pages 1630–1641, Barcelona, Spain (Online). International Committee on Computational Linguistics.

746

747

749

750

751

753

755

756

757

758

759

- Hengyi Zheng, Rui Wen, Xi Chen, Yifan Yang, Yunnan Zhang, Ziheng Zhang, Ningyu Zhang, Bin Qin, Xu Ming, and Yefeng Zheng. 2021. Prgc: Potential relation and global correpondence based joint relational triple extraction. In ACL.
- Wenxuan Zhou, Kevin Huang, Tengyu Ma, and Jing Huang. 2021. Document-level relation extraction with adaptive thresholding and localized context pooling. In Proceedings of the AAAI Conference on Artificial Intelligence.

### A Appendices

#### A.1 Dataset Statistics

Our model is evaluated on three benchmark datasets, where the statistics are shown in Table 9:

**DocRED (Yao et al., 2019)** is a large humanannotated document-level RE dataset. DocRED is constructed from Wikipedia, involving 97 relation types (including NA), 132,275 entities, and 56,354 relations. In the training set, around 97.03% entity pairs do not hold any explicit relations. Around 54.2% of the relations are intra-sentence. The others can only be extracted by considering multiple sentences. In our experiments, the performance on the test set is validated through the Leader board. <sup>2</sup>

**CDR** (Li et al., 2016) is a relation extraction dataset in the biomedical domain. The only two entity types are chemicals and diseases and the only two relations are NA and the causal relation between chemicals and disease concepts.

**GDA** (Wu et al., 2019) is also a biomedical dataset with two entity types only: diseases and genes, and one non-NA relation type only: the interactions between disease concepts and genes.

Statistics	DocRED	CDR	GDA
# Train	3053	500	23353
# Dev	1000	500	5839
# Test	1000	500	1000
# Relation types	97	2	2
# Avg.# entities per Doc	19.5	7.6	5.4
# Avg.# sentences per Doc	8.0	9.7	10.2
Percent of Intra Rel	54.2	75.7	84.7

Table 9: Statistics of the datasets in experiments. The percentage of intra-sentence relations is calculated from the development set of DocRED and calculated from the test set of CDR and GDA.

#### A.2 Implementation Details

Our model is implemented based on PyTorch and Huggingface's Transformers (Wolf et al., 2019). We use cased-BERT<sub>base</sub> (Devlin et al., 2019) and RoBERTa<sub>large</sub> as the base encoders and optimize our model using AdamW with learning rate 5e-5 for the encoder and 1e - 4 for other parameters. We adopt a linear warmup for the first 6% steps. The batch size (number of documents per batch) is set to 4 and the ratio between relation extraction and evidence extraction losses is set to 0.1. We perform early stopping based on the F1 score on the development set, with a maximum of 30 epochs. Our BERT<sub>base</sub> models are trained with one GTX 1080 Ti GPU and RoBERTa<sub>large</sub> models with one RTX A6000 GPU.

#### A.3 Error Analysis of EIDER

		Grou	nd Truth
0U		$r \in \mathcal{R}$	NA
icti	$r \in \mathcal{R}$ (Correct)	7,696 (🗸 )	3 613 ( <b>X</b> )
red	$r \in \mathcal{R}$ (Wrong)	287 ( <b>X</b> )	5,015 (🖍)
Б	NA	4,340 (X)	380,854 (🗸 )

Table 10: Statistics of one run of EIDER-RoBERTa<sub>large</sub>. " $r \in \mathcal{R}$ " means non-NA relations. We use " $\checkmark$ " or " $\checkmark$ " to denote whether the prediction is correct or wrong. For example, we have 4,340 wrong predictions where the ground truth is some  $r \in \mathcal{R}$  but the prediction is NA.

Reason	Count
Labeling Mistakes*	18
Fail in Commonsense Reasoning	11
Fail in Coreferential Reasoning	6
Fail in Multi-hop Reasoning	4
Wrong Evidence Extraction	1
Others	11

Table 11: Error types of EIDER in 50 randomly sampled error cases in DocRED. Where "Labeling Mistakes" means our model predicts correctly but the annotation is wrong.

The detailed statistics of the predictions of our model are listed in Table 10. Among all the errors, the majority is because the model wrongly predicts the non-NA relations (i.e.,  $r \in \mathcal{R}$ ) as "NA" or predicts "NA" as some non-NA relations. Only  $\frac{287}{287+4340+3613} = 3.48\%$  of the errors result from wrongly taking some non-NA relation as another.

To check the exact reason why our model makes these errors, we randomly select 50 cases from DocRED where our model predicts wrongly. We summarize the error types in Table 11 and provide one or two examples for each of the common error types in Table 12.

Our analysis shows that 18 out of 50 "error cases" are actually correct. It suggests that labeling mistakes are still prevalent in the DocRED dataset. We show an example under "*Error Type 1*" in Table 12. The annotator wrongly labels "*U.S. Route 20*", a highway, as the country of "*Capital District*".

Another common error type is failing to conduct commonsense reasoning. Some examples 801

802

803

804

11

781

782

784

790

794

\_

<sup>&</sup>lt;sup>2</sup>Results can be found at https://competitions. codalab.org/competitions/20717.

Error Type 1: Labeling Mistakes			
Ground Truth Relation: Country (X)	<b>Ground Truth Evidence Sentence(s)</b> : [1, 4, 5]	, 7]Extracted Evidence Sentence(s):[5, 7]	
<b>Document</b> : [1] Westmere is a hamlet in the [5] U.S. Route 20 (Western Avenue) bisect District's largest shopping mall is in West	town of Guilderland, Albany County, New York. ts the community and is the major thoroughfare as pere's portheastern corner.	[4] It is a suburb of the neighboring city of Albany. nd main street [7] Crossgates Mall, the Capital	
Final Prediction: NA	Prediction on Orig. Doc: NA	Prediction on Extracted Evidences: NA	
Fror Type 2: Eail in Commonsonse Pageo	ning		
Ground Truth Relation: Country	Ground Truth Evidence Sentence(s): []	Extracted Exidence Sentence(s): [1 4]	
<b>Document:</b> [1] A Route Army was a type more corps or a large number of divisions o War but was discarded as a formation type b Army. [3] Some of the more famous of the <b>Final Prediction</b> : NA (X)	of military organization during the Chinese Rep r independent brigades. [2] It was a common for y the National Revolutionary Army after 1938 (ot Route Armies were: [4] 8th Route Army: Comm <b>Prediction on Orig. Doc</b> : NA (X)	bublic, and usually exercised command over two or mation in China prior to the Second Sino-Japanese her than the 8th Route Army), in favor of the Group munist guerrilla force in North China Prediction on Extracted Evidences: NA (X)	
<b>Ground Truth Relation: Country, Locate</b> <b>Document:</b> [1] The Treaty of Edinburgh- [2] It brought an end to the First War of Scc was signed in Edinburgh by Robert the Bru on 1 May. [4] The document was written in Final Perdiction: NA (X) Prod	<b>d in Ground Truth Evidence Sentence(s)</b> : [1 –Northampton was a peace treaty, signed in 132: ottish Independence, which had begun with the E ice, King of Scotland, on 17 March 1328, and wa French, and is held by the National Archives of liction on Orig Dec: NA (X)	1, 2, 3, 4] <b>Extracted Evidence Sentence(s):</b> [1] 8 between the Kingdoms of England and Scotland. nglish invasion of Scotland in 1296. [3] The treaty is ratified by the English Parliament at Northampton Scotland in Edinburgh	
Final Frediction. NA (A) Fred		on on Extracted Evidences. Country, Located in	
Error Type 3: Fail in Coreferential Reason	ing Cround Truth Exidence Sentence(c): []	Extracted Evidence Sontange(c): [1]	
Document: [1] Monon Polletti (1740-17	Ground Huth Evidence Sentence(s). []	ing in Erange and lover of the famous womenizer	
<b>Document:</b> [1] Manon balletti $(1/40-1/6)$ was the daughter of italian actors performing in France and lover of the famous womanizer Giacomo Casanova. [2] She was ten years old when she first met him; she happened to be the daughter of Silvia Balletti, an actress of the Comédie Italienne company and younger sister of Casanova's closest friend			
Final Prediction: Child (X)	Prediction on Orig. Doc: Child (X)	Prediction on Extracted Evidences: Child $(X)$	
Error Type 4: Fail in Multi-hop Reasoning			
Ground Truth Relation: Educated at	Ground Truth Evidence Sentence(s): [4	] Extracted Evidence Sentence(s): [4]	
<b>Document:</b> [1] Ronald Leonard is an Amer and teacher [4] He was a winner of the with Leonard Rose and Orlando Colo	rican cellist . [2] He has had a distinguished care e Walter Naumburg Competition while a student	er as a soloist, chamber musician , principal cellist at the Curtis Institute of Music, where he studied	
Final Prediction: NA $(X)$	<b>Prediction on Orig. Doc:</b> NA ( <b>X</b> )	<b>Prediction on Extracted Evidences:</b> NA ( $\boldsymbol{X}$ )	

Table 12: Examples for the four most common error types. We use red, blue and green to color the head entity, tail entity and relation, respectively. The indices of extracted evidence sentences are highlighted with yellow.

822 require direct reasoning from the surface names of the head and tail entities. As shown in the first example under "Error Type 2" in Table 12, 824 humans can directly identify that "China" is the country of North China without reading the doc-826 ument, despite that there are no clue in the docu-827 ment indicates this relation. However, most DocRE 828 models, including EIDER, learn to predict the relations only based on the given document and sometimes fail in such cases. There are also examples that need background knowledge. When a human 832 checks the second example, the reasoning process could be: "Scotland" and "England" signed the "Edinburgh-Northampton Treaty". We already know that "Northampton" is in "England", so it is highly possible that "Edinburgh" locates in "Scotland". However, this kind of reasoning involves back-839 ground knowledge that does not come from the document, such as the fact that "Northampton" is in "England". Even though our prediction on extracted evidence is correct, the confidence is still not high, leading to the incorrect final prediction.

Fail in Coreferential Reasoning", human can still identify the correct relation based on the extracted evidence only. As shown in our example in Table 12, in the first sentence, the model wrongly predicts "Giacomo Casanova" as the father of "Manon Balletti", but her real father should be an "Italian actor performing in France". It shows that even the reasoning within a single sentence can be difficult. Similarly, the example in "Error Type 4" also shows that the prediction can still be wrong even if we extract the correct evidence sentences and simplify the problem to sentence-level RE. This suggests that if the performance of sentencelevel RE is improved, the performance of DocRE will also improve.

845

846

847

848

849

850

851

852

853

854

855

856

857

858