

# A sequence-to-sequence approach for document-level relation extraction

John Giorgi<sup>1,4,5,✉</sup> Gary D. Bader<sup>1,2,4,6,7,†</sup> Bo Wang<sup>1,3,5,8,†</sup>

<sup>1</sup>Department of Computer Science, University of Toronto

<sup>2</sup>Department of Molecular Genetics, University of Toronto

<sup>3</sup>Department of Laboratory Medicine and Pathobiology, University of Toronto

<sup>4</sup>Terrence Donnelly Centre for Cellular & Biomolecular Research

<sup>5</sup>Vector Institute for Artificial Intelligence

<sup>6</sup>The Lunenfeld-Tanenbaum Research Institute, Sinai Health System

<sup>7</sup>Princess Margaret Cancer Centre, University Health Network

<sup>8</sup>Peter Munk Cardiac Center, University Health Network

✉Corresponding author †Equal contribution

{john.giorgi, gary.bader}@mail.utoronto.ca

bowang@vectorinstitute.ai

## Abstract

Motivated by the fact that many relations cross the sentence boundary, there has been increasing interest in document-level relation extraction (DocRE). DocRE requires integrating information within and across sentences, capturing complex interactions between mentions of entities. Most existing methods are pipeline-based, requiring entities as input. However, jointly learning to extract entities and relations can improve performance and be more efficient due to shared parameters and training steps. In this paper, we develop a sequence-to-sequence approach, seq2rel, that can learn the subtasks of DocRE (entity extraction, coreference resolution and relation extraction) end-to-end, replacing a pipeline of task-specific components. Using a simple strategy we call entity hinting, we compare our approach to existing pipeline-based methods on several popular biomedical datasets, in some cases exceeding their performance. We also report the first end-to-end results on these datasets for future comparison. Finally, we demonstrate that, under our model, an end-to-end approach outperforms a pipeline-based approach. Our code, data and trained models are available at <https://github.com/johngiorgi/seq2rel>. An online demo is available at <https://share.streamlit.io/johngiorgi/seq2rel/main/demo.py>.

## 1 Introduction

PubMed, the largest repository of biomedical literature, contains over 30 million publications and is adding more than two papers per minute. Accurate, automated text mining and natural language processing (NLP) methods are needed to maximize discovery and extract structured information from

this massive volume of text. An important step in this process is relation extraction (RE), the task of identifying groups of entities within some text that participate in a semantic relationship. In the domain of biomedicine, relations of interest include chemical-induced disease, protein-protein interactions, and gene-disease associations.

Many methods have been proposed for RE, ranging from rule-based to machine learning-based (Zhou et al., 2014; Liu et al., 2016). Most of this work has focused on *intra*-sentence binary RE, where pairs of entities within a sentence are classified as belonging to a particular relation (or none). These methods often ignore commonly occurring complexities like nested or discontinuous entities, coreferent mentions (words or phrases in the text that refer to the same entity), inter-sentence and *n*-ary relations (see Figure 1 for examples). The decision not to model these phenomena is a strong assumption. In GENIA (Kim et al., 2003), a corpus of PubMed articles labelled with around 100,000 biomedical entities, ~17% of all entities are nested within another entity. Discontinuous entities are particularly common in clinical text, where ~10% of mentions in popular benchmark corpora are discontinuous (Wang et al., 2021). In the CDR corpus (Li et al., 2016b), which comprises 1500 PubMed articles annotated for chemical-induced disease relations, ~30% of all relations are inter-sentence. Some relations, like drug-gene-mutation interactions, are difficult to model with binary RE (Zhou et al., 2014).

In response to some of these shortcomings, there has been a growing interest in *document*-level RE (DocRE). DocRE aims to model *inter*-sentence re-

Figure 1: Examples of complexities in entity and relation extraction and the proposed linearization schema to model them. CID: chemical-induced disease. GDA: gene-disease association. DGM: drug-gene-mutation.

Complexities	Example	Comment
Discontinuous mentions	Induction by <b>paracetamol</b> of <b>bladder</b> and <b>liver tumours</b> .  <code>paracetamol @DRUG@ bladder tumours @DISEASE@ @CID@</code> <code>paracetamol @DRUG@ liver tumours @DISEASE@ @CID@</code>	Discontinuous mention of <b>bladder tumours</b> .
Coreferent mentions	Proto-oncogene <b>HER2</b> (also known as <b>erbb-2</b> or <b>neu</b> ) plays an important role in the carcinogenesis and the prognosis of <b>breast cancer</b> .  <code>her2 ; erbb-2 ; neu @GENE@ breast cancer @DISEASE@ @GDA@</code>	Two coreferent mentions of <b>HER2</b> .
$n$ -ary, inter-sentence	The deletion mutation on exon-19 of <b>EGFR</b> gene was present in 16 patients, while the <b>L858E</b> point mutation on exon-21 was noted in 10. All patients were treated with <b>gefitinib</b> and showed a partial response.  <code>gefitinib @DRUG@ egfr @GENE@ l858e @MUTATION@ @DGM@</code>	Ternary <b>DGM</b> relationship crosses a sentence boundary.

lations between coreferent mentions of entities in a document. A popular approach involves graph-based methods, which have the advantage of naturally modelling inter-sentence relations (Peng et al., 2017; Song et al., 2018; Christopoulou et al., 2019; Nan et al., 2020; Minh Tran et al., 2020). However, like all pipeline-based approaches, these methods assume that the entities within the text are known. As previous work has demonstrated, and as we show in §5.2, jointly learning to extract entities and relations can improve performance (Miwa and Sasaki, 2014; Miwa and Bansal, 2016; Gupta et al., 2016; Li et al., 2016a, 2017; Nguyen and Verspoor, 2019a; Yu et al., 2020) and may be more efficient due to shared parameters and training steps. Existing end-to-end methods typically combine task-specific components for entity detection, coreference resolution, and relation extraction that are trained jointly. Most approaches are restricted to intra-sentence RE (Bekoulis et al., 2018; Luan et al., 2018; Nguyen and Verspoor, 2019b; Wadden et al., 2019; Giorgi et al., 2019) and have only recently been extended to DocRE (Eberts and Ulges, 2021). However, they still focus on binary relations. Ideally, DocRE methods would be capable of modelling the complexities mentioned above without strictly requiring entities to be known.

A less popular end-to-end approach is to frame RE as a *generative* task with sequence-to-sequence (seq2seq) learning (Sutskever et al., 2014). This framing simplifies RE by removing the need for task-specific components and explicit negative training examples, i.e. pairs of entities that *do not* express a relation. If the information to extract is appropriately linearized to a string, seq2seq methods are flexible enough to model all complexities

discussed thus far. However, existing work stops short, focusing on intra-sentence binary relations (Zeng et al., 2018; Zhang et al., 2020; Nayak and Ng, 2020; Zeng et al., 2020). In this paper, we extend work on seq2seq methods for RE to the document level, with several important contributions:

- We propose a novel linearization schema that can handle complexities overlooked by previous seq2seq approaches, like coreferent mentions and  $n$ -ary relations (§3.1).
- Using this linearization schema, we demonstrate that a seq2seq approach is able to learn the subtasks of DocRE (entity extraction, coreference resolution and relation extraction) jointly, and report the first end-to-end results on several popular biomedical datasets (§5.1).
- We devise a simple strategy, referred to as “entity hinting” (§3.3), to compare our model to existing pipeline-based approaches, in some cases exceeding their performance (§5.1).

## 2 Task definition: document-level relation extraction

Given a source document of  $S$  tokens, a model must extract all tuples corresponding to a relation,  $R$ , expressed between the entities,  $E$  in the document,  $(E_1, \dots, E_n, R)$  where  $n$  is the number of participating entities, or *arity*, of the relation. Each entity  $E_i$  is represented as the set of its coreferent mentions  $\{e_j^i\}$  in the document, which are often expressed as aliases, abbreviations or acronyms. All entities appearing in a tuple have at least one mention in the document. The mentions that express a given relation are not necessarily contained within

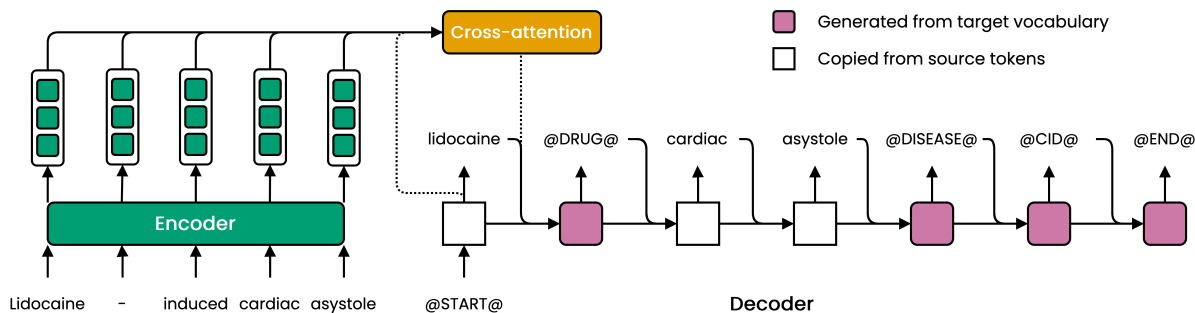


Figure 2: A sequence-to-sequence model for document-level relation extraction. Special tokens are generated by the decoder. Entity mentions are copied from the input via a copy mechanism (not shown). Decoding is initiated by a @START@ token and terminated when the model generates the @END@ token. Attention connections shown only for the second timestep to reduce clutter. CID: chemical-induced disease.

the same sentence. Commonly,  $E$  is assumed to be known and provided as input to a model. We will refer to these methods as “pipeline-based”. In this paper, we are primarily concerned with the situation where  $E$  is *not* given and must be predicted by a model, which we will refer to as “end-to-end”.

### 3 Our approach: seq2rel

#### 3.1 Linearization

To use seq2seq learning for RE, the information to be extracted must be linearized to a string. This linearization should be expressive enough to model the complexities of entity and relation extraction without being overly verbose. We propose the following schema, illustrated with an example:

$X$ : Variants in the **estrogen receptor alpha (ESR1)** gene and its mRNA contribute to risk for **schizophrenia**.

$Y$ : **estrogen receptor alpha** ; **ESR1** @GENE@ **schizophrenia** @DISEASE@ @GDA@

The input text  $X$ , expresses a gene-disease association (GDA) between **ESR1** and **schizophrenia**. In the corresponding target string  $Y$ , each relation begins with its constituent entities. A semicolon separates coreferent mentions (;), and entities are terminated with a special token denoting their type (e.g. @GENE@). Similarly, relations are terminated with a special token denoting their type (e.g. @GDA@). Two or more entities can be included before the special relation token to support  $n$ -ary extraction. Entities can be ordered if they serve specific roles as head or tail of a relation. For each document, multiple relations can be included in the target string. Entities may be nested or discontinuous in the input text. In Figure 1, we provide examples of how this

schema can be used to model various complexities, like coreferent entity mentions and  $n$ -ary relations.

#### 3.2 Model

The model follows a canonical seq2seq setup. An encoder maps each token in the input to a contextual embedding. An autoregressive decoder generates an output, token-by-token, attending to the outputs of the encoder at each timestep (Figure 2). Decoding proceeds until a special “end-of-sequence” token (@END@) is generated, or a maximum number of tokens have been generated. Formally,  $X$  is the *source* sequence of length  $S$ , which is some text we would like to extract relations from.  $Y$  is the corresponding *target* sequence of length  $T$ , a linearization of the relations contained in the source. We model the conditional probability

$$p(Y|X) = \prod_{t=1}^T p(y_t|X, y_{<t}) \quad (1)$$

During training, we optimize over the model parameters  $\theta$  the sequence cross-entropy loss

$$\ell(\theta) = - \sum_{t=1}^T \log p(y_t|X, y_{<t}; \theta) \quad (2)$$

maximizing the log-likelihood of the training data.<sup>1</sup>

The main problems with this setup for RE are: 1) The model might “hallucinate” by generating entity mentions that do not appear in the source text. 2) It may generate a target string that does not follow the linearization schema and therefore cannot

<sup>1</sup>See §4.3 for details about the encoder and decoder.

be parsed. 3) The loss function is permutation-sensitive, enforcing an unnecessary decoding order. To address 1) we use two modifications: a restricted target vocabulary (§3.2.1) and a copy mechanism (§3.2.2). To address 2) we experiment with several constraints applied during decoding (§3.2.3). Finally, to address 3) we sort relations according to their order of appearance in the source text (§3.2.4).

### 3.2.1 Restricted target vocabulary

To prevent the model from “hallucinating” (generating entity mentions that do not appear in the source text), the target vocabulary is restricted to the set of special tokens needed to model entities and relations (e.g. ; and @DRUG@). All other tokens must be copied from the input using a copy mechanism (see §3.2.2). The embeddings of these special tokens are initialized randomly and learned jointly with the rest of the model’s parameters.

### 3.2.2 Copy mechanism

To enable copying of input tokens during decoding, we use a copying mechanism (Gu et al., 2016a). The mechanism works by effectively extending the target vocabulary with the tokens in the source sequence  $X$ , allowing the model to “copy” these tokens into the output sequence,  $Y$ . Our use of the copy mechanism is similar to previous seq2seq-based approaches for RE (Zeng et al., 2018, 2020).

### 3.2.3 Constrained decoding

We experimented with several constraints applied to the decoder during test time to reduce the likelihood of generating syntactically invalid target strings (strings that do not follow the linearization schema). These constraints are applied by setting the predicted probabilities of invalid tokens to a tiny value at each timestep. The full set of constraints is depicted in Appendix A. In practice, we found that a trained model rarely generates invalid target strings, so these constraints have little effect on final performance (see §5.3). We elected not to apply them in the rest of our experiments.

### 3.2.4 Sorting relations

The relations to extract from a given document are inherently unordered. However, the sequence cross-entropy loss (Equation 2) is permutation-sensitive with respect to the predicted tokens. During training, this enforces an unnecessary decoding order and may make the model prone to overfit frequent token combinations in the training set (Vinyals

et al., 2016; Yang et al., 2019). To partially mitigate this, we sort relations within the target strings according to their order of appearance in the source text, providing the model with a consistent decoding order. The position of a relation is determined by the first occurring mention of its head entity. The position of a mention is determined by the sum of its start and end character offsets. In the case of ties, we then sort by the first mention of its tail entity (and so on for  $n$ -ary relations).

## 3.3 Entity hinting

Although the proposed model can jointly extract entities and relations from unannotated text, most existing DocRE methods provide the entities as input. Therefore, to more fairly compare to existing methods, we also provide entities as input, using a simple strategy that we will refer to as “entity hinting”. This involves prepending entities to the source text as they appear in the target string. Taking the example from §3.1, entity hints would be added as follows:

$X$ : estrogen receptor alpha ; ESR1 @GENE@ schizophrenia @DISEASE@ @SEP@ Variants in the estrogen receptor alpha (ESR1) gene and its mRNA contribute to risk for schizophrenia.

where the special @SEP@ token demarcates the end of the entity hint.<sup>2</sup> We experimented with the common approach of inserting marker tokens before and after each entity mention (Zhou and Chen, 2021) but found this to perform worse. Our approach adds fewer extra tokens to the source text and provides a location for the copy mechanism to focus, i.e. tokens left of @SEP@. In our experiments, we use entity hinting when comparing to methods that provide ground truth entity annotations as input (§5.1.1). In §5.2, we use entity hinting to compare pipeline-based and end-to-end approaches.

## 4 Experimental setup

### 4.1 Datasets

We evaluate our approach on several biomedical, DocRE datasets. We also include one non-biomedical dataset, DocRED. In Appendix B, we list relevant details about their annotations.

<sup>2</sup>Some pretrained models have their own separator token which can be used in place of @SEP@, e.g. BERT uses [SEP].

**CDR (Li et al., 2016b)** The BioCreative V CDR task corpus is manually annotated for chemicals, diseases and chemical-induced disease (CID) relations. It contains the titles and abstracts of 1500 PubMed articles and is split into equally sized train, validation and test sets. Given the relatively small size of the training set, we follow [Christopoulou et al. \(2019\)](#) and others by first tuning the model on the validation set and then training on the combination of the train and validation sets before evaluating on the test set. Similar to prior work, we filter negative relations with disease entities that are hypernyms of a corresponding true relations disease entity within the same abstract (see [Appendix C](#)).

**GDA (Wu et al., 2019)** The gene-disease association corpus contains 30,192 titles and abstracts from PubMed articles that have been automatically labelled for genes, diseases and gene-disease associations via distant supervision. The test set is comprised of 1000 of these examples. Following [Christopoulou et al. \(2019\)](#) and others, we hold out a random 20% of the remaining abstracts as a validation set and use the rest for training.

**DGM (Jia et al., 2019)** The drug-gene-mutation corpus contains 4606 PubMed articles that have been automatically labelled for drugs, genes, mutations and ternary drug-gene-mutation relationships via distant supervision. The dataset is available in three variants: sentence, paragraph, and document-length text. We train and evaluate our model on the paragraph-length inputs. Since the test set does not contain relation annotations on the paragraph level, we report results on the validation set. We hold out a random 20% of training examples to form a new validation set for tuning.

**DocRED (Yao et al., 2019)** DocRED includes over 5000 human-annotated documents from Wikipedia. There are six entity and 96 relation types, with  $\sim 40\%$  of relations crossing the sentence boundary. We use the same split as previous end-to-end methods ([Eberts and Ulges, 2021](#)), which has 3,008 documents in the training set, 300 in the validation set and 700 in the test set<sup>3</sup>.

## 4.2 Evaluation

We evaluate our model using the micro F1-score by extracting relation tuples from the decoder’s output (see [Appendix D](#)). Similar to prior work, we use a “strict” criteria. A predicted relation is considered

correct if the relation type and its entities match a ground truth relation. An entity is considered correct if the entity type and its mentions match a ground truth entity. However, since the aim of DocRE is to extract relations at the *entity*-level (as opposed to the *mention*-level), we also report performance using a relaxed criterion (denoted “relaxed”), where predicted entities are considered correct if more than 50% of their mentions match a ground truth entity (see [Appendix E](#)).

Existing methods that evaluate on CDR, GDA and DGM use the ground truth entity annotations as input. This makes it difficult to directly compare with our end-to-end approach, which takes only the raw text as input. To make the comparison fairer, we use entity hinting (§3.3) so that our model has access to the ground truth entity annotations. We also report the performance of our method in the end-to-end setting on these corpora to facilitate future comparison. To compare to existing end-to-end approaches, we use DocRED.

## 4.3 Implementation, training and hyperparameters

**Implementation** We implemented our model in PyTorch ([Paszke et al., 2017](#)) using AllenNLP ([Gardner et al., 2018](#)). As encoder, we use a pre-trained transformer, implemented in the Transformers library ([Wolf et al., 2020](#)), which is fine-tuned during training. When training and evaluating on biomedical corpora, we use PubMedBERT ([Gu et al., 2020](#)), and BERT<sub>BASE</sub> ([Devlin et al., 2019](#)) otherwise. In both cases, we use the default hyperparameters of the pretrained model. As decoder, we use a single-layer LSTM ([Hochreiter and Schmidhuber, 1997](#)) with randomly initialized weights. We use multi-head attention ([Vaswani et al., 2017](#)) as the cross-attention mechanism between encoder and decoder. Select hyperparameters were tuned on the validation sets, see [Appendix F](#) for details.

**Training** All parameters are trained jointly using the AdamW optimizer ([Loshchilov and Hutter, 2019](#)). Before training, we re-initialize the top  $L$  layers of the pretrained transformer encoder, which has been shown to improve performance and stability during fine-tuning ([Zhang et al., 2021b](#)). During training, the learning rate is linearly increased for the first 10% of training steps and linearly decayed to zero afterward. Gradients are scaled to a vector norm of 1.0 before backpropagating. During each forward propagation, the hidden state of the LSTM

<sup>3</sup><https://github.com/lavis-nlp/jerex>

Table 1: Comparison to existing pipeline-based methods. Performance reported as micro-precision, recall and F1-scores (%) on the CDR and GDA test sets. Results below the horizontal line are not comparable to existing methods. Bold: best scores.

Method	CDR			GDA		
	P	R	F1	P	R	F1
Christopoulou et al. (2019)	62.1	65.2	63.6	-	-	81.5
Nan et al. (2020)	-	-	64.8	-	-	82.2
Minh Tran et al. (2020)	-	-	66.1	-	-	82.8
Lai and Lu (2021)	64.9	67.1	66.0	-	-	-
Xu et al. (2021)	-	-	68.7	-	-	83.7
Zhou et al. (2021)	-	-	<b>69.4</b>	-	-	83.9
seq2rel (entity hinting)	68.2	66.2	67.2	84.4	85.3	<b>84.9</b>
seq2rel (entity hinting, relaxed)	68.2	66.2	67.2	84.5	85.4	85.0
seq2rel (end-to-end)	43.5	37.5	40.2	55.0	55.4	55.2
seq2rel (end-to-end, relaxed)	56.6	48.8	52.4	70.3	70.8	70.5

decoder is initialized with the mean of token embeddings output by the encoder. The decoder is regularized by applying dropout (Srivastava et al., 2014) with probability 0.1 to its inputs, and Drop-Connect (Wan et al., 2013) with probability 0.5 to the hidden-to-hidden weights. As is common, we use teacher forcing, feeding previous ground truth inputs to the decoder when predicting the next token in the sequence. During test time, we generate the output using beam search (Graves, 2012). Beams are ranked by mean token log probability after applying a length penalty.<sup>4</sup> Models were trained and evaluated on a single NVIDIA Tesla V100.<sup>5</sup>

## 5 Results

### 5.1 Comparison to existing methods

In the following sections, we compare our model to existing DocRE methods on several benchmark corpora. We compare to existing pipeline-based methods (§5.1.1), including  $n$ -ary methods (§5.1.2), and end-to-end methods (§5.1.3). Details about these methods are provided in Appendix G.

#### 5.1.1 Existing pipeline-based methods

In Table 1, we use entity hinting to compare our method to existing pipeline-based methods on CDR and GDA. We also report end-to-end performance, which is not comparable to existing pipeline-based methods but will facilitate future comparisons.

The large performance improvement when using entity hinting (+27-29%) confirms that the model

<sup>4</sup>[https://docs.allennlp.org/main/api/nn/beam\\_search/#lengthnormalizedsequencelogprobabilityscorer](https://docs.allennlp.org/main/api/nn/beam_search/#lengthnormalizedsequencelogprobabilityscorer)  
<sup>5</sup><https://www.nvidia.com/en-us/data-center/v100/>

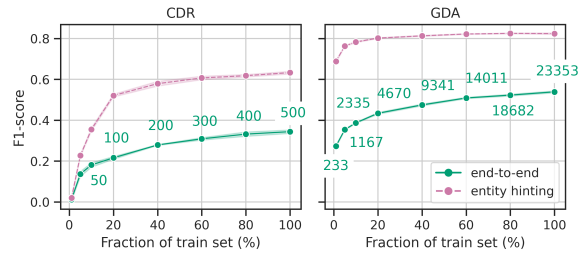


Figure 3: Effect of training set size on performance. Performance reported as the median micro F1-score obtained over five runs with different random seeds on the CDR and GDA validation sets, with and without entity hinting. Error bands correspond to the standard deviation over the five runs. The absolute number of training examples are displayed for each corpus. Some labels are excluded to reduce clutter.

exploits the entity annotations. The fact that relaxed entity matching makes a large difference in the end-to-end setting (+12-15%) suggests that a significant portion of the model’s mistakes occur during coreference resolution. Although our method is designed for end-to-end RE, we find that it outperforms existing pipeline-based methods when using entity hinting on GDA. Our method is competitive with existing methods when using entity hinting on the CDR corpus but ultimately underperforms state-of-the-art results. Given that GDA is 46X larger, we speculated that our method might be underperforming in the low-data regime. To determine if this is a contributing factor, we artificially reduce the size of the CDR and GDA training sets and plot the performance as a curve (Figure 3). In all cases besides GDA with entity hinting, performance increases monotonically with dataset size. There is no obvious plateau on CDR even when using all 500 training examples. Together, these results suggest that our seq2seq based approach can outperform existing pipeline-based methods when there are sufficient training examples but underperforms relative to existing methods in the low-data regime.

#### 5.1.2 $n$ -ary relation extraction

In Table 2 we compare to existing  $n$ -ary methods on the DGM corpus. With entity hinting, our method significantly outperforms the existing method. The difference in encoders partially explains this large performance gap. Where Jia et al. (2019) use a BiLSTM that is trained from scratch, we use PubMedBERT, a much larger model that has been pretrained on abstracts and full-text ar-

Table 2: Comparison to existing  $n$ -ary methods. Performance reported as micro-precision, recall and F1-scores (%) on the DGM validation set. Results below the horizontal line are not comparable to existing methods. Bold: best scores. † Jia et al. 2019 do not report results on the validation set, so we re-run their paragraph-level model.

Method	P	R	F1
Jia et al. (2019) †	62.9	76.2	68.9
seq2rel (entity hinting)	<b>84.0</b>	<b>84.8</b>	<b>84.4</b>
seq2rel (entity hinting, relaxed)	84.1	84.9	84.5
seq2rel (end-to-end)	68.9	65.9	67.4
seq2rel (end-to-end, relaxed)	78.3	74.9	76.6

articles from PubMedCentral.<sup>6</sup> However, this does not completely account for the improvement in performance, as recent work that has replaced the BiLSTM encoder of (Jia et al., 2019) with PubMedBERT found that it improves performance by approximately 2-4% on the task of drug-gene-mutation prediction (Zhang et al., 2021a).<sup>7</sup> Our results on the DGM corpus suggest that our linearization schema effectively models  $n$ -ary relations without requiring changes to the model architecture or training procedure.

### 5.1.3 End-to-end methods

In Table 3 we compare to an existing end-to-end approach on DocRED, JEREX (Eberts and Ulges, 2021). To make the comparison fair, we use the same pretrained encoder (BERT<sub>BASE</sub>). We find that although our model is arguably simpler (JEREX contains four task-specific sub-components, each with its own loss) it only slightly underperforms JEREX, mainly due to recall. We speculate that one reason for this is a large number of relations per document, which leads to longer target strings and, therefore, more decoding steps. The median length of the target strings in DocRED, using our linearization, is 110, whereas the next largest is 19 in GDA. Improving the decoder’s ability to process long sequences, e.g. switching the LSTM for a transformer or modifying the linearization schema to produce shorter target strings, may improve recall and close the gap with existing methods.

## 5.2 Pipeline vs. End-to-end

In §5.1.1 and §5.1.2, we provide gold-standard entity annotations from each corpus as input to

<sup>6</sup><https://www.ncbi.nlm.nih.gov/pmc/>

<sup>7</sup>The authors have not released code at the time of writing, so we were unable to evaluate this model on the DGM validation set in order to compare with our method directly.

Table 3: Comparison to existing end-to-end methods. Performance reported as micro-precision, recall and F1-scores (%) on the DocRED test set. Results below the horizontal line are not comparable to existing methods. Bold: best scores.

Method	P	R	F1
JEREX (Eberts and Ulges, 2021)	42.8	<b>38.2</b>	<b>40.4</b>
seq2rel (end-to-end)	<b>44.0</b>	33.8	38.2
seq2rel (end-to-end, relaxed)	53.7	41.3	46.7

Table 4: Comparison of pipeline-based and end-to-end approaches. Gold hints use gold-standard entity annotations to insert entity hints in the source text. Silver hints use the entity annotations provided by PubTator. Pipeline is identical to silver entity hints, except that we filter out entity mentions predicted by our model that PubTator does not predict. The end-to-end model only has access to the unannotated source text as input. Performance reported as micro-precision, recall and F1-scores (%) on the CDR test set, with strict and relaxed entity matching criteria. Bold: best scores.

	Strict			Relaxed		
	P	R	F1	P	R	F1
Gold hints	68.2	66.2	67.2	68.2	66.2	67.2
Silver hints	42.4	37.3	39.7	53.0	46.7	49.7
Pipeline	<b>45.0</b>	16.9	24.6	<b>62.5</b>	23.5	34.1
End-to-end	43.5	<b>37.5</b>	<b>40.2</b>	56.6	<b>48.8</b>	<b>52.4</b>

our model via entity hinting (referred to as “gold” hints from here on, see §3.3). This allowed us to compare to existing methods that also provide these annotations as input. However, gold-standard entity annotations are (almost) never available in real-world settings, such as large-scale extraction on PubMed. In this setting, there are two strategies: pipeline-based, where independent systems perform entity and relation extraction, and end-to-end, where a single model performs both tasks. To compare these approaches under our model, we perform evaluations where a named entity recognition (NER) system is used to determine entity hints (referred to as “silver” hints from here on) and when no entity hints are provided (end-to-end).<sup>8</sup> However, this alone does not create a true pipeline, as our model can recover from both false negatives and false positives in the NER step. To mimic error propagation in the pipeline setting, we filter any entity mention predicted by our model that was *not* predicted by the NER system. In Table 4, we

<sup>8</sup>Specifically, we use PubTator (Wei et al., 2013). PubTator provides up-to-date entity annotations for PubMed using state-of-the-art machine learning systems.

Table 5: Ablation study results. Performance reported as the micro-precision, recall and F1-scores (%) on the CDR and DocRED validation sets.  $\Delta$ : difference to the complete models F1-score. Bold: best scores.

	CDR				DocRED			
	P	R	F1	$\Delta$	P	R	F1	$\Delta$
seq2rel (end-to-end)	<b>41.0</b>	35.1	37.8	-	46.9	<b>36.1</b>	<b>40.8</b>	-
- pretraining	9.4	6.9	8.0	-29.8	18.5	7.7	10.8	-30.0
- fine-tuning	24.3	20.5	22.2	-15.6	42.4	15.5	22.7	-18.1
- vocab restriction	39.6	32.2	35.5	-2.3	45.2	35.5	39.7	-1.1
- sorting relations	36.1	29.2	32.3	-5.6	<b>52.9</b>	17.4	26.2	-14.7
+ constrained decoding	40.8	<b>35.6</b>	<b>38.0</b>	+0.2	46.8	35.9	40.6	-0.2

present the results of all four settings (gold and silver entity hints, pipeline and end-to-end) on CDR.

We find that using gold entity hints significantly outperforms all other settings. This is expected, as the gold-standard entity annotations are high-quality labels produced by domain experts. Using silver hints significantly drops performance, likely due to a combination of false positive and false negatives from the NER step. In the pipeline setting, where there is no recovery from false negatives, performance falls by another 15%. The end-to-end setting significantly outperforms the pipeline setting (due to a large boost in recall) and performs comparably to using silver hints. Together, our results suggest that performance reported using gold-standard entity annotations may be overly optimistic and corroborates previous work demonstrating the benefits of jointly learning entity and relation extraction (Miwa and Sasaki, 2014; Miwa and Bansal, 2016; Gupta et al., 2016; Li et al., 2016a, 2017; Nguyen and Verspoor, 2019a; Yu et al., 2020).

### 5.3 Ablation

In Table 5, we present the results of an ablation study. We perform the analysis twice, once on the biomedical corpus CDR and once on the general domain corpus DocRED. Unsurprisingly, we find that fine-tuning a pretrained encoder greatly impacts performance. Training the same encoder from scratch (- pretraining) reduces performance by  $\sim 30\%$ . Using the pretrained weights without fine-tuning (- fine-tuning) drops performance by 15.6-18.1%. Restricting the target vocabulary (- vocab restriction, see §3.2.1) has a small positive impact, boosting performance by 1.1%-2.3%. Deliberately ordering the relations within each target string (- sorting relations, see §3.2.4) has a large positive impact, boosting performance by 5.6%-14.7%. This effect is larger on DocRED, likely because it has more relations per document on average than CDR, so ordering becomes more impor-

tant. Finally, adding constraints to the decoding process (+ constrained decoding) has little impact on performance, suggesting that a trained model rarely generates invalid target strings (see §3.2.3).

## 6 Discussion

### 6.1 Related work

Seq2seq learning for RE has been explored in prior work. CopyRE (Zeng et al., 2018) uses an encoder-decoder architecture with a copy mechanism, similar to our approach, but is restricted to intra-sentence relations. Additionally, because CopyRE’s decoding proceeds for exactly three timesteps per relation, the model is limited to generating binary relations between single token entities. The ability to decode multi-token entities was addressed in follow-up work, CopyMTL (Zeng et al., 2020). A similar approach was published concurrently but was again limited to intra-sentence binary relations (Nayak and Ng, 2020). Most recently, GenerativeRE (Cao and Ananiadou, 2021) proposed a novel copy mechanism to improve performance on multi-token entities. None of these approaches deal with the complexities of DocRE, where many relations cross the sentence boundary, and coreference resolution is critical.<sup>9</sup>

More generally, our paper is related to a recently proposed “text-to-text” framework (Raffel et al., 2020). In this framework, a task is formulated so that the inputs and outputs are both text strings, enabling the use of the same model, loss function and even hyperparameters across many seq2seq, classification and regression tasks. This framework has recently been applied to biomedical literature to perform named entity recognition, relation extraction (binary, intra-sentence), natural language inference, and question answering (Phan et al., 2021). Our work can be seen as an attempt to formulate the task of DocRE within this framework.

### 6.2 Limitations and future work

**Permutation-sensitive loss** Our approach adopts the sequence cross-entropy loss (Equation 2), which is sensitive to the order of predicted tokens, enforcing an unnecessary decoding order on the inherently unordered relations. To partially mitigate this problem, we order relations within the

<sup>9</sup>Concurrent to our work, REBEL (Huguet Cabot and Navigli, 2021) also extends seq2seq methods to document-level RE, achieving strong performance on DocRED. However, the method was not evaluated on  $n$ -ary relations.



target string according to order of appearance in the source text, providing the model with a consistent decoding order that can be learned (see §3.2.4, §5.3). Previous work has addressed this issue with various strategies, including reinforcement learning (Zeng et al., 2019), unordered-multi-tree decoders (Zhang et al., 2020), and non-autoregressive decoders (Sui et al., 2020). However, these works are limited to binary intra-sentence relation extraction, and their suitability for DocRE has not been explored. A promising future direction would be to modify our approach such that the arbitrary order of relations is not enforced during training.

**Input length restriction** Due to the pretrained encoder’s input size limit (512 tokens), our experiments are conducted on paragraph-length text. Our model could be extended to full documents by swapping its encoder with any of the recently proposed “efficient transformers” (Tay et al., 2021). Future work could evaluate such a model’s ability to extract relations from full scientific papers.

**Pretraining the decoder** In our model, the encoder is pretrained, while the decoder is trained from scratch. Several recent works, such as T5 (Raffel et al., 2020) and BART (Lewis et al., 2020), have proposed pretraining strategies for entire encoder-decoder architectures, which can be fine-tuned on downstream tasks. An interesting future direction would be to fine-tune such a model on DocRE using our linearization schema.

## 7 Conclusion

In this paper, we extend generative, seq2seq methods for relation extraction to the document level. We propose a novel linearization schema that can handle complexities overlooked by previous seq2seq approaches, like coreferent mentions and  $n$ -ary relations. We compare our approach to existing pipeline-based and end-to-end methods on several benchmark corpora, in some cases exceeding their performance. In future work, we hope to extend our method to full scientific papers and develop strategies to improve performance in the low-data regime and in cases where there are many relations per document.

## Acknowledgements

This research was enabled in part by support provided by Compute Ontario ([www.computeontario.ca](http://www.computeontario.ca)), Compute Canada

([www.computecanada.ca](http://www.computecanada.ca)) and the CIFAR AI Chairs Program and partially funded by the US National Institutes of Health (NIH) [U41 HG006623, U41 HG003751].

## References

- Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. 2019. **Optuna: A next-generation hyperparameter optimization framework**. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2019, Anchorage, AK, USA, August 4-8, 2019*, pages 2623–2631. ACM.
- Giannis Bekoulis, Johannes Deleu, Thomas Demeester, and Chris Develder. 2018. Joint entity recognition and relation extraction as a multi-head selection problem. *Expert Systems with Applications*, 114:34–45.
- James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. 2011. **Algorithms for hyper-parameter optimization**. In *Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011. Proceedings of a meeting held 12-14 December 2011, Granada, Spain*, pages 2546–2554.
- Jiarun Cao and Sophia Ananiadou. 2021. **GenerativeRE: Incorporating a novel copy mechanism and pretrained model for joint entity and relation extraction**. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2119–2126, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Fenia Christopoulou, Makoto Miwa, and Sophia Ananiadou. 2019. **Connecting the dots: Document-level neural relation extraction with edge-oriented graphs**. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4925–4936, Hong Kong, China. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. **BERT: Pre-training of deep bidirectional transformers for language understanding**. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Markus Eberts and Adrian Ulges. 2021. **An end-to-end model for entity-level relation extraction using multi-instance learning**. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 3650–3660, Online. Association for Computational Linguistics.

- Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke Zettlemoyer. 2018. [AllenNLP: A deep semantic natural language processing platform](#). In *Proceedings of Workshop for NLP Open Source Software (NLP-OSS)*, pages 1–6, Melbourne, Australia. Association for Computational Linguistics.
- John Giorgi, Xindi Wang, Nicola Sahar, Won Young Shin, Gary D Bader, and Bo Wang. 2019. [End-to-end named entity recognition and relation extraction using pre-trained language models](#). *ArXiv preprint*, abs/1912.13415.
- Alex Graves. 2012. [Sequence transduction with recurrent neural networks](#). *ArXiv preprint*, abs/1211.3711.
- Jiatao Gu, Zhengdong Lu, Hang Li, and Victor O.K. Li. 2016a. [Incorporating copying mechanism in sequence-to-sequence learning](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1631–1640, Berlin, Germany. Association for Computational Linguistics.
- Jinghang Gu, Longhua Qian, and Guodong Zhou. 2016b. Chemical-induced disease relation extraction with various linguistic features. *Database: The Journal of Biological Databases and Curation*, 2016.
- Jinghang Gu, Fuqing Sun, Longhua Qian, and Guodong Zhou. 2017. Chemical-induced disease relation extraction via convolutional neural network. *Database: The Journal of Biological Databases and Curation*, 2017.
- Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2020. [Domain-specific language model pretraining for biomedical natural language processing](#). *ArXiv preprint*, abs/2007.15779.
- Pankaj Gupta, Hinrich Schütze, and Bernt Andrassy. 2016. [Table filling multi-task recurrent neural network for joint entity and relation extraction](#). In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 2537–2547, Osaka, Japan. The COLING 2016 Organizing Committee.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural Computation*, 9:1735–1780.
- Pere-Lluís Huguet Cabot and Roberto Navigli. 2021. [REBEL: Relation extraction by end-to-end language generation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2370–2381, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Sarthak Jain, Madeleine van Zuylen, Hannaneh Hajishirzi, and Iz Beltagy. 2020. [SciREX: A challenge dataset for document-level information extraction](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7506–7516, Online. Association for Computational Linguistics.
- Robin Jia, Cliff Wong, and Hoifung Poon. 2019. [Document-level n-ary relation extraction with multi-scale representation learning](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3693–3704, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jin-Dong Kim, T. Ohta, Yuka Tateisi, and Jun’ichi Tsujii. 2003. Genia corpus - a semantically annotated corpus for bio-textmining. *Bioinformatics*, 19 Suppl 1:i180–2.
- Po-Ting Lai and Zhiyong Lu. 2021. Bert-gt: Cross-sentence n-ary relation extraction with bert and graph transformer. *Bioinformatics*.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Fei Li, Meishan Zhang, Guohong Fu, and Donghong Ji. 2017. A neural joint model for entity and relation extraction from biomedical text. *BMC bioinformatics*, 18(1):198.
- Fei Li, Yue Zhang, Meishan Zhang, and Donghong Ji. 2016a. [Joint models for extracting adverse drug events from biomedical text](#). In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI 2016, New York, NY, USA, 9-15 July 2016*, pages 2838–2844. IJCAI/AAAI Press.
- Jiao Li, Yueping Sun, Robin J. Johnson, Daniela Sciak, Chih-Hsuan Wei, Robert Leaman, Allan Peter Davis, Carolyn J. Mattingly, Thomas C. Wieggers, and Zhiyong Lu. 2016b. [Biocreative v cdr task corpus: a resource for chemical disease relation extraction](#). *ArXiv preprint*, abs/d.
- Feifan Liu, Jinying Chen, Abhyuday Jagannatha, and Hong Yu. 2016. [Learning for biomedical information extraction: Methodological review of recent advances](#). *ArXiv preprint*, abs/1606.07993.
- Ilya Loshchilov and Frank Hutter. 2019. [Decoupled weight decay regularization](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh Hajishirzi. 2018. [Multi-task identification of entities, relations, and coreference for scientific knowledge](#)

- [graph construction](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3219–3232, Brussels, Belgium. Association for Computational Linguistics.
- Hieu Minh Tran, Minh Trung Nguyen, and Thien Huu Nguyen. 2020. [The dots have their values: Exploiting the node-edge connections in graph-based neural models for document-level relation extraction](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4561–4567, Online. Association for Computational Linguistics.
- Makoto Miwa and Mohit Bansal. 2016. [End-to-end relation extraction using LSTMs on sequences and tree structures](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1105–1116, Berlin, Germany. Association for Computational Linguistics.
- Makoto Miwa and Yutaka Sasaki. 2014. [Modeling joint entity and relation extraction with table representation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1858–1869, Doha, Qatar. Association for Computational Linguistics.
- 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, Online. Association for Computational Linguistics.
- Tapas Nayak and Hwee Tou Ng. 2020. [Effective modeling of encoder-decoder architecture for joint entity and relation extraction](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 8528–8535. AAAI Press.
- Dat Quoc Nguyen and Karin Verspoor. 2019a. [End-to-end neural relation extraction using deep biaffine attention](#). In *Advances in Information Retrieval*, pages 729–738, Cham. Springer, Springer International Publishing.
- Dat Quoc Nguyen and Karin Verspoor. 2019b. [End-to-end neural relation extraction using deep biaffine attention](#). In *European Conference on Information Retrieval*, pages 729–738. Springer.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*.
- 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.
- Long N Phan, James T Anibal, Hieu Tran, Shaurya Chanana, Erol Bahadroglu, Alec Peltekian, and Grégoire Altan-Bonnet. 2021. [Scifive: a text-to-text transformer model for biomedical literature](#). *ArXiv preprint*, abs/2106.03598.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21:1–67.
- Linfeng Song, Yue Zhang, Zhiguo Wang, and Daniel Gildea. 2018. [N-ary relation extraction using graph-state LSTM](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2226–2235, Brussels, Belgium. Association for Computational Linguistics.
- Nitish Srivastava, Geoffrey E. Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.*, 15:1929–1958.
- Dianbo Sui, Yubo Chen, Kang Liu, Jun Zhao, Xiangrong Zeng, and Shengping Liu. 2020. [Joint entity and relation extraction with set prediction networks](#). *ArXiv preprint*, abs/2011.01675.
- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. [Sequence to sequence learning with neural networks](#). In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 3104–3112.
- Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. 2021. [Long range arena : A benchmark for efficient transformers](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- 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.
- Patrick Verga, Emma Strubell, and Andrew McCallum. 2018. [Simultaneously self-attending to all mentions for full-abstract biological relation extraction](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 872–884, New Orleans, Louisiana. Association for Computational Linguistics.

- Oriol Vinyals, Samy Bengio, and Manjunath Kudlur. 2016. [Order matters: Sequence to sequence for sets](#). In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*.
- David Wadden, Ulme Wennberg, Yi Luan, and Hananeh Hajishirzi. 2019. [Entity, relation, and event extraction with contextualized span representations](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5784–5789, Hong Kong, China. Association for Computational Linguistics.
- Li Wan, Matthew D. Zeiler, Sixin Zhang, Yann LeCun, and Rob Fergus. 2013. [Regularization of neural networks using dropconnect](#). In *Proceedings of the 30th International Conference on Machine Learning, ICML 2013, Atlanta, GA, USA, 16-21 June 2013*, volume 28 of *JMLR Workshop and Conference Proceedings*, pages 1058–1066. JMLR.org.
- Yucheng Wang, Bowen Yu, Hongsong Zhu, Tingwen Liu, Nan Yu, and Limin Sun. 2021. [Discontinuous named entity recognition as maximal clique discovery](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 764–774, Online. Association for Computational Linguistics.
- Chih-Hsuan Wei, Hung-Yu Kao, and Zhiyong Lu. 2013. Pubtator: a web-based text mining tool for assisting biocuration. *Nucleic acids research*, 41(W1):W518–W522.
- 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. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- 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*, pages 272–284, Cham. Springer International Publishing.
- Benfeng Xu, Quan Wang, Yajuan Lyu, Yong Zhu, and Zhendong Mao. 2021. Entity structure within and throughout: Modeling mention dependencies for document-level relation extraction. In *AAAI*.
- Pengcheng Yang, Fuli Luo, Shuming Ma, Junyang Lin, and Xu Sun. 2019. [A deep reinforced sequence-to-set model for multi-label classification](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5252–5258, Florence, Italy. Association for Computational Linguistics.
- 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, Florence, Italy. Association for Computational Linguistics.
- Bowen Yu, Zhenyu Zhang, Xiaobo Shu, Tingwen Liu, Yubin Wang, Bin Wang, and Sujian Li. 2020. Joint extraction of entities and relations based on a novel decomposition strategy. In *ECAI 2020*, pages 2282–2289. IOS Press.
- Daojian Zeng, Haoran Zhang, and Qianying Liu. 2020. [Copymtl: Copy mechanism for joint extraction of entities and relations with multi-task learning](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 9507–9514. AAAI Press.
- Xiangrong Zeng, Shizhu He, Daojian Zeng, Kang Liu, Shengping Liu, and Jun Zhao. 2019. [Learning the extraction order of multiple relational facts in a sentence with reinforcement learning](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 367–377, Hong Kong, China. Association for Computational Linguistics.
- Xiangrong Zeng, Daojian Zeng, Shizhu He, Kang Liu, and Jun Zhao. 2018. [Extracting relational facts by an end-to-end neural model with copy mechanism](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 506–514, Melbourne, Australia. Association for Computational Linguistics.
- Ranran Haoran Zhang, Qianying Liu, Aysa Xuemo Fan, Heng Ji, Daojian Zeng, Fei Cheng, Daisuke Kawahara, and Sadao Kurohashi. 2020. [Minimize exposure bias of Seq2Seq models in joint entity and relation extraction](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 236–246, Online. Association for Computational Linguistics.
- Sheng Zhang, Cliff Wong, Naoto Usuyama, Sarthak Jain, Tristan Naumann, and Hoifung Poon. 2021a. [Modular self-supervision for document-level relation extraction](#). In *Proceedings of the 2021 Conference*

on *Empirical Methods in Natural Language Processing*, pages 5291–5302, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Tianyi Zhang, Felix Wu, Arzoo Katiyar, Kilian Q. Weinberger, and Yoav Artzi. 2021b. [Revisiting few-sample BERT fine-tuning](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.

Deyu Zhou, Dayou Zhong, and Yulan He. 2014. Biomedical relation extraction: from binary to complex. *Computational and mathematical methods in medicine*, 2014.

Wenxuan Zhou and Muhao Chen. 2021. An improved baseline for sentence-level relation extraction. *ArXiv*, abs/2102.01373.

Wenxuan Zhou, Kevin Huang, Tengyu Ma, and Jing Huang. 2021. Document-level relation extraction with adaptive thresholding and localized context pooling. In *AAAI*.

## A Constrained decoding

In [Figure 4](#), we illustrate the rules used to constrain decoding. At each timestep  $t$ , given the prediction of the previous timestep  $t - 1$ , the predicted class probabilities of tokens that would generate a syntactically invalid target string are set to a tiny value. In practice, we found that a model rarely generates invalid target strings, so these constraints have little effect on final performance (see [§3.2.3](#) and [§5.3](#)).

## B Details about dataset annotations

In [Table 6](#), we list which complexities (e.g. nested & discontinuous mentions,  $n$ -ary relations) are contained within each dataset used in our evaluations. We also report the fraction of relations in the test set that are inter-sentence. We consider a relation intra-sentence if *any* sentence in the document contains *at least one* mention of each entity in the relation, and inter-sentence otherwise. This produces an estimate that matches previously reported numbers for CDR ( $\sim 30\%$ ). In [Yao et al. \(2019\)](#), the fraction of inter-sentence relations in DocRED is reported as  $\sim 40.7\%$ . We can reproduce this value if we consider relations intra-sentence when *all* mentions of an entity exist within a single sentence and inter-sentence otherwise.

## C Hypernym filtering

The CDR dataset is annotated for chemical-induced disease (CID) relationships between the most

specific chemical and disease mentions in an abstract. Take the following example from the corpus:

**Carbamazepine**-induced **cardiac dysfunction** [...] A patient with sinus **bradycardia** and **atrioventricular block**, induced by **carbamazepine**, prompted an extensive literature review of all previously reported cases.

In this example (PMID: 1728915), only (*carbamazepine, bradycardia*) and (*carbamazepine, atrioventricular block*) are labelled as true relations. The relation (*carbamazepine, cardiac dysfunction*), although true, is not labelled as *cardiac dysfunction* is a hypernym of both *bradycardia* and *atrioventricular block*. This can harm evaluation performance, as the prediction (*carbamazepine, cardiac dysfunction*) will be considered a false positive. Therefore, we follow previous work ([Gu et al., 2016b, 2017](#); [Verga et al., 2018](#); [Christopoulou et al., 2019](#); [Zhou et al., 2021](#)) by filtering negative relations like these, with disease entities that are hypernyms of a corresponding true relations disease entity within the same abstract, according to the hierarchy in the MeSH vocabulary.<sup>10</sup>

## D Parsing the models output

At test time, our model autoregressively generates an output, token-by-token, using beam search decoding (see [§3.2](#)). In order to extract the predicted relations from this output, we apply the following steps. First, predicted token ids are converted to a string. We use the `decode()`<sup>11</sup> method of the HuggingFace Transformers tokenizer ([Wolf et al., 2020](#)) to do this. For example, after calling `decode()` on the predicted token ids, this string might look like:

```
monoamine oxidase b ; maob @GENE@ parkinson's  
disease ; pd @DISEASE@ @GDA@
```

We then use regular expressions to extract any relations from this string that match our linearization schema (see [§3.1](#)), which produces a dictionary of nested lists, keyed by relation class:

```
{  
  "GDA": [  
    [  

```

<sup>10</sup><https://meshb.nlm.nih.gov>

<sup>11</sup>[https://huggingface.co/docs/transformers/main\\_classes/tokenizer#transformers.PreTrainedTokenizerBase.decode](https://huggingface.co/docs/transformers/main_classes/tokenizer#transformers.PreTrainedTokenizerBase.decode)

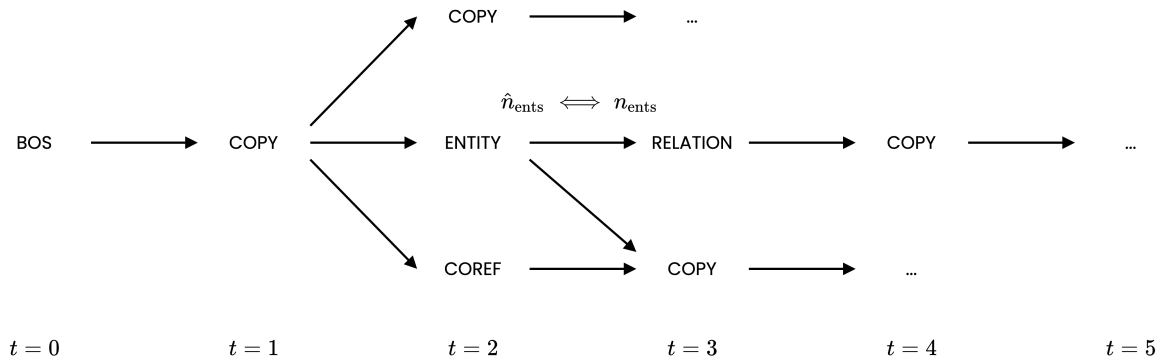


Figure 4: A diagram depicting syntactically valid predictions during decoding at each timestep  $t$ . The log probabilities of all other possible predictions are set to a tiny value to prevent the model from producing a syntactically invalid target string. BOS is the special beginning-of-sequence token, COPY denotes any token copied from the source text, and COREF is the special token used to separate coreferent mentions (i.e. ;). ENTITY is any special entity token (e.g. @GENE@) and RELATION any special relation token (e.g. @GDA@ for gene-disease association).  $\hat{n}_{\text{ents}}$  denotes the number of entities predicted by the current timestep and  $n_{\text{ents}}$  the expected arity of the relation. The special end-of-sequence token (not shown) is always considered valid and its log probability is never modified.

Table 6: Evaluation datasets used in this paper with details about their annotations. Inter-sentence relations (%) are the fraction of relations in the test set that cross sentence boundaries. We consider a relation intra-sentence if any sentence in the document contains at least one mention of each entity in the relation, and inter-sentence otherwise. \*This differs from the estimate in Yao et al. (2019), see Appendix B.

Corpus	Nested Mentions?	Discontinuous Mentions?	Coreferent mentions?	$n$ -ary relations?	Inter-sentence relations (%)
CDR (Li et al., 2016b)	✓	✓	✓	✗	29.8
GDA (Wu et al., 2019)	✓	✗	✓	✗	15.6
DGM (Jia et al., 2019)	✗	✗	✓	✓	63.5
DocRED (Yao et al., 2019)	✗	✗	✓	✗	12.5*

```

    [ ["monoamine oxidase b", "maob"], "GENE"],
    [ ["parkinson's disease", "pd"], "DISEASE"]
  ]
}

```

Finally, we apply some normalization steps to the entity mentions. Namely, we strip leading and trailing white space characters, sort entity mentions lexicographically (as their order is not important), and remove duplicate mentions. Similarly, we remove duplicate relations. These steps are applied to both target and model output strings. The F1-score can then be computed by tallying true positives, false positives and false negatives.

## E Relaxed entity matching

The aim of DocRE is to extract relations at the *entity*-level. However, it is common to evaluate these methods with a “strict” matching criteria, where a predicted entity  $\mathcal{P}$  is considered correct if and only if all its *mentions* exactly match a corresponding gold entities mentions, i.e.  $\mathcal{P} = \mathcal{G}$ . This penalizes model predictions that miss even a single coreferent mention, but are otherwise correct. A relaxed

criteria, proposed in prior work (Jain et al., 2020) considers  $\mathcal{P}$  to match  $\mathcal{G}$  if more than 50% of  $\mathcal{P}$ ’s mentions belong to  $\mathcal{G}$ , that is

$$\frac{|\mathcal{P} \cap \mathcal{G}|}{|\mathcal{P}|} > 0.5$$

In this paper, alongside the strict criteria, we report performance using this relaxed entity matching strategy, denoted “relaxed”.

## F Hyperparameters

In Table 7, we list the hyperparameter values used during evaluation on each corpus, with and without entity hinting. Select hyperparameters were tuned using Optuna (Akiba et al., 2019). The tuning process selects the best hyperparameters according to the validation set micro F1-score using the TPE (Tree-structured Parzen Estimator) algorithm (Bergstra et al., 2011).<sup>12</sup> During tuning, we use greedy decoding (i.e. beam size of one). Once opti-

<sup>12</sup><https://optuna.readthedocs.io/en/stable/reference/generated/optuna.samplers.TPESampler.html>

Table 7: Hyperparameter values used for each corpus. Hyperparameters values when using entity hinting, if they differ from the values used without entity hinting, are shown in parentheses. Tuned indicates whether or not the hyperparameters were tuned on the validation sets.

Hyperparameter	Tuned?	CDR	GDA	DGM	DocRED
Batch size	✓	4	4	4	4
Training epochs	✓	130 (70)	30 (25)	30 (45)	50
Encoder learning rate	✗	2e-5	2e-5	2e-5	2e-5
Encoder weight decay	✗	0.01	0.01	0.01	0.01
Encoder re-initialized top $L$ layers	✓	1	1 (2)	1	1
Decoder learning rate	✓	1.21e-4 (1.13e-4)	5e-4 (4e-4)	8e-4 (1.5e-5)	7.8e-5
Decoder input dropout	✗	0.1	0.1	0.1	0.1
Decoder hidden-to-hidden weights dropout	✗	0.5	0.5	0.5	0.5
Target embedding size	✗	256	256	256	256
No. heads in multi-head cross-attention	✗	6	6	6	6
Beam size	✓	3 (2)	4 (1)	3 (2)	8
Length penalty	✓	1.4 (0.2)	0.8 (1.0)	0.2 (0.8)	1.4
Max decoding steps	✗	128	96	96	400

mal hyperparameters are found, we tune the beam size (bs) and length penalty ( $\alpha$ ) using a grid search over the values  $bs = \{2...10\}$ , with a step size of 1, and  $\alpha = \{0.2...2.0\}$ , with a step size of 0.2.

## G Baselines

This section contains detailed descriptions of all methods we compare to in this paper.

### G.1 Pipeline-based methods

These methods are pipeline-based, assuming the entities are provided as input. Many of them construct a document-level graph using dependency parsing, heuristics, or structured attention and then update node and edge representations using propagation.

- [Christopoulou et al. \(2019\)](#) propose EoG, an edge-orientated graph neural model. The nodes of the graph are constructed from mentions, entities, and sentences. Edges between nodes are initially constructed using heuristics. An iterative algorithm is then used to generate edges between nodes in the graph. Finally, a classification layer takes the representation of entity-to-entity edges as input to determine whether those entities express a relation or not. We compare to EoG in the pipeline-based setting on the CDR and GDA corpora.
- [Nan et al. \(2020\)](#) propose LSR (Latent Structure Refinement). A “node constructor” encodes each sentence of an input document and outputs contextual representations. Representations that correspond to mentions and tokens on the shortest dependency path in a sentence

are extracted as nodes. A “dynamic reasoner” is then applied to induce a document-level graph based on the extracted nodes. The classifier uses the final representations of nodes for relation classification. We compare to LSR in the pipeline-based setting on the CDR and GDA corpora.

- [Lai and Lu \(2021\)](#) propose BERT-GT, which combines BERT with a graph transformer. Both BERT and the graph transformer accept the document text as input, but the graph transformer requires the neighbouring positions for each token, and the self-attention mechanism is replaced with a neighbour-attention mechanism. The hidden states of the two transformers are aggregated before classification. We compare to BERT-GT in the pipeline-based setting on the CDR and GDA corpora.
- [Minh Tran et al. \(2020\)](#) propose EoGANE (EoG model Augmented with Node Representations), which extends the edge-orientated model proposed by [Christopoulou et al. \(2019\)](#) to include explicit node representations which are used during relation classification. We compare to EoGANE in the pipeline-based setting on the CDR and GDA corpora.
- [SSAN \(Xu et al., 2021\)](#) propose SSAN (Structured Self-Attention Network), which inherits the architecture of the transformer encoder ([Vaswani et al., 2017](#)) but adds a novel structured self-attention mechanism to model the coreference and co-occurrence dependencies between an entities mentions. We compare

to SSAN in the pipeline-based setting on the CDR and GDA corpora.

- [Zhou et al. \(2021\)](#) propose ALTOP (Adaptive Thresholding and Localized cOntext Pooling), which extends BERT with two modifications. Adaptive thresholding, which learns an optimal threshold to apply to the relation classifier. Localized context pooling, which uses the pre-trained self-attention layers of BERT to create an entity embedding from its mentions and their context. We compare to ALTOP in the pipeline-based setting on the CDR and GDA corpora.

## G.2 $n$ -ary relation extraction

These methods are explicitly designed for the extraction of  $n$ -ary relations, where  $n > 2$ .

- [Jia et al. \(2019\)](#) propose a multiscale neural architecture, which combines representations learned over text spans of varying scales and for various sub-relations. We compare to [Jia et al. \(2019\)](#) in the pipeline-based setting on the  $n$ -ary DGM corpus.

## G.3 End-to-end methods

These methods are capable of performing the sub-tasks of DocRE in an end-to-end fashion with only the document text as input.

- [Eberts and Ulges \(2021\)](#) propose JEREX, which extends BERT with four task-specific components that use BERTs outputs to perform entity mention localization, coreference resolution, entity classification, and relation classification. They present two versions of their relation classifier, denoted “global relation classifier” (GRC) and “multi-instance relation classifier” (MRC). We compare to JEREX-MRC in the end-to-end setting on the DocRED corpus.