

# An Interpretable Generative Adversarial Approach to Classification of Latent Entity Relations in Unstructured Sentences

Shiou Tian Hsu, Changsung Moon, Paul Jones, Nagiza F. Samatova

North Carolina State University, Raleigh, NC, 27606

{shsu3,cmoon2,pjones}@ncsu.edu, samatova@csc.ncsu.edu

## Abstract

We propose a generative adversarial neural network model for relation classification that attempts to emulate the way in which human analysts might process sentences. Our approach provides two unique benefits over existing capabilities: (1) we make predictions by finding and exploiting supportive rationales to improve interpretability (i.e. words or phrases extracted from a sentence that a person can reason upon), and (2) we allow predictions to be easily corrected by adjusting the rationales. Our model consists of three stages: Generator, Selector, and Encoder. The Generator identifies candidate text fragments; the Selector decides which fragments can be used as rationales depending on the goal; and finally, the Encoder performs relation reasoning on the rationales. While the Encoder is trained in a supervised manner to classify relations, the Generator and Selector are designed as unsupervised models to identify rationales without prior knowledge, although they can be semi-supervised through human annotations. We evaluate our model on data from SemEval 2010 that provides 19 relation-classes. Experiments demonstrate that our approach outperforms state-of-the-art models, and that our model is capable of extracting good rationales on its own as well as benefiting from labeled rationales if provided.

## Introduction

Sentence-level entity relation classification is the task of recognizing the relationship between two entities within a sentence. It is the core function of autonomous knowledge discovery in unstructured text (such as web documents), and thus is critical for applications including information extraction, knowledge base construction, and many other higher level NLP tasks (?; ?). For example, given the following sentence:

*He had chest pains and [headaches] $e_1$  from [mold] $e_2$  in the bedrooms.*

with marked target entity  $e_1$  = “headaches” and  $e_2$  = “mold”, the goal would be to correctly identify that this sentence expresses a casual relationship from  $e_2$  to  $e_1$ , for which we use the notation Cause-Effect( $e_2, e_1$ ).

The field of entity-relation classification has greatly benefited from recent developments in neural network

models, and these benefits have propagated into many other important applications like document classification, sentiment analysis, summarization, topic discovery, word-sense disambiguation, co-reference identification and more. Initially neural models pushed classification performance to new heights, improving F1 score on relation classification to 88% (?) compared to 82.2% using an SVM with engineered features (?). However, neural models often seem to perform incomprehensible operations on the sentence to reason desired knowledge, whereas human readers usually derive relations between entities by locating key indicating words. For instance, given the following sentence, a human will likely focus only on a few keywords that explain the core relationships, such as the Component-Whole( $e_1, e_2$ ) relationship, whereas a neural model will leverage every word in the sentence based on computed scores:

**Example sentence:** *The disgusting scene was retaliation against her brother Philip who rents the [room] $e_1$  inside this apartment [house] $e_2$  on Lombard street.*

**Human interpretation:** *The disgusting scene was retaliation against her brother Philip who rents the [room] $e_1$  inside this apartment [house] $e_2$  on Lombard street.*

**Attention-based CNN models:** *The<sub>0.02</sub> disgusting<sub>0.03</sub> scene<sub>0.1</sub> was<sub>0.02</sub> retaliation<sub>0.03</sub> against<sub>0.07</sub> her<sub>0.015</sub> brother<sub>0.02</sub> Philip<sub>0.04</sub> who<sub>0.08</sub> rents<sub>0.03</sub> the<sub>0.1</sub> [room] $e_1$  inside<sub>0.13</sub> this<sub>0.03</sub> apartment<sub>0.09</sub> [house] $e_2$  on<sub>0.01</sub> Lombard<sub>0.03</sub> street<sub>0.08</sub>.*

**(value behind the words stands for the derived importance of the word for determining the relation)**

In this research, we aim to improve interpretability of neural models by imitating the logic used by human annotators through a three-stage Generator-Selector-Encoder approach. The Generator first enumerates candidate text fragments from the sentence using a Recurrent Neural Network (RNN) model, then the Selector extracts fragments that are informative for determining the relation and passes these on to the Encoder to make final relation predictions. We refer to the text fragments as **rationales** as per the definition in (?; ?; ?). Some good examples of rationales are like “inside” in the above example sentence, or “caused”, “resulted in” and “leads to” for entities of Cause-Effect relationship.

However, under most situations, a heavily annotated

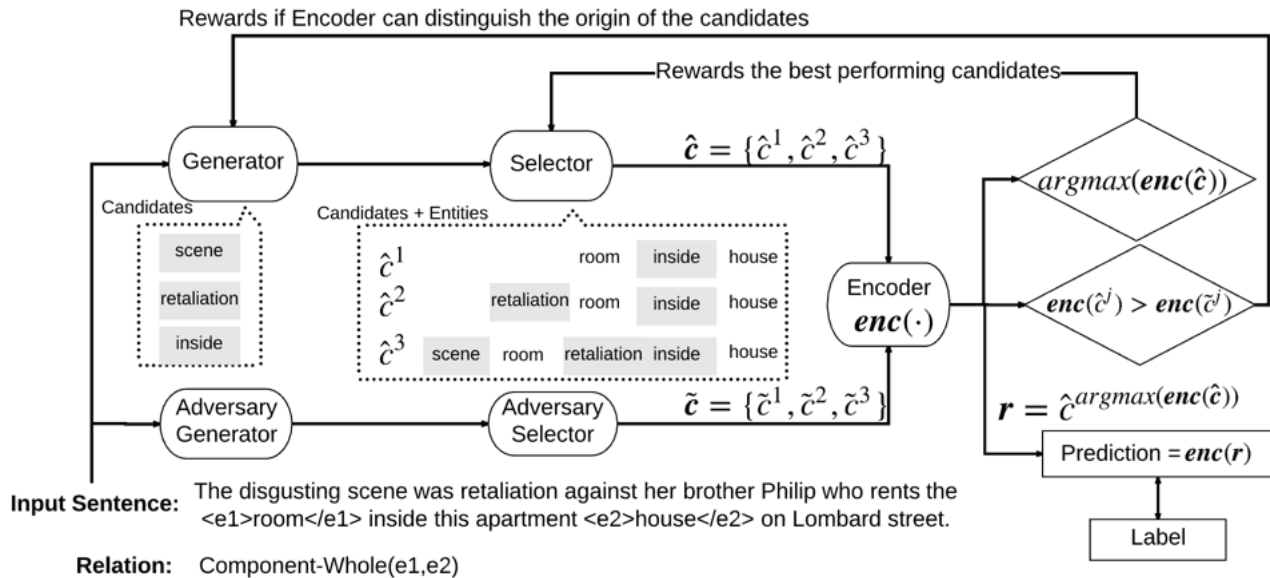


Figure 1: Model Structure with sample input sentence. Entities are given from the dataset. The Generator selects candidate rationales, and the Selector enumerates all possible combinations of candidates with entities and selects one ‘best performing candidate set’. An adversary Generator and Selector carry out the same process using a randomized approach. Candidate sets are then transformed into a vector representation by the Encoder and are evaluated against the ground-truth. Rewards are fed back to the Generator if the Encoder is able to identify candidate sets that are generated randomly, while the Selector is rewarded if the best performing rationale candidate set outperforms the one in the adversary. The best performing candidate set from Selector are finally considered as the extracted rationales, and then used in the Encoder to predict entity relation.

dataset with rationales is often too expensive to create; hence our Generator and Selector are mostly trained to identify rationales in a unsupervised manner. The Generator and Selector employ adversarial reinforcement learning where they are rewarded when their outputs are favored in the Encoder over a separate adversary set of Generator and Selector. We illustrate the process in Figure ???. Finally, as many learning systems can be manually optimized, we facilitate a refinement process in our pipeline. The refinement process provides an interface for users to examine and correct rationales to refine the predictions. After refinement, the model is then retrained to incorporate the changes. When refinements are provided, the Generator and Selector are trained in a semi-supervised manner.

Experiments demonstrate that our model produces an 89.5% F1 score on relation classification tasks, compared to a score of 88.0% for the leading state-of-the-art approach (?). While the quality of interpretability is difficult to measure without human evaluation, we make the assumption that rationales that lead to better F1 scores are more interpretable than those that lead to lower scores. In effect, we treat the Encoder as a proxy human reader. Based on this assumption, we evaluated the quality of our rationales by treating attention models as proxy rationale generating models, where rationales are the highest-scoring consecutive n-grams in the attention model. Compared to the proxy rationale models, experiments using rationales generated from our model show a 1.7~5.0% improvement in F1 score. In

addition, in-depth analysis shows that our model is capable of identifying good rationales that a human reader can use to infer relationships, and predictions can be corrected if provided with better rationales. We naively annotated some rationales using the guidelines in ? (?). We discovered that, even with minimal annotations on the dataset, the refined results showed improvement in both F1 score and training speed.

## Related Work

Interpretable neural network models have recently become of increased interest to the natural language processing research community (?). Several approaches to interpretability have been introduced, including: visualization of activation-functions within neural networks (?; ?), accrediting distributed word vectors (?), relevance propagation (?), attention based models (?), and rationale finding (?; ?; ?). Even though each model is optimized for a slightly different task, a common approach they adopted is to discriminate the input data and selectively use fragments from the input. Interpretability can then be measured by the quality of the selected fragments.

The major distinction between the attempts to demystify neural networks is thus the degree of discrimination imposed on the input data. Less-discriminating approaches like attention based models usually allow the model to use the whole input data. These models introduce weights into the models where the weights are commonly generated based on a sim-

ilarity function between signals in the sentence with the desired goal. Signals, which are often n-grams in the sentence, receive higher weights if analogous to the desired goal, and can be considered as “highlights” of the data. While the major benefit of this approach is that can be effective even with naive methods like n-gram to generate signals, weights attributed to the signals do not always correspond to how humans might approach a given problem. Also, complexity of signal generation can become exponential if considering *nonconsecutive* n-grams cases. Most recent research in entity-relation classification falls into the low-constraint category.

High-discriminating approaches, which aim to extract rationales for reasoning, pose a higher degree of constraint of selecting data to be used in the model. Models that fall under this category are usually trained to select only part of the data for reasoning, and to obscure the rest of the data. The criteria for selection can be totally data-oriented or manually curated. The ultimate perk of posing constraints on the model is the proximity to human reasoning which in turn gives better interpretability. Furthermore, it introduces methods to selectively ignore data to be considered, which can in turn helps downsize complexity for skipping n-grams cases.

We position this work on the high-constraint (rationale-finding) side comparing with the attention-based models. In our approach, we find and use nonconsecutive words as rationales in a completely unsupervised way in  $O(NM)$  where  $M \leq N$ .

## Adversarial Rationale Generation

We formalize here the general idea of rationale generation for entity relationship classification in this research. Consider an input sequence of words  $x = \{x_t\}_{t=1}^T$ , where  $e_1$  and  $e_2$  are the entities of interest and  $e_1$  and  $e_2 \subset \{x_1, x_2, \dots, x_T\}$ .  $T$  is the length of the input sequence. The goal of relation classification is to use  $(x, e_1, e_2)$  to predict the relationship  $y$  between  $e_1$  and  $e_2$ . On top of the representation learning regime this research falls into, we refer to the prediction process as encoding the matrix  $(x, e_1, e_2)$  to a target vector that is representative of the desired relationship. We use  $enc(x, e_1, e_2)$  to represent this process. The challenge in this research lies in the complexity within the encoding structure that obscures interpretability. This results in two problems: 1) unjustifiable results, and 2) difficulty in adjusting the model when errors are detected.

The general idea of rationale generation for relation classification is to extract a set of rationales  $r$  that are incorporated into the encoding process as  $enc(r, x, e_1, e_2)$ , where  $r \subset \{x_1, x_2, \dots, x_T\}$  but  $r \notin \{e_1, e_2\}$ . The goal of selecting rationales is to provide a set of terms that help to justify relationships in the sentence, and that make sense to human readers. For example, consider the following sentence:

**Example sentence:** *The disgusting scene was retaliation against her brother Philip who rents the [room] $_{e_1}$  inside this apartment [house] $_{e_2}$  on Lombard street.*

This entails a Component-Whole( $e_1, e_2$ ) relationship, hence good rationales should be words like  $\{inside, apartment\}$  which readers would use to describe the

relationship.

We split the rationale generation process into two steps: the first step is candidate generation  $gen(x, e_1, e_2)$ , which samples  $c = \{c_t\}_{t=1}^T$  from the input and  $c_t \in [0, 1]$  represents that whether  $x_t$  should be considered as a rationale ; the second step is to sample rationales  $r$  from  $c$  and is represented as  $sel(c, e_1, e_2)$ .

We can consider this as an adversarial problem where our target is not to find the best words that the Encoder can reason upon, but instead to find a set of words that is more useful for the Encoder than another set of words. Hence, the goal of the Generator and Selector is to generate rationales that  $enc(r, x, e_1, e_2)$  can outperforms another set of rationales  $enc(\tilde{r}, x, e_1, e_2)$ .

Building on the adversarial approach, the selection of rationales can be completely unsupervised. However, another desirable characteristic is to make the model capable of human intervention. While the rationales are trained unsupervised, we introduce a  $r_{re}$  factor where “re” indicates refinement, and can be used to support semi-supervised training.

A difference between this model and previous work is worth noting: in ? (?), the aim is to generate concise and reasonable summaries for customer reviews. In their research, the goal of  $gen(x)$  is to generate rationales  $r$  that ensure that  $enc(x)$  and  $enc(r)$  are at similar points in the target vector space. In other words, the goal of  $r$  in their research is to serve as a proxy of  $x$  and fool  $enc(\cdot)$ , while in our work the goal of  $r$  is to outperform all possible short enumerations of  $x$  in  $enc(\cdot)$ .

## Encoder, Generator and Selector

We now describe the details of each component in our model. We begin by describing the Encoder in order to highlight the key contributions of our model.

### Encoder

Given an input training tuple  $(x, e_1, e_2, y)$ , where  $x = \{x_t\}_{t=1}^T$ , and  $y$  is the one-hot  $m$  dimensional vector representing the relation class. The goal of the Encoder  $enc(r, x, e_1, e_2)$  is to produce a probability distribution  $\hat{y}$  that approximates to  $y$ , given the set of rationales  $r$  as explained in the previous section. The quality of the Encoder is thus determined by the closeness between  $y$  and  $\hat{y}$ , where closeness is often referred to as loss, and quantified as:

$$loss = \sum_{i=1}^{|N|} -y_i \cdot \log(\hat{y}_i) \quad (1)$$

where  $N$  is the training set and  $|N|$  is the size of the set.

The Encoder performs two computations: 1) it generates a representation of  $(r, x, e_1, e_2)$ , and 2) it generates  $\hat{y}$  by projecting the representation of  $(r, x, e_1, e_2)$  to the  $y$  space. The second part is achieved by a simple feed-forward neural network with a softmax function. The first computation consists of two parts: computing  $V_x$  which represents the encoded sentence, and  $V_r$  which represents the encoded rationals  $r$  and the entities  $(e_1, e_2)$ .

To compute  $V_x$ , we incorporate a Convolutional Neural Network (CNN) model, which aims to represent the overall sentence meaning. The choice of CNN is based on recent successes in NLP applications (?; ?). As in ? (?), the CNN model contains a set of convolutional filters  $W^\ell$  and bias terms  $b^\ell$ . These filters are applied to all possible consecutive n-grams of length  $\ell$  in the sentence so that we can learn a representation for every n-gram in the sentence. Finally, to derive a sentence representation from n-gram representations, we adopt the commonly used approach of taking only the most significant feature for each dimension of the n-gram representations - this is known as *max-pooling*.

We compute  $V_r$  through a feed-forward neural network with a weight matrix  $W_r$  and a bias term  $b_r$ . Note that the input is order sensitive. For example, if  $\mathbf{r} = \{x_{10}, x_{15}\}$ ,  $e_1 = x_4$  and  $e_2 = x_7$ , which means the model uses the 10th and 15th word of the sentence as rationale, the input should be sorted as  $(e_1, e_2, x_{10}, x_{15})$ , whereas the input will be  $(x_3, e_1, e_2, x_{10})$  if  $\mathbf{r} = \{x_3, x_{10}\}$ .

In summary, the Encoder  $enc(\mathbf{r}, \mathbf{x}, e_1, e_2)$  receives as input the selected rationales, the entities, and the target sentence. In mathematical terms:

$$\begin{aligned} \hat{y} &= enc(\mathbf{r}, \mathbf{x}, e_1, e_2) & (2) \\ &= softmax(W_{enc} \cdot [V_r \oplus V_x \oplus e_1 \oplus e_2] + b_{enc}) \\ &\quad (\oplus = \text{vector concatenation}) \\ V_r &= Relu(W_r \cdot [\mathbf{r}, e_1, e_2] + b_r) \\ V_{x_t}^\ell &= Relu(W_x^\ell \cdot x_{t:t+\ell} + b_x^\ell) \\ V_x^\ell &= [V_{x_1}^\ell, V_{x_2}^\ell, \dots, V_{x_{T-\ell}}^\ell] \\ V_x &= \{maxpooling(V_x^\ell)\}^{\ell \in L} \end{aligned}$$

## Generator

The goal of the Generator  $gen(\mathbf{x})$  is to extract pieces of text from input sentence  $\mathbf{x}$  that can be used to connect and predict the relationships between target entities in the sentence. A different way to view the role of the generator is to derive another set of binary variables  $\mathbf{c} = \{c_t\}_{t=1}^T$  where each  $c_t \in [0, 1]$  is whether  $x_t$  is chosen as a candidate rationale. We represent it as  $\mathbf{c} \sim gen(\mathbf{x})$ .

We adopted a recurrent neural network (RNN) model in this work to compute  $\mathbf{c}$  due to the sequential characteristics of  $\mathbf{x}$  and  $\mathbf{c}$ . RNNs are ideal for sequential modeling, where the output of step  $t$  is used as a part of the input for step  $t+1$ . Specifically, all inputs of the model share the same set of transformations and output functions, where these functions help to project inputs to the target dimensional space, and are fed back for use in the next step. In this work, we use a bi-directional RNN where the input is modeled from both ends, and shares the same parameters in both directions. In the formal description in Equation (3), both  $W$  and  $U$  are weight matrices, and  $b$  are bias terms.

*RecurrentUnit* : (3)

$$\begin{aligned} z_t &= \sigma(W_z \cdot x_t + U_z \cdot h_{t-1} + b_z) \\ o_t &= \sigma(W_r \cdot x_t + U_r \cdot h_{t-1} + b_r) \\ \hat{h}_t &= \tanh(W_h \cdot x_t + U_h \cdot (h_{t-1} \circ o_t) + b_h) \\ h_t &= z_t \circ \hat{h}_t + (1 - z_t) \circ h_{t-1} \\ &\quad (\circ = \text{elementwise-multiplication}) \\ \overleftarrow{h}_t &= RecurrentUnit(x_t, \overleftarrow{h}_{t+1}) \\ \overrightarrow{h}_t &= RecurrentUnit(x_t, \overrightarrow{h}_{t-1}) \\ s_t &= sigmoid(W_s \cdot [\overleftarrow{h}_t, \overrightarrow{h}_t, e_1, e_2] + b_s) \end{aligned}$$

In the most simplistic generator, the probability that  $x_t$  is chosen is conditionally independent of all other  $\mathbf{x}$ . In other words, we can sample candidate  $\mathbf{c}$  from a uniform distribution using  $\{s_t\}_{t=1}^T$ . However, we observed from the data set that our target rationales often consist only very limited words - for example “*room inside apartment*”. We thus added the following Equation (4) for  $\mathbf{c}$  to limit the number of rationales selected to be at most  $J$ .

$$\begin{aligned} c_t &= \begin{cases} 1, & \text{sort}(\{(s_t > rand(0, 1)) \cdot s_t\}_{t=1}^T)[1 : J] \\ 0, & \text{otherwise} \end{cases} & (4) \\ \mathbf{c} &= \{c_t\}_{t=1}^T \end{aligned}$$

## Selector

The last component of the model is the Selector  $sel(\mathbf{c})$ , where the goal of the Selector is to choose rationales  $\mathbf{r}$  from  $\mathbf{c}$ . The Selector is a feed-forward neural network can be considered as  $\mathbf{r} \sim sel(\mathbf{c})$ , and defined as follows:

$$\begin{aligned} \hat{c}_t^j &= \begin{cases} 1, & \text{sort}(\{c_t \cdot s_t\}_{t=1}^T)[1 : j], j \in [1 : J] \\ 0, & \text{otherwise} \end{cases} & (5) \\ \hat{\mathbf{c}}^j &= \{\hat{c}_t^j\}_{t=1}^T, \text{ where } \hat{c}_t^j = 1 \\ score(\hat{\mathbf{c}}^j) &= W_c \cdot enc(\hat{\mathbf{c}}^j, \mathbf{x}, e_1, e_2) \\ scores^j &= softmax(\{score(\hat{\mathbf{c}}^j)\}_{j=1}^J)^j \\ \hat{scores}^j &= \begin{cases} 1, & j = argmax(scores) \\ 0, & \text{otherwise} \end{cases} \\ \mathbf{r} &= \{x_t\}, \text{ where } \hat{c}_t^j = 1, j = argmax(scores) \end{aligned}$$

Notice that, in order to generate scores for each candidate, the model adopts the same weight functions in the Encoder section. A more intuitive way to understand this is that the Selector scores each candidate by their expected projection and contribution to the  $y$ -space. For instance, if using two rationales ( $j = 2$ ) outperforms using only one rationale ( $j = 1$ ), the Selector will choose to use two rationales.



## Joint Objective and Adversarial Training

In order to jointly train the three components, we formulate a joint loss function that binds the components together. However, the major challenge lies in the fact that both the Generator and Selector are not provided with ground-truth labels, and hence are unsupervised. In order to facilitate this model, we formulate an adversarial training regime.

Our adversarial training model is similar to that used in the GAN model (?). In previous work, the example application was generating images based on sentences, where the Generator generates a fake image that tries to fool the Encoder. After training, the Generator becomes effective at generating fooling images, and the Encoder become discerning at identifying these fakes. We applied this same concept in our model but with one key modification: all the components of our model work together and compete with an adversary Generator and Selector.

In our model, the Generator competes with an adversary Generator, for which we use a randomized approach. The adversary Generator will sample  $\tilde{c}$  randomly from  $\mathbf{x}$ . The adversary rationales  $\tilde{r}$  will then be selected using  $\tilde{c}$  as per Equation (5).  $\tilde{r}$  will be passed to the Encoder to generate a projection onto the target dimension, which we consider as  $\tilde{y}$ . We compare  $\tilde{y}$  with  $\hat{y}$  in terms of their similarity to  $y$  and is denoted by  $\mathcal{D} \in \{0, 1\}$ .  $\mathcal{D} = 1$  when  $\hat{y}$  is closer to  $y$  comparing to  $\tilde{y}$  and is noted as ‘model success case’, meanwhile  $1 - \mathcal{D} = 1$  for the opposite situation and is noted as ‘adversary success case’. We further split  $\hat{y}$ ,  $\tilde{y}$  and  $\mathcal{D}$  into  $\hat{y}^j$ ,  $\tilde{y}^j$ ,  $\mathcal{D}^j$  and  $j \in [1 : J]$ , where  $\mathcal{D}^j$  represents that  $\hat{y}^j$  outperforms  $\tilde{y}^j$  when using  $j$  rationales. Finally, we introduce a penalizing factor  $P^j$  when  $j \neq \text{argmax}(\text{scores})$  to penalize gradients from less-performing cases. The objective for each component reacts to the two different cases differently during training, which we summarize as follows:

$$\begin{aligned}
 f(c_t^j) &= -\log(s_t) \cdot c_t^j - \log(1 - s_t) \cdot (1 - c_t^j) \\
 f_y(\hat{y}_i^j) &= -y_i \cdot \log(\hat{y}_i^j) \\
 \mathcal{L}_{generator} &= \sum_{j=1}^J \sum_{t=1}^T \mathcal{P}^j \left[ \mathcal{D}^j f(\hat{c}_t^j) + (1 - \mathcal{D}^j) f(\tilde{c}_t^j) \right] \\
 \mathcal{L}_{selector} &= \sum_{j=1}^J -\mathcal{D}^j * \text{scores}^j * \log(\text{scores}^j) \\
 \mathcal{L}_{encoder} &= \sum_{j=1}^J P^j \left[ \mathcal{D}^j f_y(\hat{y}^j) - (1 - \mathcal{D}^j) f_y(\tilde{y}^j) \right]
 \end{aligned} \tag{6}$$

The final joint objective and expected cost is defined as:

$$\mathcal{L}(\mathbf{r}, \mathbf{x}, e_1, e_2, y) = \mathcal{L}_{generator} + \mathcal{L}_{selector} + \mathcal{L}_{encoder} \tag{7}$$

$$\min_{\theta_e, \theta_g, \theta_s} \sum_{(\mathbf{x}, e_1, e_2, y) \in N} \mathbb{E}_{\mathbf{r} \sim \text{sel}(e), c \sim \text{gen}(\mathbf{x})} \mathcal{L}(\mathbf{r}, \mathbf{x}, e_1, e_2, y) \tag{8}$$

where  $N$  is the training set and  $\theta_e$ ,  $\theta_g$ ,  $\theta_s$  are parameters used in the Encoder, Generator and Selector respectively.

Notice that the Generator and Encoder handle losses from adversary success cases differently. The Generator exploits the adversary success cases as a potential learning resource, and uses  $\tilde{r}$  to train the RNN model. However, the Encoder is trying to make good predictions based only on rationales from the Generator, so it will treat good rationales from the adversary Generator as false cases and try to maximize adversary losses.

Finally, if ground-truth labels or human interventions for rationales  $r_{re}$  are provided, then we change  $r$  in Equation (2) and  $\hat{c}$  in Equation (6) to  $r_{re}$ . This can be understood as a simple sequential discriminatory model where the Generator and the Selector is guided by  $r_{re}$ . However, while keeping the reminder of the model the same, introducing  $r_{re}$  does not guarantee the Selector and Encoder will use  $r_{re}$  while reasoning.

## Experiments

### Dataset

We performed our evaluation experiments on SemEval 2010 Task 8 data. A collection of sentences is provided, where each sentence is labeled with the target entities and the relationship between the entities.

The dataset contains 10,717 sentences, and 9 types of relationship plus an ‘‘Other’’ class. Relationship types (except for ‘‘Other’’) are expressed bi-directionally, making 19 classes in total. Following SemEval 2010’s protocol, we used the macro-F1 metric in the experiment, excluding all cases of the ‘‘Other’’ class.

### Experimental Setup

Here we describe the detailed parameters used in the experiment. We initialized word vectors using the GoogleNews<sup>1</sup> vector that is learned by word2vec algorithm. Word vector dimensionality is set to 300. We naively set the the recurrent unit size in the Generator to 50, and the CNN filter size in the Encoder 50 for each  $\ell$ . We use at most 3 rationales and set the penalizing factor  $P^j$  to 0.5 based on our analysis on the training set.

We used the Adagrad optimizer in our experiments - the learning rate is initially set to 0.01, and reduces by 10% after every iteration. During training, we sampled 50 times from the Generator, and re-sampled 50 times when training the Selector and Encoder.

### Results

We present the comparison between our model and other models in Table ???. Dependency tree models are neural network models that utilize a grammar parser to predict relation classes. Different models utilize the dependency tree differently. For instance in MVRNN (?) they used the whole dependency tree, while in SPTree (?) they experimented using both the whole tree and only nodes on the path that connected the target entities. We refer to other models that work without using a dependency tree as independent models.

<sup>1</sup><https://code.google.com/p/word2vec>

Classifier	F1
<i>Manually Engineered Models</i>	
SVM(?)	82.2
<i>Dependency Tree Models</i>	
RNN (?)	77.6
MVRNN (?)	82.4
FCM (?)	83.0
Hybrid FCM (?)	83.4
DepNN (?)	83.6
DRNNs (?)	85.6
SPTree (?)	84.5
<i>Independent Models</i>	
CNN & softmax (?)	82.7
Stackforward(?)	83.4
Vote-bidirect (?)	84.1
Vote-backward (?)	84.1
ATT-Input-CNN (?)	87.5
ATT-Pooling-CNN (?)	88.0
Our model	<b>89.5</b>

Table 1: Comparisons with benchmarking models

Besides the general goal of improving F1 score, we also evaluated the generated rationales qualities. Due to the size of the dataset, we were not able to annotate all samples manually - in order to handle this, we make an assumption that an Encoder will perform better if provided with rationales that are more interpretable, that is, rationales which are semantically closer to the entities and the desired goal. In effect, we utilize a set of identically-initialized Encoders, and evaluate rationales from different rationale finding approaches through a naive way by simply comparing the corresponding F1 scores from each Encoder.

However, since there are few comparable rationale models in this field, we created a proxy rationale model for comparison. We adapted a CNN model similar to ? (?), and added attention weights calculated by the similarity with ( $e_1, e_2$ ). To extract rationales from the proxy model, we first trained the CNN model, then extracted the top-weighted text fragments in the model as rationales. The proxy rationales were finally sent to the Encoder to compare with our model. We experimented this idea with 1, 2 and 3 rationales extracted from the proxy model, and results are shown in Table ??.

Classifier	F1
Proxy attention model	87.5
<i>Proxy rationales + Encoder</i>	
One rationale from proxy	84.5
Two rationales from proxy	85.3
Three rationales from proxy	87.8
Our model	<b>89.5</b>

Table 2: Comparison with proxy rationale models.

Finally, we summarized a few samples of rationales chosen from the test set in Table ??. From the table, we see that the selected rationales are similar to those that a human

reader might use to infer the relationship between entities.

## Results with Ground Truth for Rationales

Finally we evaluate our model when human annotators are included to provide ground-truth for rationales. While we were not able to perform detailed labeling, we naively generated ground-truth for rationales following the guidelines provided by SemEval 2010 task 8. We list the words that were used as ground-truth for rationales in Table ??. The Generator and Selector were trained in a supervised fashion if a ground-truth rationale was contained in a training sentence. Training samples that did not contain ground-truth rationales were trained in an unsupervised manner.

We injected ground-truth at three different training stages of the model: 1) converged, 2) converging, 3) from the beginning. Providing ground-truth to converged and converging models is intended to simulate models at production stage, and where human annotators are included only for final adjustment. Intervention performed at the beginning is designed to see how the model performs when it is closer to a fully supervised model.

We show the results in Figure 2. We can see that, when providing rationale labels at the beginning, the model learns much faster. Note that there is a performance drop in the converged and converging models immediately after ground-truth is provided at iteration 15 and 25 respectively. This result is expected since the input to the Encoder changes after introducing the labeled rationales, and this component requires a few iterations to assimilate the changes. Eventually, models with ground-truth labels achieve slightly better performance, typically a 1-1.5% improvement on F1 score.

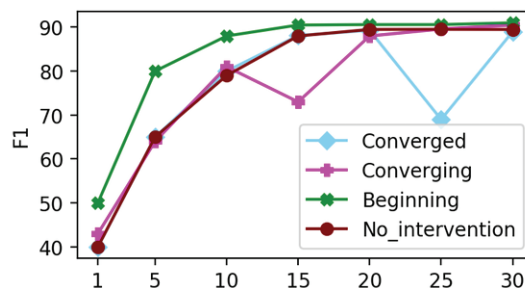


Figure 2: F1 results using different intervention strategy.

Finally, we present a few example cases that were originally predicted incorrectly but that were subsequently corrected after manual labeling of rationales. Rationales are denoted in bold below:

**Example sentence:** Instrument-Agency( $e_2, e_1$ ):

**Before intervention** - The [generator] $e_1$  **creates** electricity using much the same [principle] $e_2$  as the alternator on your car ( depending on the turbine type )” (predicted as Cause-Effect( $e_2, e_1$ ))

**After intervention** - The [generator] $e_1$  **creates** electricity **using** much the same [principle] $e_2$  as the alternator on your car ( depending on the turbine type )”

Relation Class	Selected rationales with target entities
Cause-Effect	congestion <sub>e1</sub> caused delays <sub>e2</sub> , bombing <sub>e1</sub> ... resulted in ... deaths <sub>e2</sub> devastation <sub>e1</sub> caused by the ... tsunami <sub>e2</sub> , anxiety <sub>e1</sub> caused by the accident <sub>e2</sub>
Component-Whole	carriages <sub>e1</sub> of ... train <sub>e2</sub> , title <sub>e1</sub> of the book <sub>e2</sub> , the kitchen <sub>e1</sub> has a fridge <sub>e2</sub> horn <sub>e1</sub> has ... end <sub>e2</sub> , this ... machine <sub>e1</sub> has ... reels <sub>e2</sub> , the shank <sub>e1</sub> of the hook <sub>e2</sub>
Content-Container	object <sub>e1</sub> was ... bag <sub>e2</sub> , wallet <sub>e1</sub> was inside ... locker <sub>e2</sub> , a mouse <sub>e1</sub> ... found in ... loaf <sub>e2</sub> bottle <sub>e1</sub> ... of ... dust <sub>e2</sub> , basket <sub>e1</sub> was full ... faces <sub>e2</sub> , the backpack <sub>e2</sub> contained a ... computer <sub>e2</sub>
Entity-Destination	yeast <sub>e1</sub> into ... tubes <sub>e2</sub> , niece <sub>e1</sub> moved into... apartment <sub>e2</sub> ... he put ... forkful <sub>e1</sub> ... mouth <sub>e2</sub> , hero <sub>e1</sub> enters into the building <sub>e2</sub>
Entity-Origin	artwork <sub>e1</sub> ... from ... evening <sub>e2</sub> , crisis <sub>e1</sub> originated in ... sector <sub>e2</sub> takara ... is a ... plum <sub>e1</sub> wine <sub>e2</sub> , the soul <sub>e1</sub> departs from the body <sub>e2</sub>
Instrument-Agency	marking birds <sub>e1</sub> ... researchers <sub>e2</sub> , this technique <sub>e1</sub> ... used ... websites <sub>e2</sub> technician <sub>e1</sub> uses ... emulator <sub>e2</sub> , shaman <sub>e1</sub> cured ... with herbs <sub>e2</sub> , man <sub>e1</sub> attacked ... with a ... bat <sub>e2</sub>
Member-Collection	chief <sub>e1</sub> of ... police <sub>e2</sub> , dean <sub>e1</sub> of the faculty <sub>e2</sub> , father mother <sub>e1</sub> was ... member ... team <sub>e2</sub> clowder <sub>e1</sub> of cats <sub>e2</sub> , an anthology <sub>e1</sub> of ... poems <sub>e2</sub> , the core ... body <sub>e1</sub> of ... traders <sub>e2</sub>
Message-Topic	book <sub>e1</sub> provides ... role <sub>e2</sub> , this figure <sub>e1</sub> illustrates ... results <sub>e2</sub> , the talk <sub>e1</sub> was about ... operations <sub>e2</sub> details <sub>e1</sub> about ... interview <sub>e2</sub> , the ... episode <sub>e1</sub> ... documented on video <sub>e2</sub>
Product-Producer	message <sub>e1</sub> from ... friend <sub>e2</sub> , paper <sub>e1</sub> co-authored by ... staff <sub>e2</sub> , blisters <sub>e1</sub> are caused by antibodies <sub>e2</sub> cooper <sub>e1</sub> makes leak ... barrels <sub>e2</sub> , the researchers <sub>e1</sub> produced a report <sub>e2</sub>

Table 3: List of words/phrases that are selected as rationales by our model. Additional words are denoted as ... if they are located between entities and rationales but are not selected as rationales.

Relation Class	Selected Words
Cause-Effect	from, cause, caused, causes
Component-Whole	of
Content-Container	in
Entity-Destination	to, into
Entity-Origin	from
Instrument-Agency	use, uses, used, using
Member-Collection	of, is
Message-Topic	of, about
Product-Producer	make, makes, made

Table 4: List of ground-truth rationales used

**Example sentence:** Instrument-Agency( $e_2, e_1$ ):  
**Before intervention** - *The [manufacturer]<sub>e1</sub> assembles the order using [parts]<sub>e2</sub> supplied by his preferred supplier , and ships the order to the retailer.*  
(predicted as Product-Producer( $e_2, e_1$ ))  
**After intervention** - *"The [manufacturer]<sub>e1</sub> assembles the order **using** [parts]<sub>e2</sub> supplied by his preferred supplier , and ships the order to the retailer."*

## Discussion

### Rationale Labeling

In the previous section we chose to label the rationales using a naive approach. The major drawback was that human language is complicated and often there exists several ways to express a relation class. While our naive approach only considers only very basic cases, future research would benefit from detailed annotated rationales for better generalizability.

### End-to-End Training and Combining Components

At early stages of this research, we experimented with combining the Selector and Generator, and also tried end-to-end training. However we noticed that these often led to worse performance. We reason that this is because, in both setups, the selected rationales are heavily biased by the first few training epochs. However, this might not apply to cases where ground-truth for rationales is provided.

### No-Rationale Cases

We notice that there are few cases where the model might provide ambiguous insight due to no good rationales being present, such as in the following example. The selected rationale, which is not useful, is shown in bold:

**Component-Whole( $e_2, e_1$ )** - *At the bottom of the [church]<sub>e1</sub> [steps]<sub>e2</sub> **were** three brown parishioners; two more were perched precariously on the railing of the deck.*

We were focussed only on selecting rationales in this work; however it may be possible to handle the no-rationale case in future work by introducing a mechanism that turns the Generator and Selector off under certain circumstances.

### Potential Uses

Besides relation classification, another potential use-case for our model is in generating rules for relation classification. It should be possible to automatically discover sets of rules for relation classification by supplying the Generator and Selector with carefully chosen test cases. This may be useful for knowledge base construction since building extraction rules manually becomes infeasible as the number of possible relation classes increases.

## Conclusion

In this paper, we proposed a novel generative adversarial approach for entity relation class prediction. In our model, we first generate a set of candidate rationales inside a generative RNN, and then predict the relationship between entities using a discriminative model that exploits the rationales. The generative process in our model can be turned into a semi-supervised model that is capable of incorporating human annotations. This not only boosts prediction accuracy but also provides a way for users to easily interpret and adjust the model. Quantitative analysis demonstrated that our model provides better performance than state-of-the-art dependency-tree and independent models. In addition, we showed that our model is able to select rationales that make sense to humans when combined with target entities.

## Acknowledgments

We thank the AAAI reviewers for their detailed and insightful feedback. This material is based upon work supported in whole or in part with funding from LAS. Any opinions, findings, conclusions, or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the LAS and/or any agency or entity of the United States Government.

## References

- Bahdanau, D.; Cho, K.; and Bengio, Y. 2015. Neural machine translation by jointly learning to align and translate. In *ICLR*.
- Faruqui, M.; Dodge, J.; Jauhar, S. K.; Dyer, C.; Hovy, E.; and Smith, N. A. 2015. Retrofitting word vectors to semantic lexicons. In *NAACL HLT*, 1606–1615.
- Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2014. Generative adversarial nets. In *NIPS*, 2672–2680.
- Gormley, M. R.; Yu, M.; and Dredze, M. 2015. Improved relation extraction with feature-rich compositional embedding models. In *EMNLP*, 1774–1784.
- Hendrickx, I.; Kim, S. N.; Kozareva, Z.; Nakov, P.; Ó Séaghdha, D.; Padó, S.; Pennacchiotti, M.; Romano, L.; and Szpakowicz, S. 2009. Semeval-2010 task 8: Multi-way classification of semantic relations between pairs of nominals. In *NAACL HLT*, 94–99.
- Kalchbrenner, N.; Grefenstette, E.; and Blunsom, P. 2014. A convolutional neural network for modelling sentences. In *ACL*, 655–665.
- Kim, Y. 2014. Convolutional neural networks for sentence classification. In *EMNLP*, 1746–1751.
- Lei, T.; Barzilay, R.; and Jaakkola, T. 2016. Rationalizing neural predictions. In *EMNLP*, 107–117.
- Liu, Y.; Wei, F.; Li, S.; Ji, H.; Zhou, M.; and Wang, H. 2015. A dependency-based neural network for relation classification. In *ACL*, 285–290.
- Miwa, M., and Bansal, M. 2016. End-to-end relation extraction using lstms on sequences and tree structures. In *ACL*, 1105–1116.
- Montavon, G.; Samek, W.; and Müller, K.-R. 2017. Methods for interpreting and understanding deep neural networks. In *CoRR*. [abs/1706.07979](https://arxiv.org/abs/1706.07979).
- Nguyen, T. H., and Grishman, R. 2015. Combining neural networks and log-linear models to improve relation extraction. In *IJCAI Workshop on Deep Learning for Artificial Intelligence (DLAI)*.
- Nguyen, A.; Dosovitskiy, A.; Yosinski, J.; Brox, T.; and Clune, J. 2016. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. In *NIPS*, 3387–3395.
- Niu, F.; Zhang, C.; Ré, C.; and Shavlik, J. W. 2012. DeepDive: Web-scale knowledge-base construction using statistical learning and inference. *VLDS* 12:25–28.
- Ribeiro, M. T.; Singh, S.; and Guestrin, C. 2016. Why should i trust you?: Explaining the predictions of any classifier. In *Knowledge Discovery and Data Mining (KDD)*, 1135–1144.
- Rink, B., and Harabagiu, S. 2010. Utd: Classifying semantic relations by combining lexical and semantic resources. In *SemEval*, 256–259.
- Samek, W.; Binder, A.; Montavon, G.; Lapuschkin, S.; and Müller, K.-R. 2016. Evaluating the visualization of what a deep neural network has learned. *IEEE transactions on neural networks and learning systems* 28:2660–2673.
- Simonyan, K.; Vedaldi, A.; and Zisserman, A. 2013. Deep inside convolutional networks: Visualising image classification models and saliency maps. *CoRR*. [abs/1312.6034](https://arxiv.org/abs/1312.6034).
- Socher, R.; Huval, B.; Manning, C. D.; and Ng, A. Y. 2012. Semantic compositionality through recursive matrix-vector spaces. In *EMNLP*, 1201–1211.
- Wang, L.; Cao, Z.; de Melo, G.; and Liu, Z. 2016. Relation classification via multi-level attention cnns. In *ACL*, 1298–1307.
- Xu, K.; Feng, Y.; Huang, S.; and Zhao, D. 2015. Semantic relation classification via convolutional neural networks with simple negative sampling. In *EMNLP*, 536–540.
- Yu, M.; Gormley, M.; and Dredze, M. 2014. Factor-based compositional embedding models. In *NIPS*, 95–101.
- Zaidan, O.; Eisner, J.; and Piatko, C. D. 2007. Using annotator rationales to improve machine learning for text categorization. In *NAACL HLT*, 260–267.
- Zeng, D.; Liu, K.; Lai, S.; Zhou, G.; and Zhao, J. 2014. Relation classification via convolutional deep neural network. In *COLING*, 2335–2344.
- Zhang, Y.; Marshall, I.; and Wallace, B. C. 2016. Rationale-augmented convolutional neural networks for text classification. In *EMNLP*, 795–804.