



ICML

International Conference
On Machine Learning

Neura/NERE: Neural Named Entity Relationship Extraction for End-to-End Climate Change Knowledge Graph Construction

Prakamya Mishra¹

Rohan Mittal¹

¹Independent Researcher

Outline

1. Motivation
2. *SciDCC* Dataset
3. *Neural*NERE
4. Conclusion & Future Work

Motivation

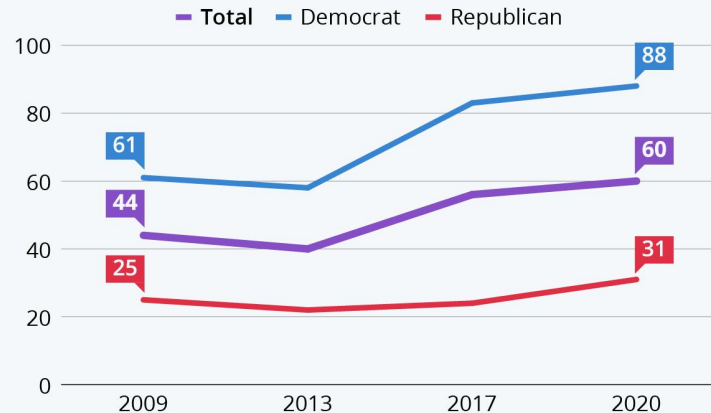
Why?

- News outlets have increased awareness about climate change through their news articles.
- Volume of news articles have been increasing rapidly.
- Difficult to extract useful information.

Algorithms that can extract and organize climate change information by condensing the relevant knowledge directly from a large collection of noisy and redundant news articles could prove to be highly valuable.

Earth Day: Climate Change Awareness Grows

Percentage of U.S. adults who say climate change is a major threat



Source: Pew Research Center



statista

Source: <https://www.statista.com/chart/21415/climate-change-awareness-earth-day/>

Motivation

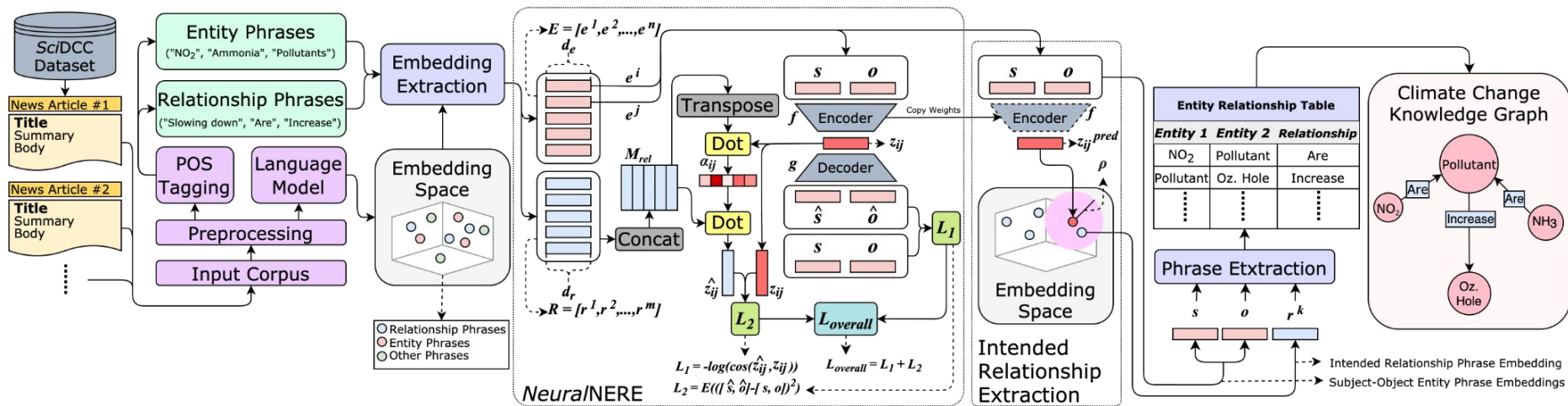
Knowledge Graph Construction

- Supervised or Semi-Supervised [Bollacker et al., 2008; Carlson et al., 2010; Dong et al., 2014]
- Neural models [Bosselut et al., 2019; Balazevic et al., 2019; Xiong et al., 2018; Trivedi et al., 2017; Garcí'a-Duran et al., 2018)]

The problem with the existing approaches is that they all use some prior knowledge in the form of Knowledge Bases (KB) to learn to predict relationships between the subject-object entity phrases for constructing KGs.

Our Solution

- We propose *Neura/NERE*, an end-to-end Neural Named Entity Relationship Extraction model for constructing climate change knowledge graphs directly from the raw text of relevant news articles.
- Additionally, we introduce a new climate change news dataset (called SciDCC dataset) containing over 11k news articles scraped from the Science Daily website.



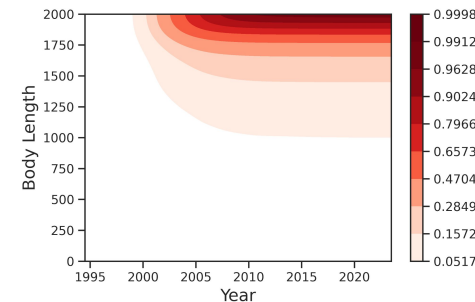
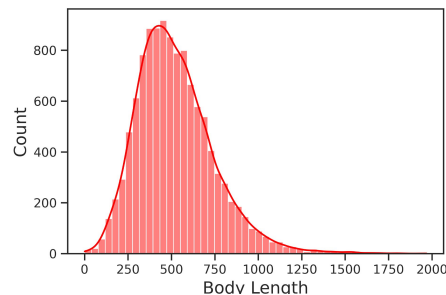
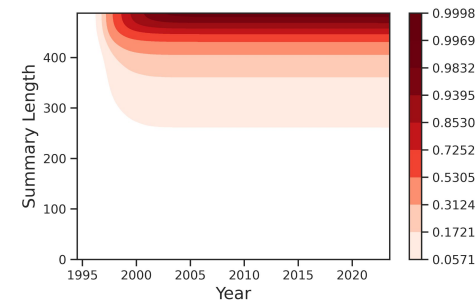
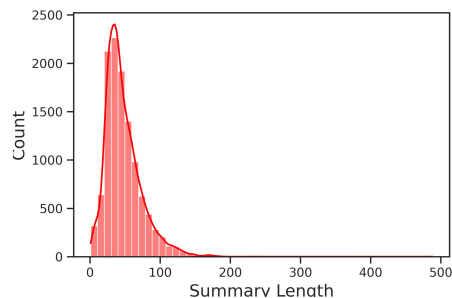
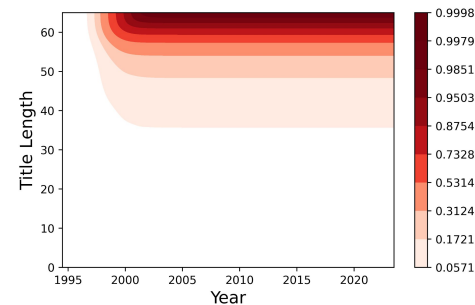
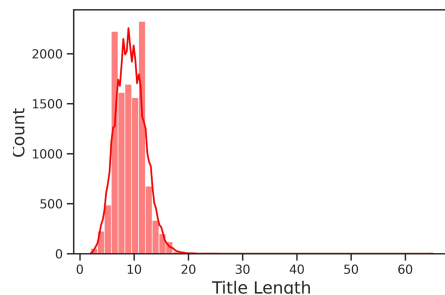
SciDCC Dataset

Table 1. Key statistics of the SciDCC dataset.

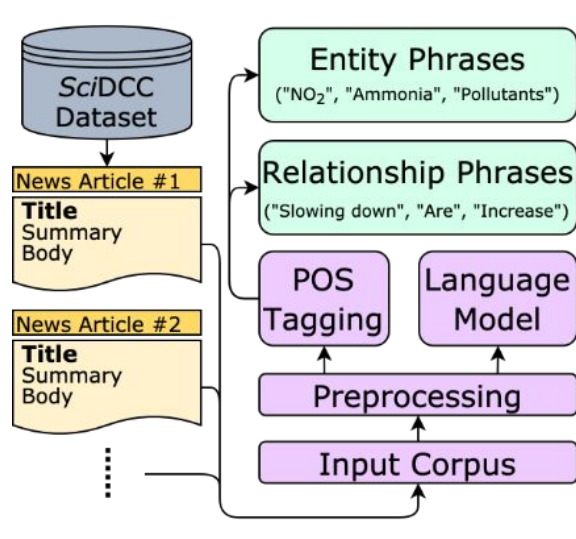
# NEWS ARTICLES	11,539	# NEWS CATEGORY	20
AVG. TITLE LEN.	9.32	MAX. TITLE LEN.	65
AVG. SUMMARY LEN.	47.28	MAX. SUMMARY LEN.	488
AVG. BODY LEN.	523.18	MAX. BODY LEN.	1968

Table 2. News Category Statistics

NO.	CATEGORY	# NEWS ARTICLES
1	EARTHQUAKES	986
2	POLLUTION	945
3	GENETICALLY MODIFIED	914
4	HURRICANES CYCLONES	844
5	AGRICULTURE & FOOD	844
6	ANIMALS	758
7	WEATHER	719
8	ENDANGERED ANIMALS	701
9	CLIMATE	700
10	OZONE HOLES	623
11	BIOLOGY	620
12	NEW SPECIES	527
13	ENVIRONMENT	478
14	BIOTECHNOLOGY	460
15	GEOGRAPHY	407
16	MICROBES	398
17	EXTINCTION	356
18	ZOOLOGY	210
19	GEOLOGY	28
20	GLOBAL WARMING	21



NeuralNERE



- We first create an input corpus by extracting the raw text from the summary part and body part of all the articles present in the SciDCC dataset.
- This input corpus is first preprocessed (tokenization, lower-casing, stemming, lemmatization) and then used for: (1) Fine-tuning a language model for learning the word embedding representations corresponding to every word present in the corpus; (2) Extracting all the named entity phrases as well as all the possible relationship phrases using noun & verb phrase chunking.

NeuralNERE: Intended Relationship Representation

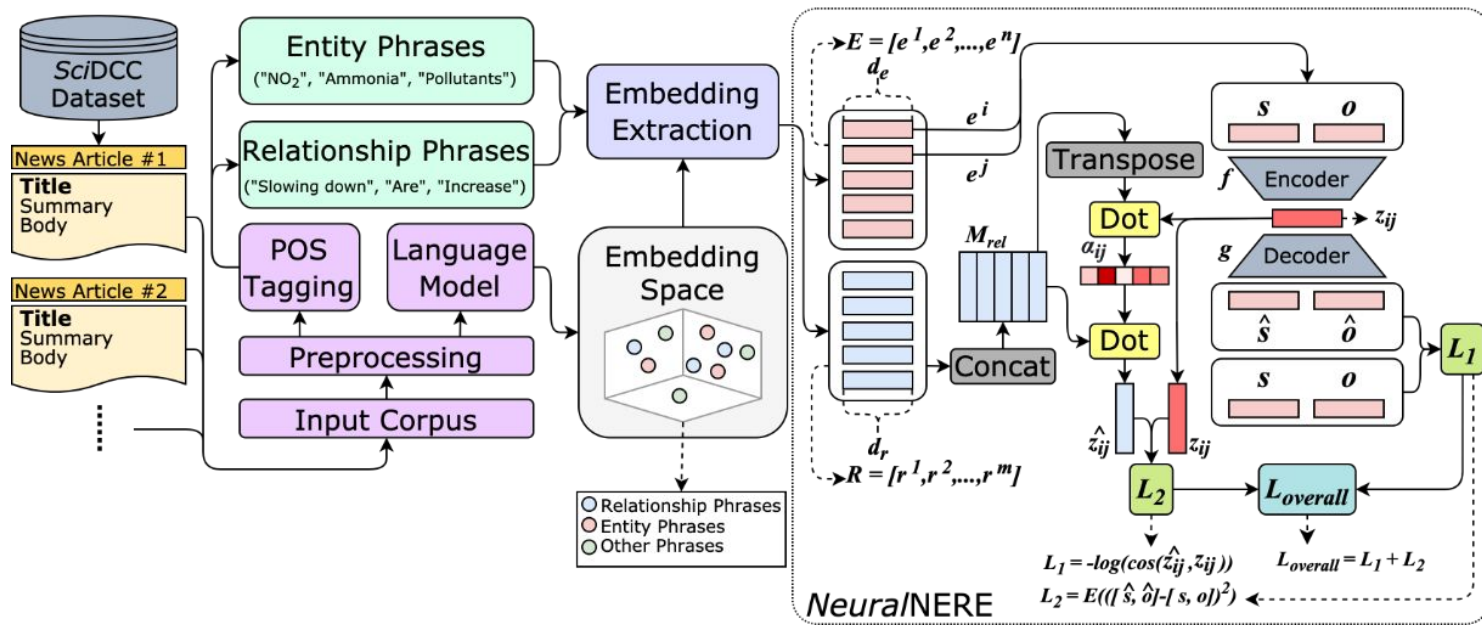
For training, NeuralNERE uses:

- A pair of entity phrase representations represented by (s, o) , where the $s \in \mathbb{R}^{d_e}$ is the entity phrase representation of the i^{th} entity phrase (in E) which acts as the subject in the subject-object relationship, and $o \in \mathbb{R}^{d_e}$ is the entity phrase representation of the j^{th} entity phrase (in E) which acts as the object in the subject-object relationship.
- A relationship phrase matrix $M_{rel} \in \mathbb{R}^{d_r \times m}$, which is basically a matrix constructed by concatenating (\otimes) all the m relationship phrase representations together from the relationship phrase list.

$$E = [e^1, \dots, e^n]; R = [r^1, \dots, r^m] \quad d_e, d_r \in \mathbb{N} \quad (1)$$

$$M_{rel} = r^1 \otimes r^2 \otimes \dots \otimes r^m; M_{rel} \in \mathbb{R}^{d_r \times m}, r^i \in \mathbb{R}^{d_r} \quad (2)$$

Neura/NERE: Intended Relationship Representation



$$E = [e^1, \dots, e^n]; R = [r^1, \dots, r^m] \quad d_e, d_r \in \mathbb{N} \quad (1)$$

$$M_{rel} = r^1 \otimes r^2 \otimes \dots \otimes r^m; M_{rel} \in \mathbb{R}^{d_r \times m}, r^i \in \mathbb{R}^{d_r} \quad (2)$$

NeuralNERE: Intended Relationship Representation

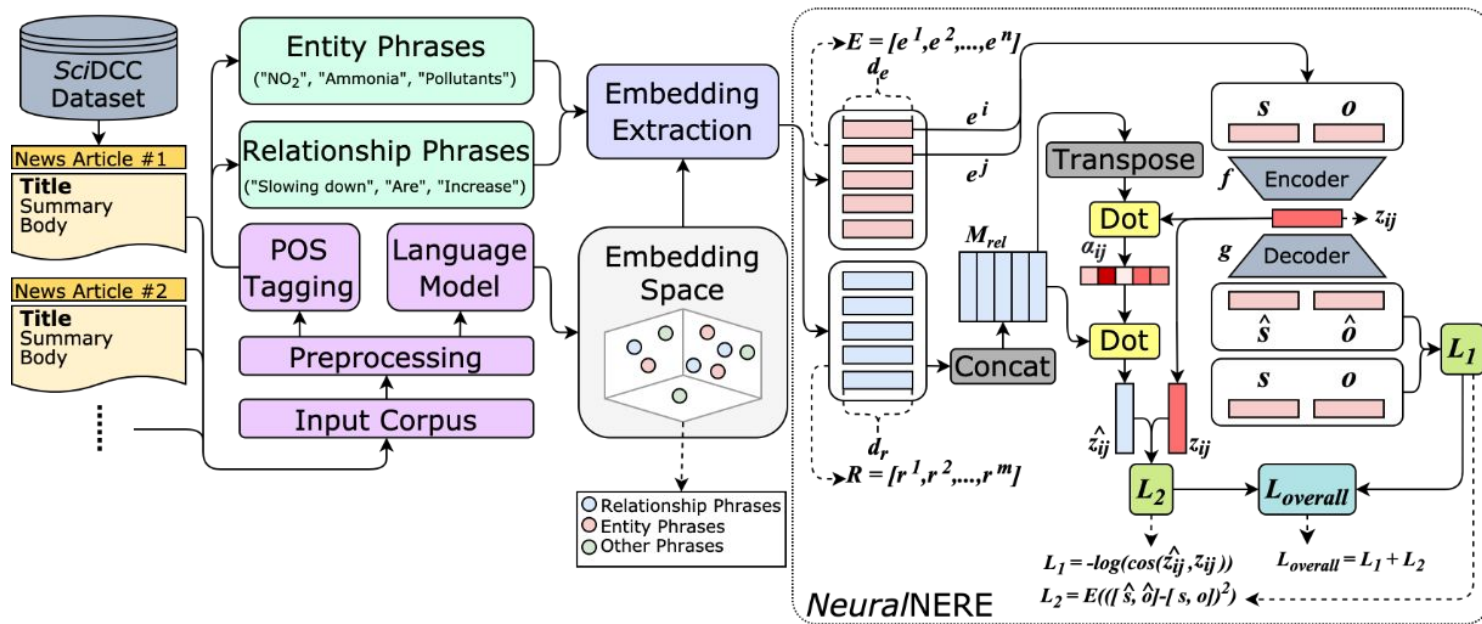
- Next, NeuralNERE uses an encoder-decoder network for encoding the relationship between the subject entity phrase and object entity phrase into an encoded representation $z_{ij} \in \mathbb{R}^{d_r}$, having the same embedding representation size as that of the relationship phrases.
- This encoded vector z_{ij} represents the embedding representation of the intended relationship phrase between subject-object entity phrase pairs that NeuralNERE is trying to learn.
- Although we want the encoded vector z_{ij} to capture the relationship between phrases represented by s and o , but in reality, we don't really know much about the nature of information being captured in the encoded vector. In order to force z_{ij} to capture such relationship-based information, NeuralNERE uses the relationship phrase matrix M_{rel} which contains embedding representations of all the existing relationship phrase.

$$\hat{s}, \hat{o} = g(z_{ij}); \quad z_{ij} = f(s, o) \quad \hat{s}, \hat{o} \in \mathbb{R}^{d_e} \quad (3)$$

$$\alpha_{ij} = [\alpha^1, \dots, \alpha^m]^T = D_{rel} \bullet M_{rel}^T \bullet z_{ij} \quad (4)$$

$$\hat{z}_{ij} = M_{rel} \bullet \alpha_{ij} \quad (5)$$

Neura/NERE: Intended Relationship Representation



$$\hat{s}, \hat{o} = g(z_{ij}); \quad z_{ij} = f(s, o) \quad \hat{s}, \hat{o} \in \mathbb{R}^{d_e} \quad (3)$$

$$\alpha_{ij} = [\alpha^1, \dots, \alpha^m]^T = D_{rel} \bullet M_{rel}^T \bullet z_{ij} \quad (4)$$

$$\hat{z}_{ij} = M_{rel} \bullet \alpha_{ij} \quad (5)$$

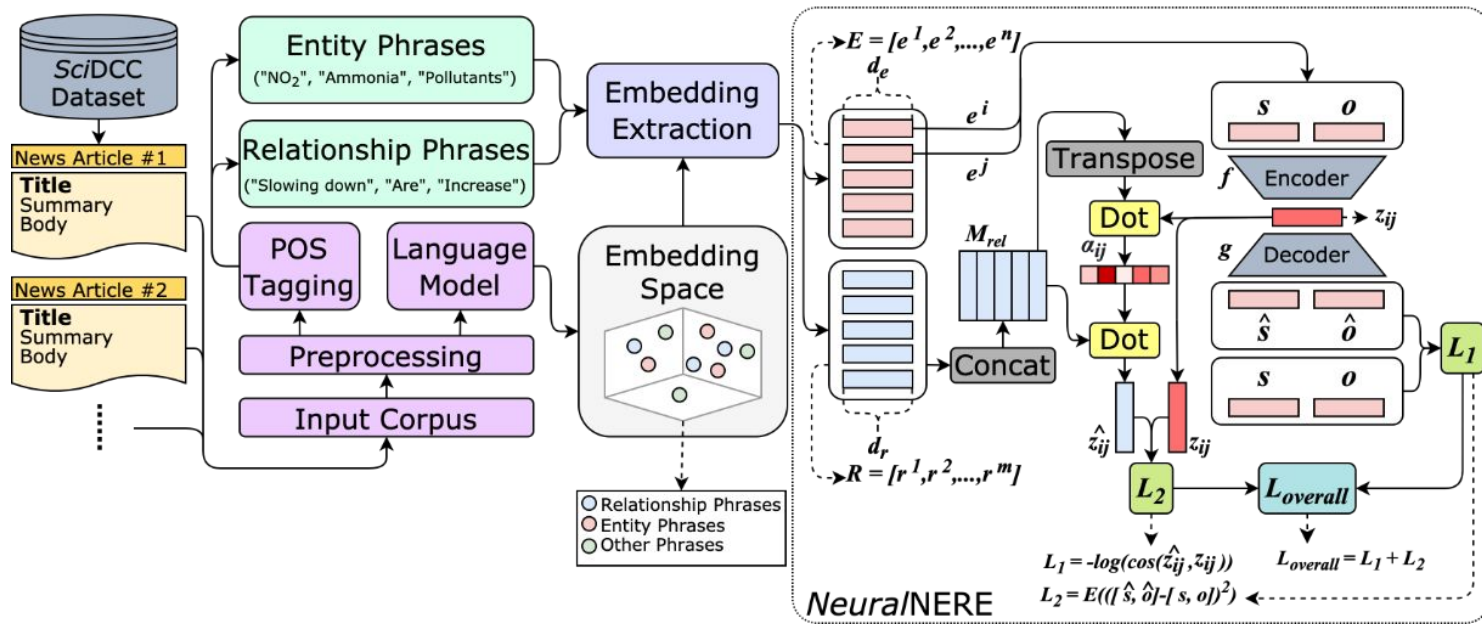
NeuralNERE: Loss Function

- Now we will use the above generated encoded vector \hat{z}_{ij} to enforce the encoded vector z_{ij} to capture relationship-based information. We will do so by modifying the overall loss function.
- The loss function of the proposed model will consist of two terms:
 - The first term will be the reconstruction loss represented by L_1 in Equation (6), which will ensure the reconstruction of input in the encoder-decoder network.
 - Second term will be the cosine similarity loss ($-\log \cos()$) between the two encoded vectors \hat{z}_{ij} & z_{ij} represented by L_2 in Equation (6), which will ensure the learned encoded representation to capture the relationship-based information from the existing relationship phrase representations
- The overall loss function ($L_{overall}$) of the NeuralNERE model will be the addition of the above mentioned individual losses, as shown in Equation (7).

$$L_1 = E(([\hat{s}, \hat{o}] - [s, o])^2); L_2 = -\log \cos(\hat{z}_{ij}, z_{ij}) \quad (6)$$

$$L_{overall} = L_1 + L_2 \quad (7)$$

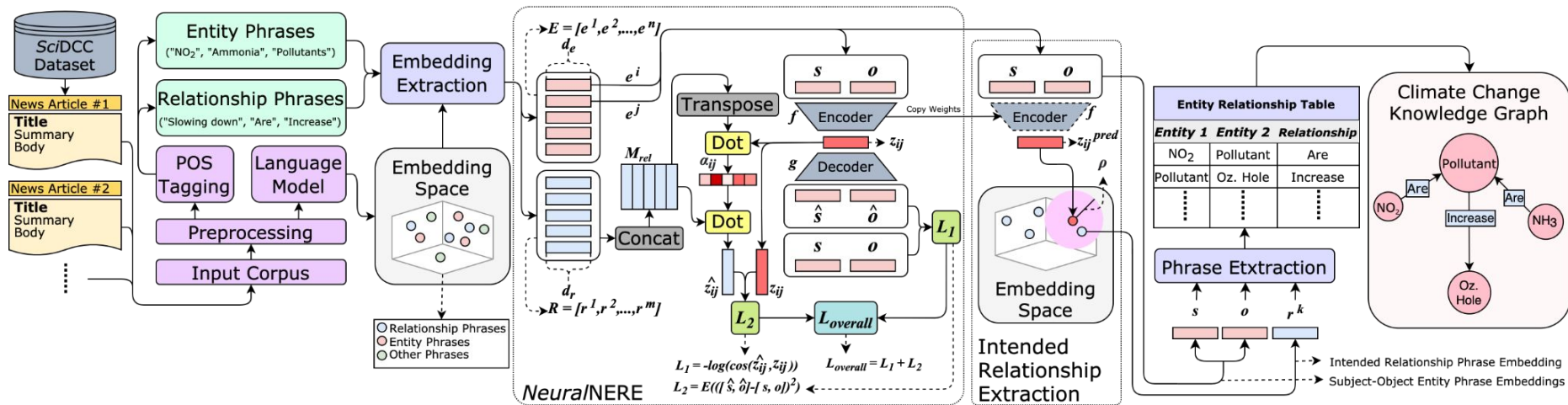
Neura/NERE: Loss Function



$$L_1 = E(\| [\hat{s}, \hat{o}] - [s, o] \|^2); L_2 = -\log \cos(\hat{z}_{ij}, z_{ij}) \quad (6)$$

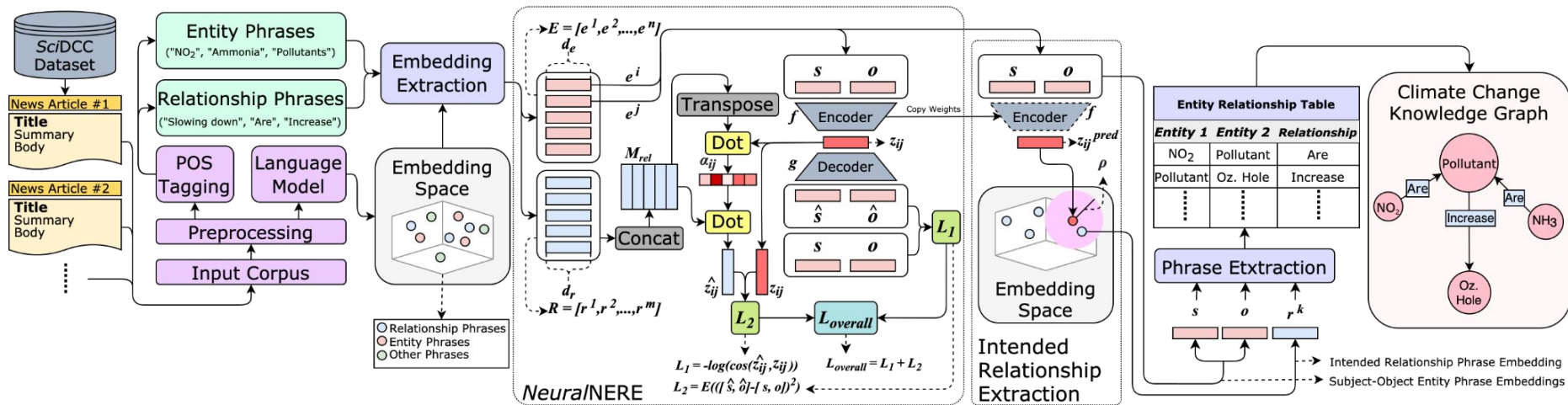
$$L_{overall} = L_1 + L_2 \quad (7)$$

NeuraNERE



- Now we use the trained encoder network f to predict the embedding representation z_{ij}^{pred} of the intended relationship phrase for all the subject-object entity phrase pairs.
- Then the relationship phrase corresponding to the r^k , which has the highest cosine similarities with z_{ij}^{pred} is chosen as the intended relationship phrase. We don't choose any if all the computed cosine similarity values fall below a threshold ρ . Such a threshold keeps the model in check and prohibits the generation of useless relationships

NeuraNERE



- Finally, these triplets comprising of a subject entity phrase, an object entity phrase, and the predicted relationship phrase (from NeuraNERE) are used to construct the climate change knowledge graph.

Conclusion & Future Work

Using the proposed *SciDCC* dataset and *Neura*/NERE model we aim to give industry leaders, analyst, and policymakers a tool for:

1. Extracting and organizing climate change information from a large collection of news articles.
2. Analyzing relationships between different factors responsible for climate change.
3. Gathering insight/reasoning about the pivotal events for more informed climate change policy making.

Thank You