

Location-Specific Information Extraction from Unstructured Biodiversity Literature

Outline

1. Overview

2. Current Progress

3. Next Activities

Why Location-Based Biodiversity Information

- Identifying Species distribution and regional habitat conditions
- more relevant area-specific conservation strategies compared to species-focused data
 - may guide efforts on endangered species and ecosystems on particular locations
- Links geographical contexts to biodiversity data to address further concerns like habitat loss and climate change impacts

Extracting such unstructured information can be achieved through Natural Language Processing

Given that we want to obtain biodiversity information with a focus on location, we want to:

- Extract meaningful entities - through Named Entity Recognition (NER)
- Determine relations among such entities - through Relation Extraction (RE)

Previous studies dealt with information extraction on Unstructured Biodiversity data, where most have explored NER and several attempts on RE. Also, most Information Extraction primarily focuses on other specific data, such as species occurrence extraction.

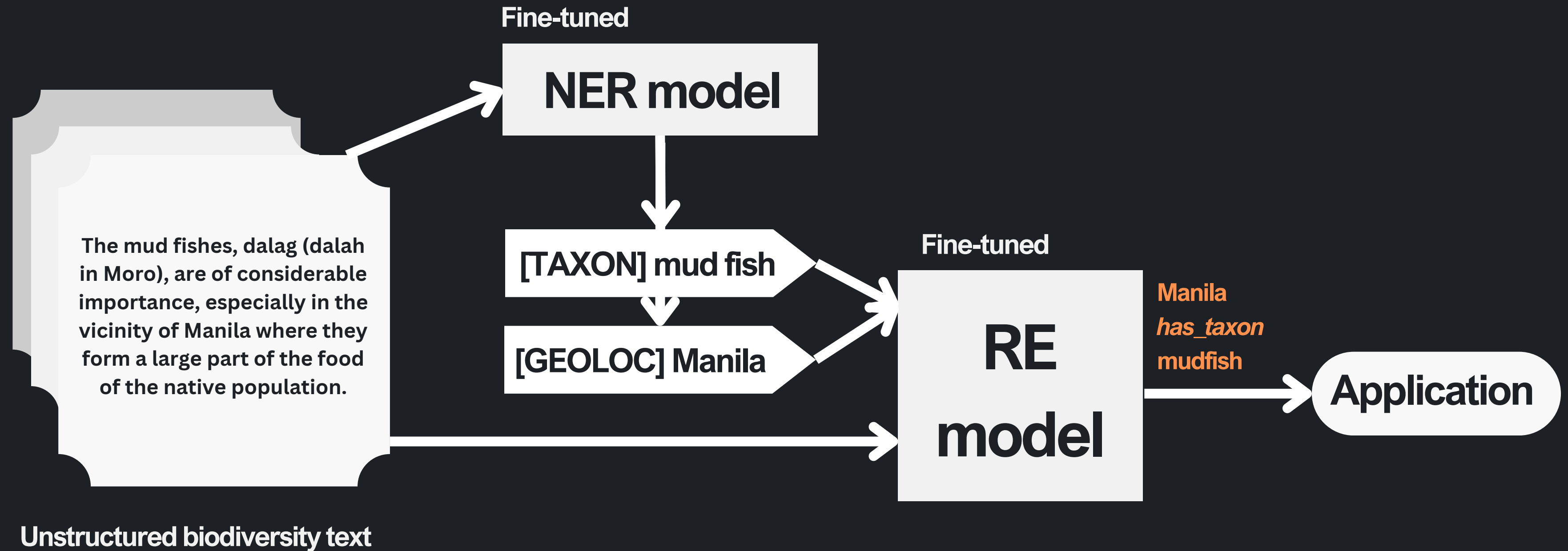
The lack of annotated datasets for RE is a challenge that invites further exploration and experimentation.

Hence, this study focuses on such gaps, particularly on location-based Information Extraction on biodiversity texts.

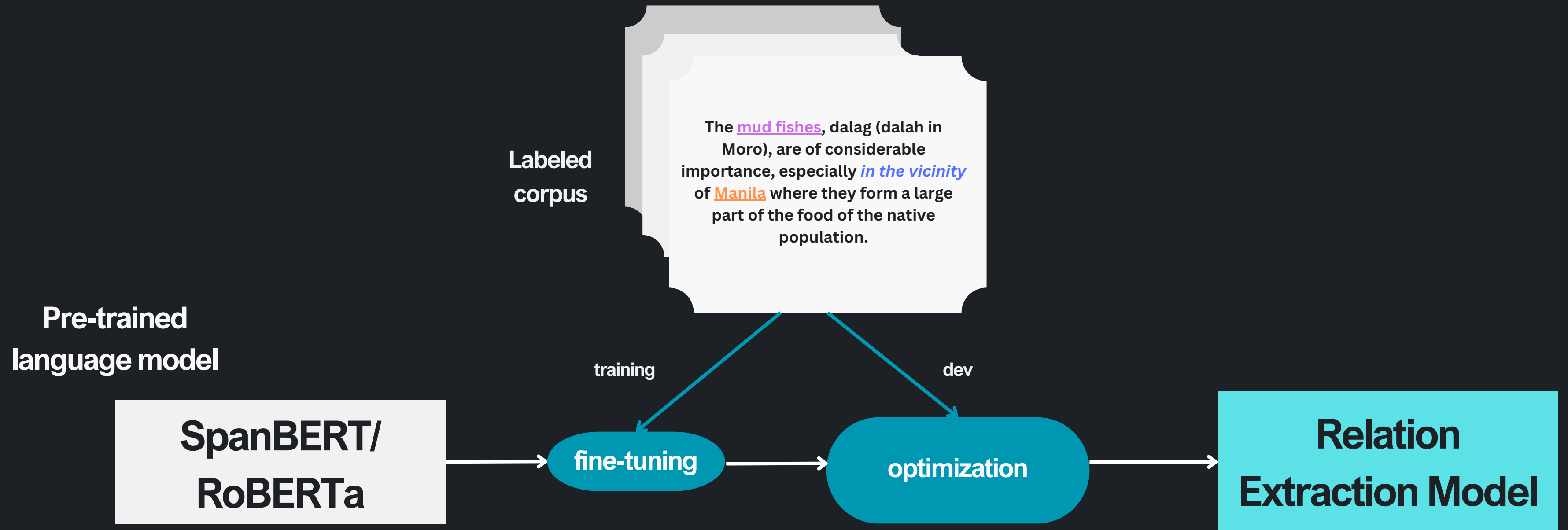
Objectives

1. Annotate existing NER-labeled biodiversity corpus for supervised relation extraction task
2. Train various RE models given the annotated dataset
3. Evaluate trained RE model performance and integrate model for application

Pipeline



Model Training



PLMs

Transformer Architectures - (faster and deeper bidirectional context understanding) - good performance on sequence-to-sequence problems

BiodivBERT (NER)

- Followed the same pre-training procedure with BERT, but considered a great amount of biodiversity-related text.
- Approx. 5 Gb of training data
- Has shown good performance for biodiversity NER, where this study would continue to further optimize its hyperparameters (2024)

SpanBERT (RE)

- Aimed to improve BERT by using span-based masking (BERT uses token-based masking) for its masked language modeling (MLM) on pre-training
- Optimizes MLM by contiguous approach on masking with the Span Boundary Objective
- trained on Wikipedia and BookCorpus like BERT (16Gb)

RoBERTa (RE)

- Aimed to improve BERT by removing Next Sentence Prediction component.
- Trained on a larger corpus, approx. 160Gb of data compared to BERT with 16 Gb, and trained on longer sequences

SpanBERT

$$\begin{aligned}\mathcal{L}(\text{football}) &= \mathcal{L}_{\text{MLM}}(\text{football}) + \mathcal{L}_{\text{SBO}}(\text{football}) \\ &= -\log P(\text{football} \mid \mathbf{x}_7) - \log P(\text{football} \mid \mathbf{x}_4, \mathbf{x}_9, \mathbf{p}_3)\end{aligned}$$

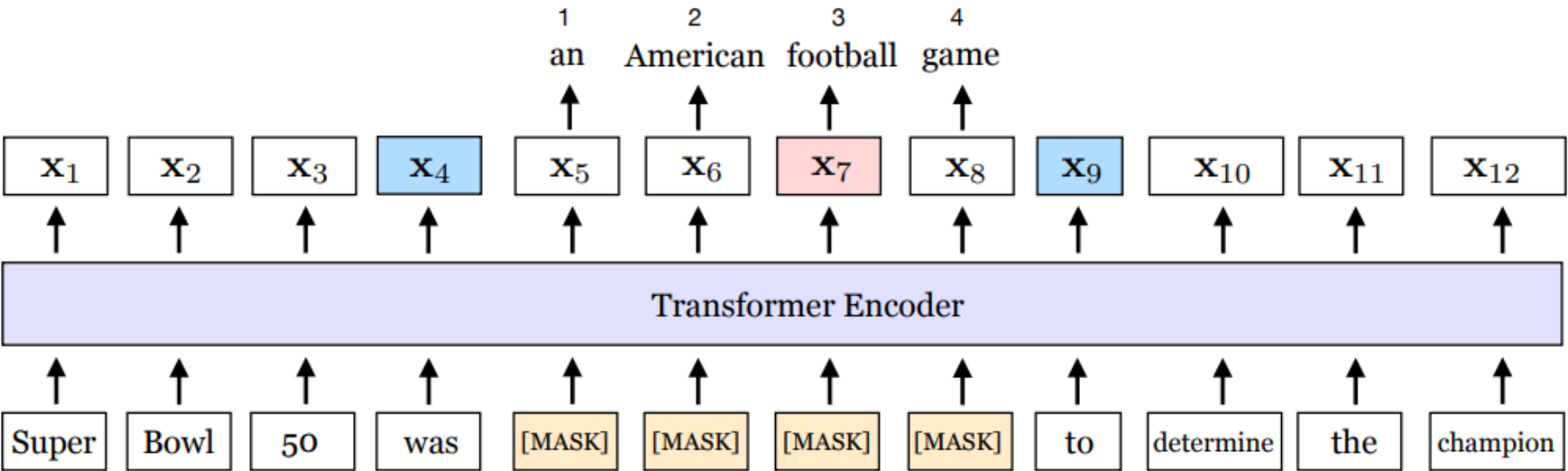


Figure 1: An illustration of SpanBERT training. The span *an American football game* is masked. The span boundary objective (SBO) uses the output representations of the boundary tokens, x_4 and x_9 (in blue), to predict each token in the masked span. The equation shows the MLM and SBO loss terms for predicting the token, *football* (in pink), which as marked by the position embedding p_3 , is the *third* token from x_4 .

This approach forces the model to learn relationships between spans and their boundaries, which is critical for tasks like relation extraction.

On pre-training

$$\begin{aligned}\mathcal{L}(\text{football}) &= \mathcal{L}_{\text{MLM}}(\text{football}) + \mathcal{L}_{\text{SBO}}(\text{football}) \\ &= -\log P(\text{football} \mid \mathbf{x}_7) - \log P(\text{football} \mid \mathbf{x}_4, \mathbf{x}_9, \mathbf{p}_3)\end{aligned}$$

A span boundary objective (SBO) is introduced that involves predicting each token of a masked span using only the representations of the observed tokens at the boundaries.

I love [*chocolate and*] vanilla ice cream
I love [MASK MASK] vanilla ice cream

chocolate is predicted by

- Span Boundary Objective:
- Left boundary: Embedding of "love"
 - Right boundary: Embedding of "vanilla"
 - Positional embedding: The relative position of "chocolate" in the span).

MLM full sentence context - embeddings of all tokens in the sequence

SpanBERT

2.4. Relation Extraction

	p	R	F1
BERT _{EM} (Soares et al., 2019)	—	—	70.1
BERT _{EM} +MTB*	—	—	71.5
Google BERT	69.1	63.9	66.4
Our BERT	67.8	67.2	67.5
Our BERT-1seq	72.4	67.9	70.1
SpanBERT	70.8	70.9	70.8

Test performance on the TACRED relation extraction benchmark

SpanBERT exceeds the reimplementation of BERT by 3.3% F1 and achieves close to the current state of the art.

Person

Billy Mays, the bearded, boisterous pitchman who, as the undisputed king of TV yell and sell,

per:city_of_death

City

became an unlikely pop culture icon, died at his home in Tampa, Fla, on Sunday.

Person

Pandit worked at the brokerage Morgan Stanley for about 11 years until 2005, when he and some

org:founded_by

Organization

Morgan Stanley colleagues quit and later founded the hedge fund Old Lane Partners.

Person

He

no_relation

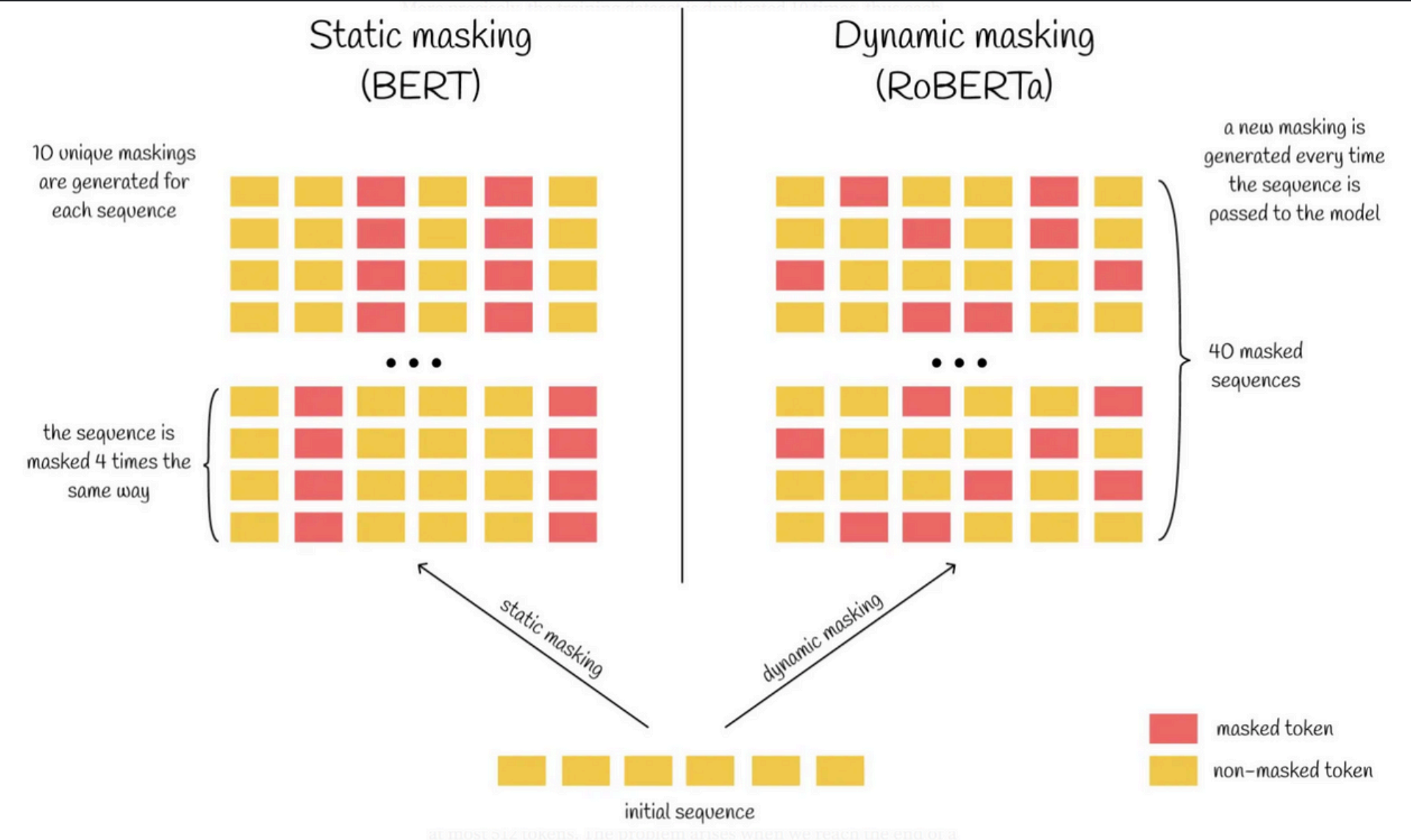
no_relation

Organization

admission to graduate school at the University of Maryland in College Park.

TACRED: Large Scale RE dataset - 41 relation types

RoBERTa



Dynamic masking

tokens to be masked are chosen randomly in every training epoch. The model is exposed to a wider range of masking patterns for the same input sentences, making it more robust to varied contexts; better generalizations

Dataset: COPIOUS (2019)

A gold-standard corpus specifically designed to extract species occurrences from biodiversity literature. Primarily aligned for entity extraction tasks, the corpus includes manually annotated entities such as (1) taxon names, (2) geographical locations, (3) habitats, (4) temporal expressions, and (5) person names.

```
107004_English_78294_34354172_1912.txt X
107004_English_78294_34354172_1912.txt
1 20 HOLLISTER.
2 Kerivouta whitehead! Thomas.
3 1894. Kerivoula whiteheadi Thomas, Ann. and Mag. Nat. Hist, VI, 14, 460.
4 Type locality. – Isabel, N. E. Luzon. Luzon (Thomas). Family MOLOSSID.
5 Genus CHAEREPHON Dobson.
6 1874. Chaerephon Dobson, Journ. Asiatic Soc. Bengal, 43, pt. 2, 144. Type. – Nyctinomiis johorensis
7 Bats with ears connected by a low band across crown; tragus very small; tail long, free from in
8 Chserephon plicatus (Buchanan).
9 1800. Vespertilio plicatus Buchanan, Trans. Linn. Soc, 5, 261. 1907. Chaerephon plicatus Ander
10 Type locality. – Puttahaut, Bengal. Philippine Islands (Dobson) ; Luzon (Elera)
11 Order CARNIVORA.
12 Family MUSTELID.
13 Genus MARTES Pinel.
14 1792. Martes Pinel, Actes Soc. Hist. Nat, Paris, 1, 55, footnote. Type. – Martes domestica Plne
15 Martens, the single supposed Philippine species about the size of a domestic cat; head and body
16 Martes henricii (Westerman).
17 1848. Mustela (Martes) henricii Westerman, Bijdragen tot de Dier-kunde, 1, 13.
18 Type locality. – Java. Sulu (Trouessart). Sanchez, Los Mamiferos de Filipinas, does not allow
19 |
```

```
107004_English_78294_34354172_1912.ann
1 T1 Taxon 290 312 Nyctinomiis johorensis
2 T2 Taxon 651 670 Chserephon plicatus
3 T3 Taxon 1068 1085 Plnelmarfes foia
4 T4 Taxon 1444 1459 Martes henricii
5 T5 Taxon 14 33 Kerivouta whitehead
6 T6 Taxon 49 69 Kerivoula whiteheadi
7 T7 Taxon 49 90 Kerivoula whiteheadi Thomas, Ann. and Mag
8 T8 Geographicallocation 133 140 Isabel
9 T9 Geographicallocation 148 153 Luzon
10 T10 Geographicallocation 155 160 Luzon
11 T11 Taxon 178 186 MOLOSSID
12 T12 Taxon 194 211 CHAEREPHON Dobson
13 T13 Taxon 219 236 Chaerephon Dobson
14 T14 Taxon 219 243 Chaerephon Dobson, Journ
15 T15 Taxon 290 319 Nyctinomiis johorensis DOBSON
16 T16 Taxon 651 681 Chserephon plicatus (Buchanan)
17 T17 Taxon 689 709 Vespertilio plicatus
18 T18 Taxon 689 719 Vespertilio plicatus Buchanan
19 T19 Taxon 689 726 Vespertilio plicatus Buchanan, Trans
20 T20 Taxon 753 772 Chaerephon plicatus
```

Dataset: COPIOUS (2019)

- Only has NER annotations for training, hence the project would also require and aim to obtain relation annotations for supervised training of a relation extraction model
- (*GeographicalLocation-Taxon*, *Geographical-Habitat*) relations

107004_English_78294_34354172_1912.txt X

107004_English_78294_34354172_1912.txt

```
1 20 HOLLISTER.
2 Kerivouta whitehead! Thomas.
3 1894. Kerivoula whiteheadi Thomas, Ann. and Mag. Nat. Hist, VI, 14, 460.
4 Type locality. – Isabel, N. E. Luzon. Luzon (Thomas). Family MOLOSSID.
5 Genus CHAEREPHON Dobson.
6 1874. Chaerephon Dobson, Journ. Asiatic Soc. Bengal, 43, pt. 2, 144. Type. – Nyctinomiis johorensis
7 Bats with ears connected by a low band across crown; tragus very small; tail long, free from in
8 Chserephon plicatus (Buchanan).
9 1800. Vespertilio plicatus Buchannan, Trans. Linn. Soc, 5, 261. 1907. Chaerephon plicatus Ander
10 Type locality. – Puttuhaut, Bengal. Philippine Islands (Dobson) ; Luzon (Elera)
11 Order CARNIVORA.
12 Family MUSTELID.
13 Genus MARTES Pinel.
14 1792. Martes Pinel, Actes Soc. Hist. Nat, Paris, 1, 55, footnote. Type. – Martes domestica Plne
15 Martens, the single supposed Philippine species about the size of a domestic cat; head and body
16 Martes henricii (Westerman).
17 1848. Mustela (Martes) henricii Westerman, Bijdragen tot de Dier-kunde, 1, 13.
18 Type locality. – Java. Sulu (Trouessart). Sanchez, Los Mamiferos de Filipinas, does not allow
19
```

107004_English_78294_34354172_1912.ann

```
1 T1 Taxon 290 312 Nyctinomiis johorensis
2 T2 Taxon 651 670 Chserephon plicatus
3 T3 Taxon 1068 1085 Plnelmarfes foia
4 T4 Taxon 1444 1459 Martes henricii
5 T5 Taxon 14 33 Kerivouta whitehead
6 T6 Taxon 49 69 Kerivoula whiteheadi
7 T7 Taxon 49 90 Kerivoula whiteheadi Thomas, Ann. and Mag
8 T8 Geographicallocation 133 140 Isabel
9 T9 Geographicallocation 148 153 Luzon
10 T10 Geographicallocation 155 160 Luzon
11 T11 Taxon 178 186 MOLOSSID
12 T12 Taxon 194 211 CHAEREPHON Dobson
13 T13 Taxon 219 236 Chaerephon Dobson
14 T14 Taxon 219 243 Chaerephon Dobson, Journ
15 T15 Taxon 290 319 Nyctinomiis johorensis DOBSON
16 T16 Taxon 651 681 Chserephon plicatus (Buchanan)
17 T17 Taxon 689 709 Vespertilio plicatus
18 T18 Taxon 689 719 Vespertilio plicatus Buchannan
19 T19 Taxon 689 726 Vespertilio plicatus Buchannan, Trans
20 T20 Taxon 753 772 Chaerephon plicatus
```

Current Progress

- **Dataset Pre-processing for Annotations**
- **Dataset Annotation**
- **Model Training Preparation**

Current Progress

CS 198 - CVMIG



```
from google.colab import drive
drive.mount('/content/drive')

# folder path
folder = "/content/drive/My Drive/CS_198/copious_published/copious_published/"
subfolders = ["train", "dev", "test"]

Mounted at /content/drive

[2] ## Single sentence
...
import os
import pandas as pd
import nltk

# Ensure NLTK has the required tokenizer
nltk.download('punkt_tab')

# Initialize storage for all entities
all_entities = []
```

Dataset Preparation

- Prepare for RE annotation
- Explore Dataset Distribution of entity and relations
 - stored in .xlsx files
- Annotating the existence of a binary relation between two entities is suggested as the approach is simple

Current Progress

Dataset Distribution: Entities

Category	Train	Dev	Test
Taxon	9,357	1,548	1,322
Geographical Location	8,121	992	878
Person	2,479	180	230
Temporal Expression	1,800	157	253
Habitat	1,308	91	115

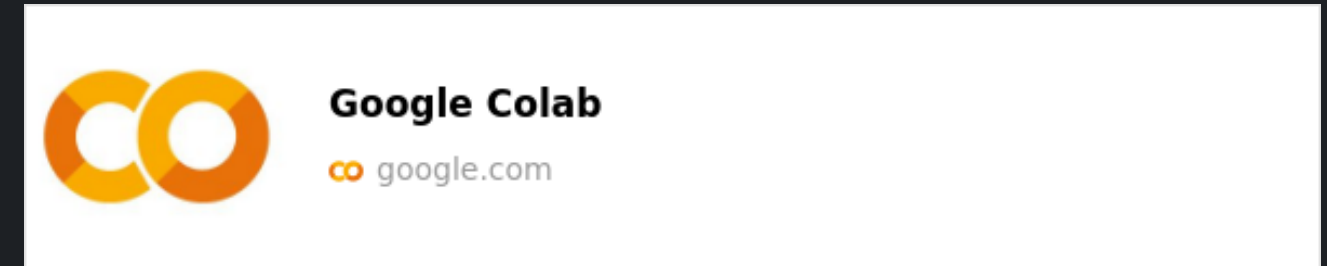
Current Progress

Dataset Distribution: Relation Pairs

Subfolder	Geolocation-Taxon	Geolocation-Habitat
train	2437	691
dev	106	30
test	313	64

Current Progress

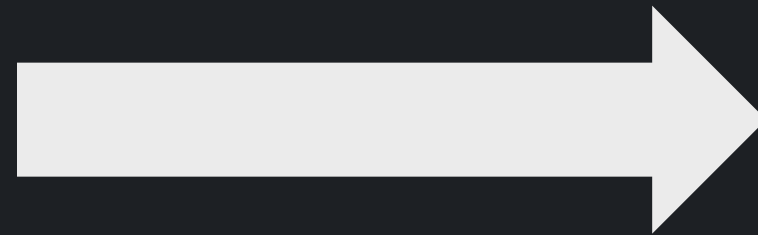
Dataset pre-processing for relation labeling



.txt file



.ann file



Spreadsheet format for relation labeling

Current Progress

Dataset pre-processing for relation labeling

.txt file

A BIOLOGICAL AND SYSTEMATIC STUDY OF PHILIPPINE PLANT GALLS 1

By Leopoldo B. Uichanco (From the College of Agriculture, University of the Philippines, Los Banos)

FIFTEEN PLATES INTRODUCTION

Galls are abnormal growths on the stems, leaves, roots, or other parts of plants, caused by the action of insects, arachnids, or fungi, or by unknown agencies. Just how these peculiar structural developments are brought about is still open to discussion and speculation, experimental proofs being, up to the present, too deficient to warrant our drawing any definite conclusion. These malformations have been ascribed to various causes, the more commonly accepted, in the absence of more reasonable, explanations being the following: 2 1, a severe mechanical injury to certain parts of the plant ; 2, a continuous mechanical irritation; 3, secretion of chemical stimulus by the causal animal or fungus. One, or a combination of two or all ...

.ann file

```
T1 GeographicalLocation 151 160 Los Banos
T2 Taxon 1314 1336 Aulax papaveris Perris
T3 Taxon 296 303 insects
T4 Taxon 305 314 arachnids
```

3-sentence chunks

(S1, S2, S3), (S2, S3, S4), (S3, S4, S5), ...



Spreadsheet rows for each entity, entity type, and chunks

File	Chunk	Entity ID	Entity Type	Entity Text	Spans
111121_English_36025311_1908	484SMITH. THE SULU ARCHIPELAGO. The Sulu Archipelago	T1	GeographicalLocation	Sulu Archipelago	36 52
111121_English_36025311_1908	THE SULU ARCHIPELAGO. The Sulu Archipelago is practicall	T2	Person	Becker	215 221
111121_English_36025311_1908	THE SULU ARCHIPELAGO. The Sulu Archipelago is practicall	T3	GeographicalLocation	Marongas Island	270 285

Current Progress

Dataset pre-processing for relation labeling

Spreadsheet rows for each entity, entity type, and chunks

File	Chunk	Entity ID	Entity Type	Entity Text	Spans
111121_English_36025311_1908	484SMITH. THE SULU ARCHIPELAGO. The Sulu Archipelago	T1	GeographicalLocation	Sulu Archipelago	36 52
111121_English_36025311_1908	THE SULU ARCHIPELAGO. The Sulu Archipelago is practicall	T2	Person	Becker	215 221
111121_English_36025311_1908	THE SULU ARCHIPELAGO. The Sulu Archipelago is practicall	T3	GeographicalLocation	Marongas Island	270 285



Filtering entities, pairing desired entities, readying a column for relation labeling

File	Entity 1	Entity 1 Type	Entity 2	Entity 2 Type	Chunk	relation
104421_English_33451123_1922-26.	Asia	GeographicalLocation	pantropic	Habitat	118, Loher 7213. Planted about r	
104421_English_33451123_1922-26.	Luzon	GeographicalLocation	GRAMINEAE	Taxon	i VOL GRAMINEAE 49	
104421_English_33451123_1922-26.	Luzon	GeographicalLocation	open grasslands	Habitat	i VOL GRAMINEAE 49	
104421_English_33451123_1922-26.	Abra	GeographicalLocation	open grasslands	Habitat	Luzon (Abra, Ilocos Norte, Pangas	
104421_English_33451123_1922-26.	Ilocos Norte	GeographicalLocation	open grasslands	Habitat	i VOL GRAMINEAE 49	
104421_English_33451123_1922-26.					Luzon (Abra, Ilocos Norte, Pangas	

Current Progress

Dataset Feature

- Hence, sentences are further grouped into chunks to at least contain two desired entities to classify a binary relation
- One such observation for the COPIOUS dataset is that two entities are not usually captured within a single sentence

	A	B	C	D	E	F	G	H	I	
influence	1	...	soil	density	varies	in	response	to	earth	
	-	O	B-Quality	I-Quality	O	O	O	O	O	B-Ph
have	1	functional	richness	...	on	ecosystem	states	and	proce	
	-	B-Quality	I-Quality	O	O	B-Environment	O	O	O	O
occur_in	1	...	and	diversity	of	chemolithoautotrophic	bacteria	in	saline	
	-	O	O	O	O	B-Organism	I-Organism	O	O	B-En
have	1	ant	abundance	and	diversity	associated	with	natural	habita	
	-	B-Quality	I-Quality	O	O	O	O	B-Environment	I-Env	
	0	for	understanding	how	microbial	communities	degrade	plant	bioma	
	-	O	O	O	B-Organism	I-Organism	O	B-Quality	I-Qua	
	0	density	of	soil	invertebrates	in	response	to	earth	
	-	B-Quality	O	O	O	O	O	O	O	B-Ph

BiodivNERE dataset

Sentence-level relations and annotations, short spans

The relevant predictions from Heaney et al. (1989) derived from our studies on Leyte and **Negros** are (1) in forested habitats, species richness and relative abundance of fruit bats should be highest in the **lowlands** and decline with increasing elevation, (2) species richness of fruit bats should be lower in **agricultural areas**, where geographically widespread species should predominate, (3) non-volant mammals should not exhibit highest species richness in **lowland forest**, (4) abundance of non-volant mammals should increase with elevation, and (5) endemic species of fruit bats and non-volant mammals should occur in primary forest, whereas non-endemics should predominate in disturbed habitats.

COPIOUS dataset

Entity pairs are not usually captured in a single sentence; longer spans

Current Progress

CS 198 - CVMIG

Dataset Annotation

- Orientation for relation labeling was conducted today (November 25, 2024)
 - Results are expected next week (December 2) to start model training

Relation Labels

has_taxon: Geolocation-Taxon pairs

has_habitat: Geolocation-Habitat pairs

Binary Classification of entity pair relation

Determine whether the sentence explicitly or implicitly established a **biological or ecological connection** with the entities

- *Tilapia* thrives in lakes near *Manila*. → *has_taxon* (*Manila*, *Tilapia*)
- *Tilapia* sold in *China* markets is imported from *Manila*. → *no_relation* (*China*, *Tilapia*)

Current Progress

Dataset Annotation

- Two groups (A and B) of five CWTS students are annotating the data

- Group A & B (member 1)
- Group A & B (member 2)
- ⋮
- Group A & B (member 5)

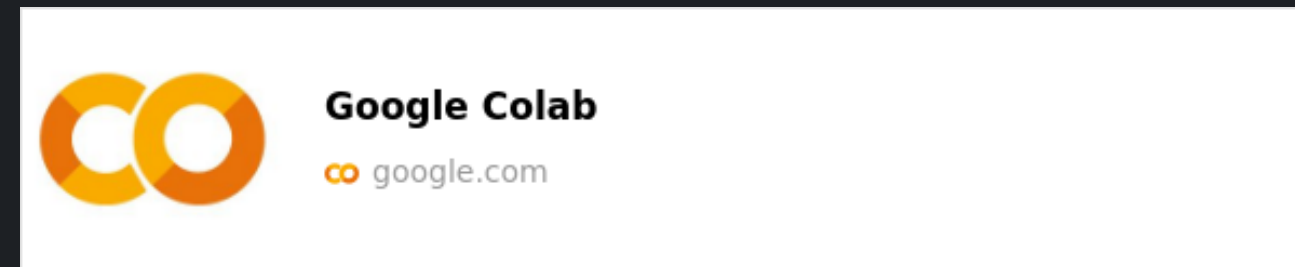
train.xlsx	dev.xlsx	test.xlsx
2–625	2–27	2–75
626–1250	28–54	76–150
1251–1875	55–81	151–225
1876–2500	82–108	226–300
2501–3128	109–137	301–378

- Two people would work independently on the same chunks per file to compare annotations
- for consolidation of final labels upon evaluation

Current Progress

CS 198 - CVMIG

Model Training Preparation



RE labeled dataset

—	—	—	—	—
—	—	—	—	—
—	—	—	—	—
—	—	—	—	—
—	—	—	—	—
—	—	—	—	—
—	—	—	—	—
—	—	—	—	—

Training code is being polished,
prepared for labeled inputs and
ready for training



```
# Preprocessing
def preprocess_data(batch, tokenizer, model_name, label_to_id, max_seq_length):
    inputs = []
    for chunk, entity1, entity2, entity_type1, entity_type2 in zip(
        batch["chunk"], batch["entity1"], batch["entity2"], batch["entity_type1"], batch["entity_type2"]
    ):
        if "spanbert" in model_name:
            # SpanBERT-specific input formatting: Similar for BERT formattings
            input_text = chunk.replace(entity1, f"[{entity_type1}]").replace(entity2, f"[{entity_type2}]")
        else:
            # RoBERTa-specific input formatting: not really different since this was updated to follow the procedure for input formatting
            input_text = chunk.replace(entity1, f"[{entity_type1}]").replace(entity2, f"[{entity_type2}]")
        inputs.append(input_text)

    encodings = tokenizer(inputs, padding=True, truncation=True, max_length=max_seq_length)
    encodings["labels"] = [label_to_id[label] for label in batch["relation"]]
    return encodings
```

Next Activities

1. Annotation Evaluation and Consolidation - By December 2
2. Model Training and Fine-tuning - By December 3
3. Concerns with the amount of training data and model performance
 - Add more training data
 - Other models
4. Start building the NER segment for implementation purposes