

WILFRID LAURIER UNIVERSITY

DEPARTMENT OF MATHEMATICS

MA695 - MAJOR RESEARCH PROJECT

---

# Enhancing Low-Resource Named Entity Recognition Results with Data Augmentation via Abstractive Summarization

---

*Author*  
Douglas BOWEN

*Supervisors:*  
Dr. Yang LIU  
Dr. Xu (Sunny) WANG

*Second Readers:*  
Ilias S. KOTSIREAS

July 24th, 2022



## Abstract

News and other articles can provide valuable info - but need to be read, analyzed, and processed. In the past, this was done manually and could be costly and time-consuming. In recent years, Natural Language Processing (NLP) has become mainstream in its ability to extract a variety of information from structured and unstructured text, without human intervention. This provides a much quicker and more cost-effective tool than manual labour.

Named Entity Recognition (NER) is one such information extraction method, allowing for the classification of named entities to be made automatically. Due to its unique characteristics and the large variance in different domains, NER has been widely studied in the past few decades.

For models to extract named entities from articles, a large amount of NER labelled data is required to learn from and sufficiently make predictions. However, finding data that is labelled can be challenging. Though some common tag sets and domains have plenty of data available to train on already that have been compiled by the public (e.g. Wikipedia, Twitter, News), custom NER tags and niche domains often have little to no data available, making training quite challenging. As a result, steps must be taken to create a large dataset from scratch or a small low-resource set. Data augmentation (DA) is a technique that attempts to alleviate this. Through various approaches, individuals can create artificial data to supplement the original and help increase NER model performance without the need for costly human data curation.

This project examines various approaches to data augmentation in low-resource domains wherein artificially produced training data is required for decent model training. The dataset simulates a low-resource scenario for data augmentation, and a new generation approach using abstractive summarization is used to create sentences with completely new structures from the original data.

Techniques to create a sufficient sample set of augmented NE data include (but are not limited to) weighting source article sentences by entity count prior to shuffling, cutting sentences with no entities, and shuffling sentences randomly to create new samples.

Results show that utilizing an abstractive summarization technique to augment data provides a significant boost over the original unaugmented data for low-resource data. Summary results are similar to paraphrased results, but lower than rules-based results. As the data size increases, the boost provided diminishes for all methods as expected.

## Acknowledgements

I would like to extend my sincere thanks to everyone who provided support during this project. Firstly, to my advisors Dr. Xu (Sunny) Wang and Dr. Yang Liu, for supporting my efforts and for providing their expertise as well as helping guide my research with valuable suggestions and insights. Secondly, to my colleagues Yanan Chen, Ian Fraser, Jason Lechter, and Herteg Kohar who worked with me, helped in researching, and provided support in tackling issues throughout the project. Lastly, I would also like to acknowledge the funding that was provided by the Natural Science and Engineering Research Council (NSERC), the MITACS Accelerate Program, and The Quantum Insider (TQI).

# Contents

<b>1</b>	<b>Introduction</b>	<b>6</b>
<b>2</b>	<b>Literature Review</b>	<b>8</b>
2.1	NER-Focused Approach . . . . .	8
2.2	Rule-Based Approaches . . . . .	9
2.3	Basic Generation Approach . . . . .	10
2.4	Cross-Domain Generation Approach . . . . .	11
2.5	Paraphrasing . . . . .	12
<b>3</b>	<b>Background</b>	<b>13</b>
3.1	Evaluation Metrics . . . . .	13
3.2	Common Model Structure . . . . .	13
3.3	Transformers . . . . .	15
3.4	Model Hyper-Parameter Choices . . . . .	17
3.5	Known Abstractive Summarization Issues . . . . .	21
<b>4</b>	<b>Methodology</b>	<b>25</b>
4.1	Proposed Summary Generation Approach . . . . .	25
4.2	Overview of Process Pipeline . . . . .	26
4.3	Models for Abstractive Summarization . . . . .	26
4.4	Named Entity Recognition Models . . . . .	30
4.5	Data Source Adjustments . . . . .	30
4.6	Model Evaluation & Comparison . . . . .	35
<b>5</b>	<b>Experimental Results and Analysis</b>	<b>36</b>
5.1	The WikiGold Dataset . . . . .	36
5.1.1	Dataset Requirements & Selection . . . . .	36
5.1.2	Low-Resource Mimicry Selection . . . . .	37

5.1.3	Descriptive Statistics . . . . .	37
5.2	Data Augmentation Method Results . . . . .	39
5.3	NER Model Settings . . . . .	39
5.4	Model Results . . . . .	40
<b>6</b>	<b>Future Work</b>	
	<b>and Conclusions</b>	<b>42</b>
6.1	Future Work . . . . .	42
6.2	Project Conclusion . . . . .	43
<b>7</b>	<b>Bibliography</b>	<b>45</b>
<b>8</b>	<b>Appendices</b>	<b>48</b>
8.1	Code Files . . . . .	48
8.1.1	Main Code Files . . . . .	48
8.1.2	Secondary Code Files . . . . .	91

# List of Tables

3.1	Example of Summary Degeneration . . . . .	22
4.1	Summary Output Success Rates . . . . .	29
4.2	Example of Article Entity Replacements . . . . .	31
4.3	Example of Sliding n-gram Mapping . . . . .	32
5.1	Metrics by Training Test Split . . . . .	38
5.2	Model Results for Tested Methods ( $\alpha = 0.05, n=10$ ) . . . . .	41

# List of Figures

2.1	Rule-Based Method Examples [1]	10
3.1	Basic Seq2Seq Model Architecture [20]	14
3.2	Encoder-Decoder Details [20]	15
3.3	Transformer Encoder-Decoder Architecture [20]	16
3.4	Greedy Search Generation [18]	18
3.5	Beam Search Generation [18]	19
3.6	Sampling Methods [18]	20
5.1	Distribution of Sentence Count per Article	38



# Chapter 1: Introduction

In the digital age news comes out at an incredible rate. Be it live tweets, online articles, published research, or anything in-between, so much information is coming in at one time and has become impossible for humans to process it all.

In some domains (such as quantum technology), the amount of data has steadily increased with the number of news releases growing every day. Subsequently, information from these releases can often require domain expertise and thus cannot be quickly understood without sufficient background research.

Given the timely and costly manual labour required to review and research for a proper understanding in these domains, along with the general size of training data required for NLP tasks, data augmentation (DA) is a growing field of research that aims to lessen said burden. DA utilizes a variety of techniques that create "new" artificially created data points from the original data. However, Named Entity Recognition (NER) is one such information extraction task that has more stringent requirements in both the structure of the data and the information required from models. For NER tasks, the structure of a sentence is considered and from it, words are designated either as an entity or not. The sequence of the sentence provides the information to designate an entity and is thus important to maintain. As a result, general DA techniques that have been developed for other tasks tend not to apply as well to NER ones due to the contextual requirements.

This project aims to examine data augmentation techniques that can be utilized for NER models in low-resource domains. Though entity classifications are not complicated, the information surrounding them can be confusing for computers. Examining the following snippet from a quantum technology article, "demonstrating the early application of NISQ so-

lutions to the energy sector", a model would likely consider NISQ to be a company rather than scientific acronym if it was not trained in the domain. An augmented version using generation could alter the sentence to the form "NISQ solutions demonstrated promising results within the energy sector" which is still coherent in structure as required by NER models and yet provides new sentence structure to learn on. By examining and applying DA techniques, the time and subsequent cost that arise from requiring individuals with sufficient expertise to provide labelled data can be reduced significantly. With this information, inferences and changes in the industry can be more readily tracked without constant monitoring.

The new augmentation method being examined relies on generating abstractive summaries as a means of creating sentences with new structure as compared to the original source article. Furthermore, approaches are taken to attempt giving priority to entities in the generated sentences. By utilizing summarization models, the generated text will hopefully contain a proper coherent structure with entities from the original text that is sufficiently unique when compared to the other sentences from the original article. To compare the success of this approach, a baseline will be constructed using other researched methods that have been proven successful in augmenting NER sentences. These are discussed in detail in the literature review.

The research performed here will highlight the performance of abstractive summarization as compared to the original dataset, as well as compared to baselines via rule-based approach and a paraphrasing generative method. Results show that performance of abstractive summarization is quite similar to that of paraphrasing, but below that of the rule-based approach for X-Small, Small, and Medium size batches while the Large batch size has similar performance for all three.

The remainder of this thesis is organized in the following order: Chapter 2, which covers the literature review on relevant research done by others; Chapter 3, which provides background information on underlying mechanisms; Chapter 4, which discusses the specific approach to be taken; Chapter 5, which presents the results of the research; Chapter 6, which discusses potential further work; and Chapter 7, which summarizes the paper.

## Chapter 2: Literature Review

This section will address the literature and published research currently available that makes use of data augmentation (DA) techniques on some Named Entity Recognition (NER) process. Though the focus is on utilizing DA for in-domain, low-resource data, the literature of utilizing DA for other purposes is also examined.

### 2.1 NER-Focused Approach

Though many DA techniques have been proposed for general NLP tasks, the application of these DA techniques is not always possible in the case of NER. This is due to the challenge that NER models present, wherein the labelling process is based on each individual token within sentences.

Critical to understanding NER DA methods, it is important to familiarize oneself with the required structure for NER data. There are three primary points of importance: (1) The sentence of a structure, used to gain insight through learned patterns, (2) The words or "tokens" within a sentence that may or may not be entities, and (3) the associated BIO-labels, indicating either the (B)eginning/(I)nside of a named entity, or (O)utside, for non-entities. Of course, there is also the variety of named-entity categories (entity types) associated with the aforementioned labels, however these can be selected as one sees fit and are non-standard across NER tasks. As a result of these tags, techniques that might transform a token have to avoid also changing their labelling which presents a unique challenge for DA in NER tasks.

## 2.2 Rule-Based Approaches

The simplest approach is one that focused on tweaking traditional rules-based NLP approaches to DA so that they may be applied to NER situations [1]. Four underlying rules-based methods for augmenting NER data were constructed, for which combinations of these methods could also be utilized. Each of the following methods utilizes a probability distribution to determine whether or not the suggested transformation should take place on the specified token/segment. A success rate  $p$  is utilized in the binomial distribution (specifically, the bernoulli distribution for  $n=1$  trial).

- The first method discussed is Label-Wise Token Replacement (LWTR) [1] wherein each token within a sentence is evaluated randomly for success or failure. If the randomly generated number outputs a success, the token is replaced with a randomly selected token designated with the same label. Otherwise, no replacement occurs and the next token is evaluated in the same way.
- The second method discussed is Synonym Replacement (SR) [1]. In a very similar method as to LWTR, each token within a sentence is evaluated randomly. In the case of SR however, if the randomly generated result indicates a success, the token is replaced with a synonym of itself. For synonyms that are multi-word, the associated label is applied in the standard B-Entity, I-Entity order.
- The third method discussed is Mention Replacement (MR) [1]. If a token is marked replacement, an entry with the same entity type is swapped into its place. The BIO-label for the new token being replaced should be maintained (i.e. the label should also be replaced).
- The final method discussed is substantially different from LWTR [1], SR, and MR. Shuffle within Segment (SIS) evaluates each continuous segment of labels within a sentence randomly for success or not. In the case of success, the segment will be randomly shuffled to alter the order of tokens, but not labels.

Lastly, it is possible to use any combination/mixture of the above methods for a more comprehensive approach to data augmentation. Figure 2.1

highlights these methods with an example provided by the original researchers.

	Instance												
None	She O	did O	not O	complain O	of O	headache B-problem	or O	any B-problem	other I-problem	neurological I-problem	symptoms I-problem	. O	
LwTR	L. O	One O	not O	complain O	of O	headache B-problem	he O	any B-problem	interatrial I-problem	neurological I-problem	current I-problem	. O	
SR	She O	did O	non O	complain O	of O	headache B-problem	or O	whatsoever B-problem	former I-problem	neurologic I-problem	symptom I-problem	. O	
MR	She O	did O	not O	complain O	of O	neuropathic B-problem	pain I-problem	syndrome I-problem	or O	acute B-problem	pulmonary I-problem	disease I-problem	. O
SiS	not O	complain O	She O	did O	of O	headache B-problem	or O	neurological B-problem	any I-problem	symptoms I-problem	other I-problem	. O	

Figure 2.1: Rule-Based Method Examples [1]

## 2.3 Basic Generation Approach

Another recent approach to augmenting NER data utilizes generational language models that are trained on linearized labelled sentences [2]. This method was applied to both supervised and semi-supervised data wherein both resulted in performances above baseline. Unlike the rules-based approach which simply applies adjustments to data without significantly altering them, generation-based approaches attempt to create sentences that have completely new structures with different words, contexts, phrases, etc., while still remaining coherent to the human eye. Thus the generated sentences still share similar patterns and trends to the original data that a model can find.

Prior to utilizing a language model, sentences have tags added to the beginning and end of sentences, denoted as [BOS] and [EOS] respectively. Furthermore, named entities have their labels added in front of their respective tokens to linearize the sentence. After the completion of this process, a language model can be utilized. In this instance, a one-layer Long-

Short Term Memory (LSTM) Recurrent Neural Network (RNN) is chosen as the language model. To train the generation model, each token is fed into the model in a linearized sentence order such that the first entry is [BOS] and the last entry is [EOS]. Each token then encounters a dropout layer, LSTM layer, another dropout later, and finally produces a prediction for the next token that should occur. With the trained model, synthetic data is generated by first feeding a [BOS] tag into the LSTM RNN, which then outputs a prediction. This prediction is then fed into the model as the next token in the sentence, and so on until sentence completion. This generated sentence is then de-linearized, with labels being added back into the proper order and [BOS]/[EOS] being removed.

Another generation approach, albeit a bit more involved, is one that utilizes a Sequence-to-Sequence (Seq2Seq) Language Model to back-translate a sentence [4]. As a simple overview, the process involves first splitting a sequence of tokens so as to break it into continuous label segments. Then, segments of 3 tokens or more that do not correspond to an entity type are fed into the model for translation from English to German and back again. The result of this back-translation is a sentence slightly different from the original in most case due to how each model handles the translation.

## 2.4 Cross-Domain Generation Approach

One method takes a different approach from the aforementioned methods. Wherein those methods focus on data augmentation in low-resource scenarios that are common in niche fields such as Quantum Technology, this approach attempts to leverage data from a high-resource domain and subsequently project it into the low-resource domain by utilizing semantics and patterns inherent to all text, even if textual patterns differ [3]. This approach builds off the work done in the generation approach.

Similarly to the generation approach, this method begins by linearizing sentences, but then pairs a "source" domain sentence with a "target" domain sentence. Each sentence then has noise added through shuffling, dropouts, or masking. Once this is done, the domain-sentence pairs are fed into the model to output into the paired domain, which allows for the models to learn a "mapping" between domains.

## 2.5 Paraphrasing

A final method is the use of paraphrasing (expressing a sentence's meaning with different word structure) to augment NER data. A Bidirectional Encoder Representations from Transformers (BERT) model is utilized for entity recognition, while the underlying model utilized for the paraphrasing is not present in the published article. Researchers first replaced tokens with their respective entity tags (not "O" tags) and then performed both paraphrasing and back-translation on the sentence to generate newly augmented sentences. Both methods (paraphrasing vs. back-translation) were compared against one another, with paraphrasing proving successful at improving the BERT model at *very* low resource data levels [13].

## Chapter 3: Background

This chapter provides a generalized overview of concepts not formally found in papers, but that are relevant to the research conducted henceforth.

### 3.1 Evaluation Metrics

#### Precision, Recall, & F1

These are standard metrics used for evaluating mistakes made by models. Precision measures the rate of false-positives, recall measures the rate of false-negatives, and F1 showcases the performance taking into account precision and recall trade-offs.

#### Recall-Oriented Understudy for Gisting Evaluation

More commonly known as "ROUGE", this is another important metric that will be used, albeit only in baseline construction and testing of the pre-trained Seq2Seq model. ROUGE evaluates the effectiveness of a summary or translation as compared to a human-produced reference.

### 3.2 Common Model Structure

Seq2Seq models are also commonly referred to as encoder-decoder models. Encoders receive inputs (for example, tokens within a sentence) and use said inputs to find relationships and acquire an "understanding" of the input. This understanding comes in the form of numerical outputs, which



are then sent to the decoder and called "context vectors". Once encoding is complete, the context vectors along with a starting token are passed into a decoder. Decoders then read the context vector and starting token inputs and try to predict outputs, token by token. Figure 3.1 and 3.2 display a very basic Seq2Seq overview. The encoder and decoder segments are in actuality often both built on recurrent neural networks, which are either Gated Recurrent Unit (GRU) or Long-Short Term Memory (LSTM) networks. However, this is not the case for transformers.

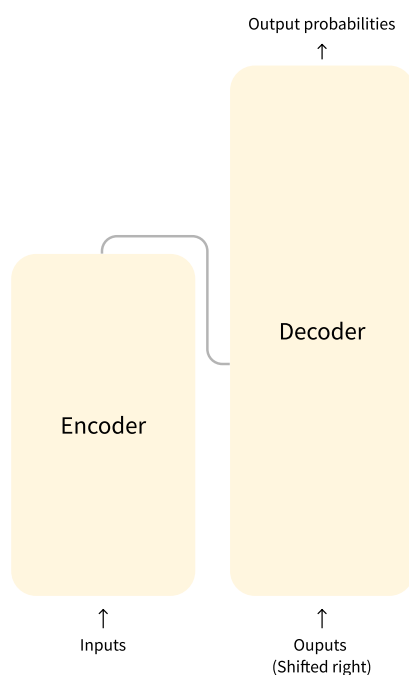


Figure 3.1: Basic Seq2Seq Model Architecture [20]

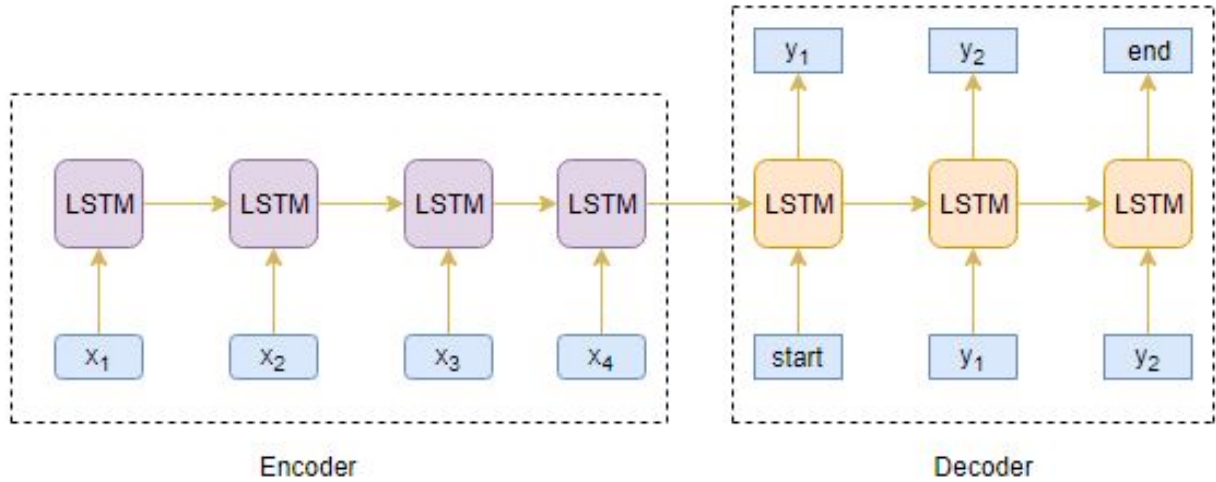


Figure 3.2: Encoder-Decoder Details [20]

### 3.3 Transformers

Transformers are built with "attention" in mind. Attention is a method wherein encoders and decoders are fed only the relevant inputs. To denote these inputs, inputs are assigned weightings where higher weights denote higher importance. These weightings are adjusted over time using feed-forward neural networks. Overall, transformers try to predict an output using only the important parts of the sentence, therefore giving a higher degree of "attention" to important input words over the others. The transformer process is very involved and thus more easily visualized as in Figure 3.3.

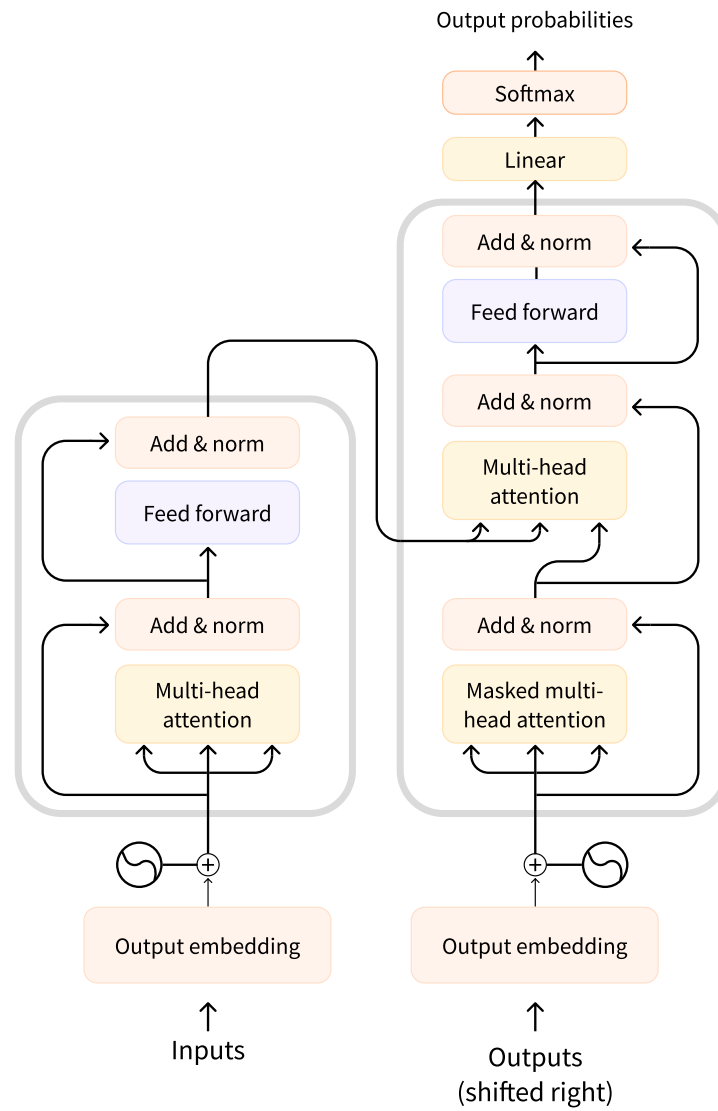


Figure 3.3: Transformer Encoder-Decoder Architecture [20]

## 3.4 Model Hyper-Parameter Choices

HuggingFace provides a strong API with a large array of customization choices [18] for users to control. Part of these customization choices are the numerous parameters that can be tuned to impact a model's generation. This section will discuss the main parameters that are utilized to tweak generated outputs, as well as discuss parameters that can reduce repetition or increase entity inclusion rates. It should be noted that though there are many parameters to adjust, they often only have a small impact on the generated output.

### **Common Adjustments:**

- **Minimum & Maximum Length:** Sets the minimum and maximum lengths that the generated output can be. While the maximum length is a hard cut-off, the minimum length is not always reached if there are no more suggested tokens for generation. Furthermore, setting a minimum length can also force the model to make longer generations, which is ideal for our goals.
- **Temperature:** A numerical input that increases or decreases the confidence a model has in what the most likely response is. A higher value makes the model less confident. For example, a model that has a low temperature examining the following sentence "The mouse ate some \_\_\_\_" might consider "Cheese" to be the correct word with 95% confidence and 5% confidence for "Pizza". If the temperature were set higher, it might begin to equalize the confidence wherein "Cheese" would drop to 75% and "Pizza" rise to 25%. As temperature rises higher, the rates would equalize further.

- **Greedy Searching:** This method selects the word with the highest probability, conditional on all prior words, as the next entry. This method is quite simple but can lead to poor and unvaried results as words with very high conditional probabilities can be skipped over if they're masked by earlier words with low probabilities. In Figure 3.4, the starting word prior to generation is "The". As we follow along the lines or beams, we find that the next words could be "dog", "nice", or "car" with 40%, 50%, and 10% probability respectively. And so when utilizing a greedy search, the highest probability is selected and the beams not selected are discarded.

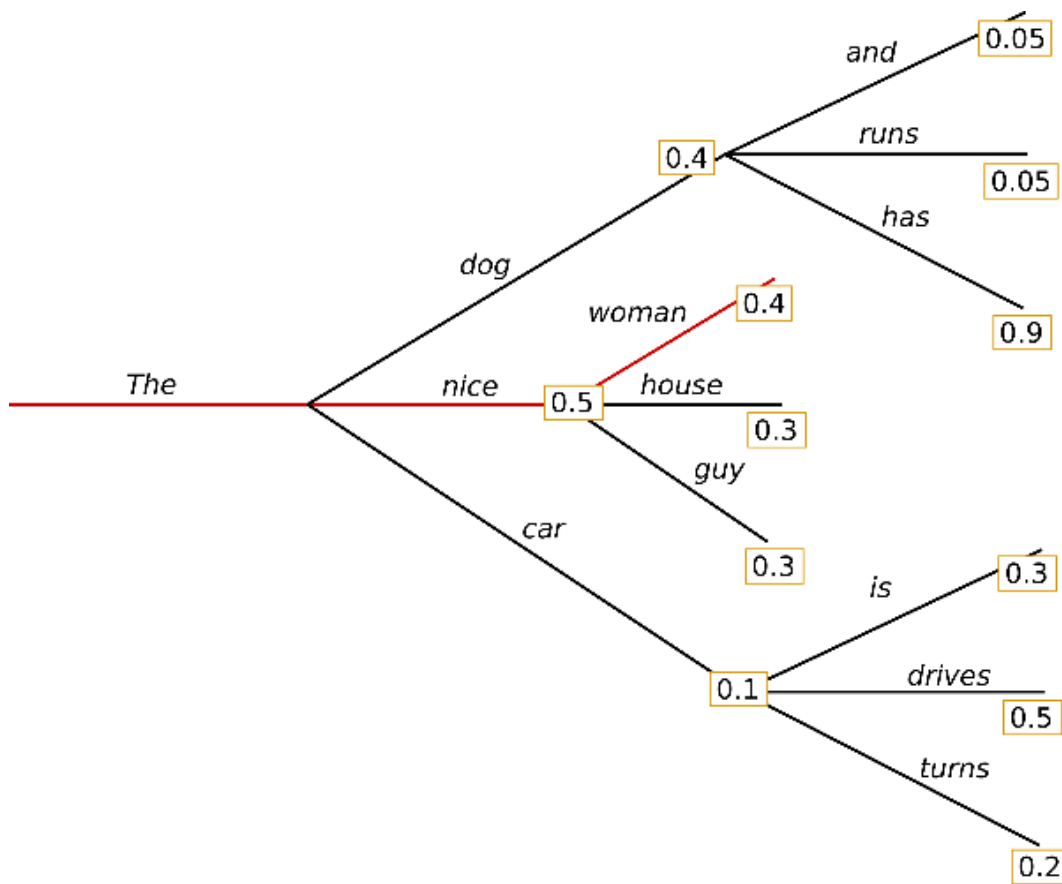


Figure 3.4: Greedy Search Generation [18]

- Beam Searching:** Beam searching provides a solution to the masking that occurs when utilizing a greedy search. Though it is more costly, this method keeps track of the probabilities for  $n$ -beams and then selects the words giving the highest probability. Figure 3.5 is a simple repeat of the one from the greedy searching section, but now has an additional line (going upwards instead of straight ahead after "The"). If the number of beams in a beam search were set to 2, the generative model would explore both the highest and second highest probability beams, in this case "Nice" and "Dog". At the next time step, it would then find the sequence "The dog has" to have an overall higher probability and thus would select this beam instead of the original, "The nice woman".

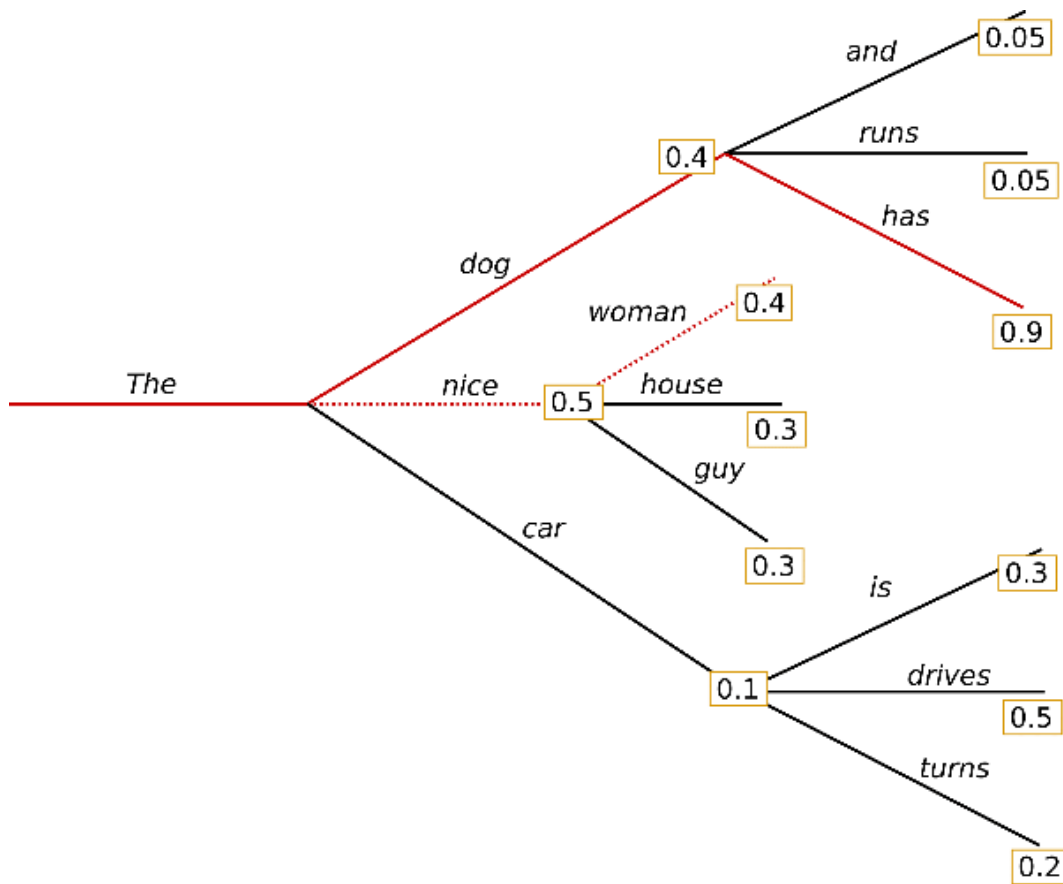
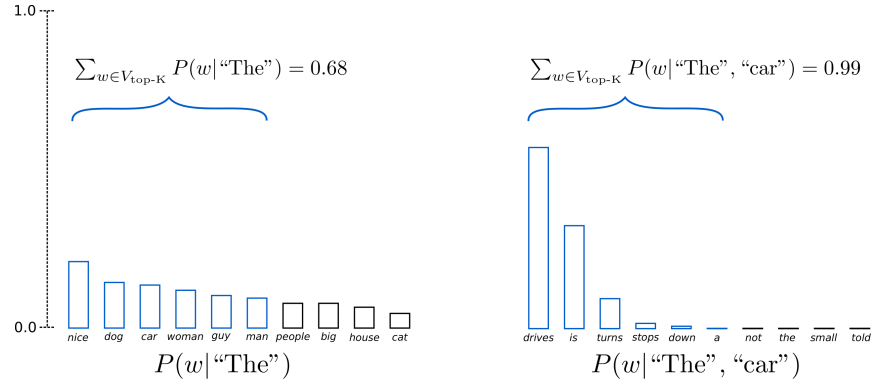
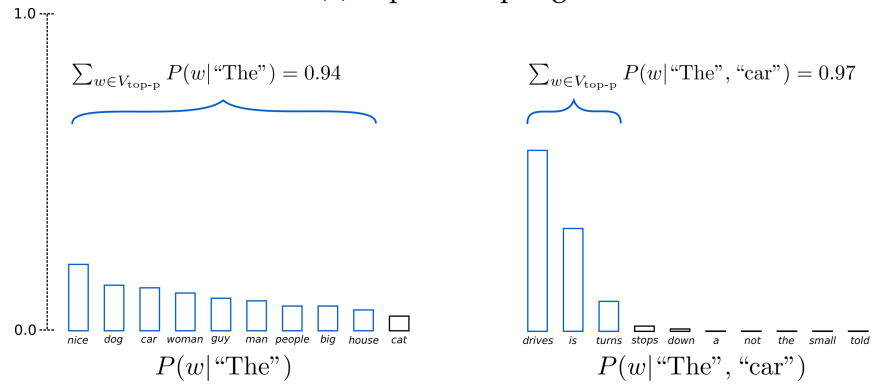


Figure 3.5: Beam Search Generation [18]

- **Sampling:** This method randomly selects the next word based on the conditional probability distribution of the prior words.
  - **Top-K Sampling:** By utilizing Top-K sampling, the  $k$  most likely words are selected, while the rest of the potential words are removed, thus changing the probability distribution of the remaining  $k$  words.
  - **Top-P (Nucleus) Sampling:** In a very similar process to Top-K sampling, Top-P or "Nucleus" sampling selects the smallest group of words whose total probability exceeds some rate  $p$ . Thus it provides a more dynamic approach as compared to Top-K.



(a) Top-K Sampling



(b) Top-P Sampling

Figure 3.6: Sampling Methods [18]

### **Repetition Solutions:**

- **Repetition Penalty:** One common issue with generated outputs is that tokens can tend to repeat themselves. As a quick overview, setting a penalty value reduces the value given to tokens that have already appeared. There is also a setting that can prevent  $n$ -gram repetition by preventing any outputs that would have an  $n$ -gram appear twice.

## **3.5 Known Abstractive Summarization Issues**

This section will provide some background on known issues with abstractive summarization generation, with examples from the dataset used presented. There are two issues that are well-known and well-researched in the NLP field, and various solutions have been proposed to mitigate them. One of the common issues in text generation is repetitive outputs, while the other is that generated text can create new entities that were not present in the original text and can even go on tangents related to these entities. Due to the use of pre-trained models for this project which prevents the altering of the underlying Seq2Seq encoder-decoder structure, some of the researched solutions could not be implemented, but are still discussed.

As an example used in this research, Table 3.1 displays one of the original articles from the WikiGold dataset [19] and two of the subsequent summarizations provided by the selected summarization model. Though the cases with repetition and entity hallucination are not necessarily frequent, different paragraph structures and hyper-parameter configurations can influence rates of occurrence.



Type	Output
Original	010 is the tenth album from Japanese Punk Techno band The Mad Capsule Markets . This album proved to be more commercial and more techno-based than Osc-Dis , with heavily synthesized songs like Introduction 010 and Come . Founding member Kojima Minoru played guitar on Good Day , and Wardanceis cover of a song by UK post punk industrial band Killing Joke . XXX can of This had a different meaning , and most people did n't understand what the song was about . it was later explained that the song was about Cannabis ( ' can of this ' sounding like Cannabis when said faster ) it is uncertain if they were told to change the lyric like they did on P.O.P and HUMANITY . UK Edition came with the OSC-DIS video , and most of the tracks were re-engineered .
Repetition Loop	Japanese Punk Techno band The Mad Capsule Markets' second album Osc-Dis was met with a lot of controversy after the band changed the lyrics of their song XXX can of This to make it sound more like XXX can of This, make it sound more like XXX can of This, make it sound more like XXX can of This, make it sound more like XXX can of This
Entity Hallucination	Japanese Punk Techno band The Mad Capsule Markets have released their second album.010 in the UK, following the success of their first album Osc-Dis, which was released in the UK and Ireland in May of this year and became one of the best-selling albums of all time in Ireland.

Table 3.1: Example of Summary Degeneration

In the paper "The Curious Case of Neural Text *De*Generation" [14], the authors explored hyper-parameter tuning options and their impact on both repetition and entity hallucination rates. While this paper was for sequential text generation as opposed to abstractive summarization, it still provides valuable insights.

### **Repetition Loops**

In a further expansion on the paper by Holtzman et al. exploring text degeneration, the likelihood objective function was examined to determine if it had an impact on degeneration [15]. After examination, it was found that the likelihood function was giving too much weight to generated outputs that contained repetition. To solve this issue, researchers adjusted the likelihood objective function to penalize unlikely tokens further than it otherwise might have at both a sequence and token level.

Furthermore, though the above techniques were tested on sequential text generation, this method of "unlikelihood" training was also extended for abstractive summarization [16], with some modifications made. As a baseline without any tweaks, the unlikelihood method boosted ROUGE scores and reduced repetition. The key modification applied to further boost scores and reduce repetition used a variation on a coverage mechanism that penalized the model if the decoder's cross attention mechanism examines the same token multiple times. This variation acted on a token level, wherein the model is penalized for giving high probabilities to the same tokens multiple times.

### **Entity Hallucination/Factual Consistency**

One such approach to handling entity hallucination hypothesized that the issue of entity hallucination is embedded within the training dataset. The researchers [17] proposed three new metrics to quantify if entities are consistent across both source and generated texts.

- Precision-Source or  $\text{Prec}(S) = N(G \cap S) / N(G)$   
Where  $N$  refers to the entity set for the source document ( $S$ ) and generated summary ( $G$ ). This metric examines how many entities from the generated summary are a part of the source document.

- Precision-Target or  $\text{Prec}(T) = N(G \cap T) / N(G)$   
Where T refers to the human-made summary. This metric examines how many entities from the generated summary are a part of the human-made summary.
- Recall-Target or  $\text{Recall}(T) = N(G \cap T) / N(T)$   
This metric examines how many entries from the human-made summary are *not* present in the generated summary.

To mitigate entity hallucination, these metrics were utilized to subsequently filter out the dataset based on specific precision-recall thresholds. Specifically, if entities within the generated summary were not present in the source doc, the sentence in the summary is removed. If the summary is only one sentence, then the source document-human summary pair is entirely removed from the training data.

Along with this data filtration method was a proposed change to encoder structure, wherein the encoder is trained to classify summary-"worthy" entities that are in the source document and human summary. The filtering of data combined with summary-"worthy" entity classification resulted in significant improvements in curbing entity hallucination.

As the goal of this research is creating augmented data for named entity recognition, hallucination in summary outputs is un-ideal unless a method can be found to properly tag these newly added entities. Though this method was successful in reducing hallucination, it requires a training set with both source document and human-made summary. Unfortunately, no dataset exists that has (a) entity tags (b) source articles, and (c) human-made summaries of the articles. As a result, this method could not be employed. In an attempt to still mitigate the issue of entity hallucination, ROUGE score filtering is utilized and discussed in more detail in the Methodology chapter.

## Chapter 4: Methodology

The methodology section aims to highlight the approach that will be taken stemming from the aforementioned methods and information discussed in the literature review and background sections. Some of these reviewed methods will serve as a baseline for comparison to the new approach examined. While rules-based data augmentation approaches have been fairly well studied, even within the NER field where the task is somewhat more challenging, generative data augmentation approaches have significantly less exposure and testing.

### 4.1 Proposed Summary Generation Approach

There are two types of summarization methods, these being extractive and abstractive. Extractive summarization is a method wherein each sentence within the overall text is evaluated for importance. These sentences are then compared amongst one another, at which point only the subset of sentences classified as important are used to construct a summary. Meanwhile, abstractive summarization evaluates the text and then generates a brand new summary, both taking portions from sentences as well as reorganizing, changing phrases, and adding/removing words.

Since the goal of DA is to artificially construct *new* sentences, extractive summarization provides no benefit. Any summary constructed would be a combination of sentences taken verbatim from the original data. Thus, we consider abstractive summarization for this project.

## 4.2 Overview of Process Pipeline

This section will quickly outline the steps to the process of generating sufficient abstractive summaries for data augmentation and subsequent testing. Section 4.3 through to Chapter 5 will go into more detail.

1. Pre-Trained Model
  - (a) Abstractive Summary Models
  - (b) Named Entity Recognition Models
2. Source Data Adjustments
  - (a) Tag Linearization
  - (b) Tag Replacement Variants
  - (c) Named Entity Weighted Order
  - (d) Article Stemming
  - (e) Shuffling
  - (f) ROUGE Scoring
3. Train-Test Splitting
4. Model Training
5. Model Evaluation and Comparison

## 4.3 Models for Abstractive Summarization

For summarization, there are currently 4 well-known options that can provide abstractive outputs (as opposed to extractive). These are:

- BART [7]
- T5 [9]
- GPT-2 [10]
- PEGASUS (X-Sum Variant) [11]

Of the above models, BART and T5 are Seq-2-Seq structured transformers, while GPT-2 is a decoder and PEGASUS has an underlying T5 model but was trained on data specifically for abstractive summarization. BART, T5, and GPT-2 can be utilized for text generation, translation, Q&A, and abstractive summarization, among other things, while PEGASUS is geared more towards only abstractive summarization. The following sections briefly highlight the unique differences in models.

## **BART**

The BART transformer structure was created by Facebook. The structure is an expansion on BERT, which at its core is a bi-directional encoder which was trained via masked language modelling. BART furthers this by adding in an autoregressive decoder to provide further functionality over BERT and access to standard encoder-decoder tasks. Training of the model first corrupted source text via noising function and then learned to reconstruct the corrupted text as close to source-level as possible.

## **T5**

The T5 transformer structure is very similar to BART and was created by Google and released within the same week that BART was. It was also an expansion from BERT and thus has the same bi-directional encoder as BART, and further adds its own auto-regressive decoder. The primary differences in T5 and BART are the structure of layers within the decoder. One other slight difference is that the training method used by T5 is no longer corrupt/fill-in-the-blank but instead has a mix of variations used to train for different tasks.

## **GPT-2**

GPT-2 was created by Open-AI and is a transformer structure comprised of only an autoregressive decoder. Training for GPT-2 occurred on web-texts, wherein it attempted to predict new tokens with only information of the prior ones given to it. The decoder method is similar to T5 and BART, but the actual structure of the model is different. Due to lack of support via HuggingFace for GPT-2 summarization, this model was not examined in the paper.

## PEGASUS

PEGASUS was created by Google and based off the original T5 Seq2Seq transformer structure. Its goal was to fix other models' weaknesses in abstractive summarization generation. PEGASUS pre-trains the Seq2Seq model on a large text corpus wherein important sentences get removed/masked from the input source and then generated utilizing the remaining sentences. As a result, this pre-trained model is quite strong at summarization specifically, but cannot be utilized as well for other Seq2Seq tasks.

### Final Decision

Based on the training information for the above models, PEGASUS seems an ideal choice. To quantitatively verify this, a small sample of 25 articles from the WikiGold dataset were passed into all three models (BART, T5, PEGASUS) simultaneously, with each model using the same configuration of hyper-parameters. The output of each model was then examined by human eye for sentence structure and fidelity. The ideal (or "sufficient", these two words may be used interchangeably) output was one that was coherent, stayed mostly faithful to the original article, and also provided a new structure that was not simply re-utilizing a sentence directly from the article. Sufficient length was also a factor. Multiple hyper-parameter combinations were examined, after which point it became clear that the PEGASUS X-Sum model (herein referred to simply as PEGASUS) was able to provide the best results consistently (as seen in Table 4.1. The other models tended to do one of three things most frequently:

1. Provide frequent incoherent, or repetitive outputs
2. Provide extractive (or near-extractive) outputs as opposed to abstractive ones
3. Provide excessively short outputs (~3-7 words maximum, even for articles that were 300+ words).

On the other hand, the PEGASUS model tended to provide outputs that were often coherent, of significant length (long or multiple sentences), and sufficiently unique from the sentences within articles.

As a result, the PEGASUS model was selected for this research. After this selection, other variants of the PEGASUS model were also examined (such as those tuned specifically on CNN Daily Mail and WikiHow datasets), but the X-Sum variant proved most consistent. Table 4.1 displays parameter settings by model and generated output issue rates on 25 articles.

Parameters	Model	Satisfactory	Extractive	Repetitive	Incoherent	Hallucination
Default	BART	0/25	19/25	0/25	6/25	0/25
	T5	1/25	19/25	0/25	5/25	0/25
	PEGASUS	17/25	7/25	0/25	1/25	0/25
MinLength=Article	BART	0/25	22/25	0/25	3/25	0/25
	T5	1/25	22/25	0/25	2/25	0/25
	PEGASUS	18/25	0/25	5/25	0/25	2/25
MinLength=Article Top P=0.9	BART	0/25	22/25	0/25	3/25	0/25
	T5	0/25	24/25	0/25	1/25	0/25
	PEGASUS	11/25	0/25	6/25	0/25	8/25
MinLength=Article Num Beams=32	BART	0/25	22/25	0/25	3/25	0/25
	T5	0/25	23/25	0/25	2/25	0/25
	PEGASUS	17/25	0/25	6/25	0/25	2/25
MinLength=Article Num Beams=32 RepetitionPenalty=2.0	BART	0/25	22/25	0/25	3/25	0/25
	T5	0/25	23/25	0/25	2/25	0/25
	PEGASUS	17/25	0/25	4/25	0/25	4/25

Table 4.1: Summary Output Success Rates

Clearly, PEGASUS heavily outperforms the T5 and BART models. As for the impact of parameters on the PEGASUS model, the minimum length being set to the article length seems to provide the largest boost to results. Further adding in a fixed number of 32 beams, we see a slight uptick in repetition. Adding in a repetition penalty ends up increasing hallucination. The parameters chosen need to strike a balance between providing sufficiently new sentences while mitigating repetition and especially hallucination. In the end, the parameters selected were:

- Minimum Length = Article Length
- Number of Beams = 32

This pairing did have more repetitions than when the number of beams was set to none, but these repetitions proved to be much smaller token-wise and only at the end of the sentence for a few words, as where in the case of no beams the repetition was much larger (e.g. the entire sentence



would be repeating several times versus a few words repeating at the end of the sentence twice).

## 4.4 Named Entity Recognition Models

Unlike abstractive summarization, named entity recognition models are consistently strong at performing the task required of them, with little training data required to fine-tune for decent results.

A common choice for NER tasks, the "BERT" [12] model was selected. This model was pre-trained on a large corpus of text data in the English language. The model was unsupervised in that no labels were provided by humans. The model training was performed utilizing a masked language model (MLM) wherein for each sentence 15% of the words within were masked or "hidden" from the model, at which point it had to estimate what the actual word was. Feeding the entire sentence at once allowed the model to learn bi-directionally. Similarly, the model was also tasked in predicting if two sentences followed one another or not.

Two key variations on the BERT model are the cased and un-cased versions. The difference is simple, in that the cased variant was trained on un-edited text, while the un-cased variant had all of the text converted to lower-case. Based on the structure of the WikiGold dataset and given that the entities are generally speaking all capitalized (when not numerical), the cased version was selected over the uncased one.

## 4.5 Data Source Adjustments

Data source adjustments refer to changes made to the source article text in an attempt to improve summarization output.

### Linearization

One common approach in NER tasks for preparing generational models is linearization, wherein words with entity tags have said tags added into the sentence itself in some way. Commonly, the tag is moved in front of or behind the word associated with it. Sometimes the tag goes both in front

and behind, or has "B"/"E" prefix and suffix to denote the beginning and end of an entity.

## Tag Replacement

In some cases, entities can be quite long and constructed of words that are commonly not entities. As a result, models can struggle to identify when certain non-entity words should be grouped together as entities instead. For example, "The 10th Battalion of Slayer's Creek Pontamac Group" is a very long entity that a model might struggle with understanding contextually within an article and thus might result in incoherent generations being output. To help the model understand sentence structure and context more easily, three tweaks were made to the original article structure, resulting in the four following potential input texts for every article. Replacement did not include "O" tags.

1. Original: Unaltered
2. Full Tag Replacement: Every token was replaced with their respective entity tags.
3. One Tag Replacement: Every entity was replaced with a single variant of their entity tag.
4. Unique Tag Replacement: Every entity was replaced with a single entity tag that had a unique number identifier appended to it.

See Table 4.2 as an example that compares the adjusted variants. Note that the unique tag replacement would mark subsequent unique PER/MISC tags as PER2, PER3, and so on.

Variation	Output
Original	Oleg Gazmanov is a Russian singer .
Reference Tags	PER PER            O O MISC   O    O
Full Tag Replacement	PER PER is a MISC singer .
One Tag Replacement	PER is a MISC singer .
Unique Tag Replacement	PER1 is a MISC1 singer .

Table 4.2: Example of Article Entity Replacements

## Generated Summary Entity Mapping

There are two primary ways of mapping tokens in the generated outputs.

1. **Sliding n-gram Approach:** The sliding n-gram approach is required for generated summaries based off of the original un-altered article. This approach maps all n-gram segments of entities from the source article. Then, the generated summary is initialized as all "outside" tags. Lastly, each n-gram within the summary is looped through and the tag overwritten if a mapping is found. Though this approach is not perfect, it is fairly accurate. Table 4.3 illustrates how the mapping function iterates on a sample text.

Iteration	Output					
Original	Oleg Gazmanov is a Russian singer .					
Reference Tags	PER	PER	O	O	MISC	O O
Initialization	O	O	O	O	O	O
1-Gram Iteration	O	O	O	O	MISC	O O
2-Gram Iteration	PER	PER	O	O	MISC	O O

Table 4.3: Example of Sliding n-gram Mapping

2. **Substitution:** As an attempt to improve entity mapping and prevent any errors, the tag replacement article variations were utilized. In the case of the One-Tag replacement scheme, it was sufficient to simply substitute one of the entities associated with the tag as this would not cause any incoherence within sentences. In the case of the Unique-Tag replacement, the exact entity was substituted back in.

## Named Entity Weighting

Weighted ordering is a method of shuffling the sentences within the original article. Sentences are given a weighting based on the proportion of entity words within the entire sentence. Once each sentence is weighted, the article has its sentences re-ordered in either ascending or descending order of entity weighting. This technique was devised from the fact that, generally speaking, the first few sentences of an article tend to be the most

important in summarizations done by humans. By proxy, one can assume that training data for summarizations would also follow this rule and as a result the PEGASUS model might inherently prioritize the first few sentences in its generation. In an examination of twenty-five articles, twenty-one of the article summaries (84%) referred to content from the first three sentences, often starting the summary the same as the first few words or phrases from the article.

### **Article Stemming**

Also relying on the entity weighting at a sentence level, this method simply removes any sentences that had a weighting below some defined threshold (e.g. 5%, 10%, etc.). While this approach might help increase the number of entities in the generated summary, it also has a large downside. Often, the generated outputs will combine two or three sentences into one fluent sentence that pulls details from all of the other sentences. By removing some sentences, we also remove filler information that the model could use to transition contextually from one entity to another. As a result, the summaries could become more extractive in nature, which is not necessarily ideal.

### **Shuffling**

One of the issues with abstractive summarization generation is that results are near identical for repeated runs on the same article. This provides a unique issue wherein it becomes challenging to get a sufficient number of generated samples from which to augment with. A simple approach of shuffling the sentences within the articles is taken. For each article, the sentences are split into a list. This list is then randomly shuffled, recorded, and recombined into a full article string. For subsequent shuffles, they are checked against the previously recorded shuffles to ensure that said shuffle is not mimicking the one already utilized. This list of newly shuffled articles is then fed into the model to generate summaries. Though shuffling does not always result in new summary generations, there are enough permutations existing that this is no longer an issue.

## ROUGE Scoring & Summary Filtering

The ROUGE metric stands for "Recall-Oriented Understudy for Gisting Evaluation". It provides a quantifiable output that scores the degree of similarity between one text and another. There are a few variations on ROUGE that adjust the level of granularity. For example, ROUGE-N measures how many N-grams overlap between two texts, while ROUGE-L compares via the longest sequence match and ROUGE-S compares via ordered pair similarities.

This project focuses on the ROUGE-N metric, for which recall refers to the percentage of n-grams from the reference text that are present in the generated text, and precision refers to the percentage of n-grams in the generated text are also present in the reference text. An F1 score also exists as a balancing metric.

In the context of this project, ROUGE scores provide a metric that can assess the adequacy of the generated summaries as compared to the original article. However, ROUGE scores do not *necessarily* provide insight into if a summary is fluent and void of repetition. As a result, there is some degree of post-processing that is required by a human. Given that the goal is to generate numerous summaries for data augmentation, human post-processing would require extensive examination of generated outputs, which is less than ideal.

To get around this issue, we can utilize the knowledge that *most* summaries generated using the PEGASUS model are sufficiently fluent and that *most* lack repetition. In previous human examination of a subset of article summaries (Table 4.1), less than 8% were incoherent and though 25% were repetitive, this only occurred at the very end of summaries. Furthermore, though higher ROUGE scores do not directly indicate the degree of fluency or repetition, higher scores tend to be more suitable.

As a result, we can filter out samples with lower scores (those with F1-scores below 20%) and simply take the top samples. For the number of samples selected, we simply take the article sentences and multiply this by 3 to generate 3 samples per sentence in the article. Though this approach does not guarantee that all the generated summaries utilized are sufficient, it provides more confidence than simply selecting randomly.

## 4.6 Model Evaluation & Comparison

The baseline approaches test other well-researched in an attempt to provide valid comparison for the new summarization approach. The two baseline approaches (both discussed prior in Section 2) being utilized are:

1. **Rule-Based Approach:** As a rule-based baseline, we examine three of the four data augmentation techniques; Label-Wise Token Replacement, Shuffle in Segments, and Synonym Replacement. Mention Replacement is not included due to the use of IO tag scheme over the required BIO scheme.
2. **Paraphrase Approach:** As a generational model baseline, we examine the method of paraphrasing on a sentence level to provide newly augmented sentences.

# Chapter 5: Experimental Results and Analysis

## 5.1 The WikiGold Dataset

### 5.1.1 Dataset Requirements & Selection

To carry out this research focused on utilizing abstractive summarization for Named Entity Recognition data augmentation, the dataset utilized had two key requirements.

1. The dataset needed to contain named entity tags for every token.
2. The dataset needed to contain sentences that were organized and coherent, with numerous relating to a single topic to form a paragraph/article. Individual unrelated sentences would not suffice for summarization.

The only dataset that adhered to the above two conditions was the "WikiGold" dataset [19]. The WikiGold dataset is a manually annotated corpus for named entity recognition and made up of a small sample of varied Wikipedia articles. Each word within the dataset (tokenized) was labelled with one of five tags which were taken from the CONLL-03 dataset. These tags were:

- O - The standard tag indicating a non-entity or "Outside" entity.
- LOC - An indication that the associated word is a Location.
- PER - An indication that the associated word is a Person.

- ORG - An indication that the associated word is an Organization.
- MISC - An indication that the associated word is an entity, but one that does not fall into the above categories.

Similarly to the CONLL03 dataset, a simple IO tagging format was utilized. Though this could easily be converted into an IOB system, for the purposes of abstractive summarization this change was not necessary and would hinder sentence structure for some techniques utilized.

### 5.1.2 Low-Resource Mimicry Selection

To mimic the varying degrees of data resources available and to properly examine the augmented data's impact on NER model performance, different groupings of 50, 100, 250, and 500 sentences are selected for testing purposes to perform augments upon.

To select these batches, the sentences within each article were counted. Next, articles with less than 2 and more than 40 sentences were filtered out to prevent both low sentence counts that cannot generate a sufficient number of shuffles and high sentences that would fill the entire batch size with just one or two articles and thus cause training data to be concentrated on one context/topic. With this filtering done, articles were then randomly selected until the total sentences combined from them were within 10% of the batch size (e.g. for batch size 50 between 45 and 55 sentences). Articles that were not selected for the training set were left for testing. This was then repeated 10 times with different random selections for each batch size to provide sufficient replication.

### 5.1.3 Descriptive Statistics

As a high-level overview of the WikiGold dataset, some statistics are provided to the reader. The dataset can be described with the following metrics:

- # of Articles: 145
- # of Sentences: 1,768
- # of Tokens: 39,007



- # of Non-Alphanumeric Tokens: 4,893
- # of Entities: 6,431
- Entity Categories: O, ORG, LOC, PER, MISC

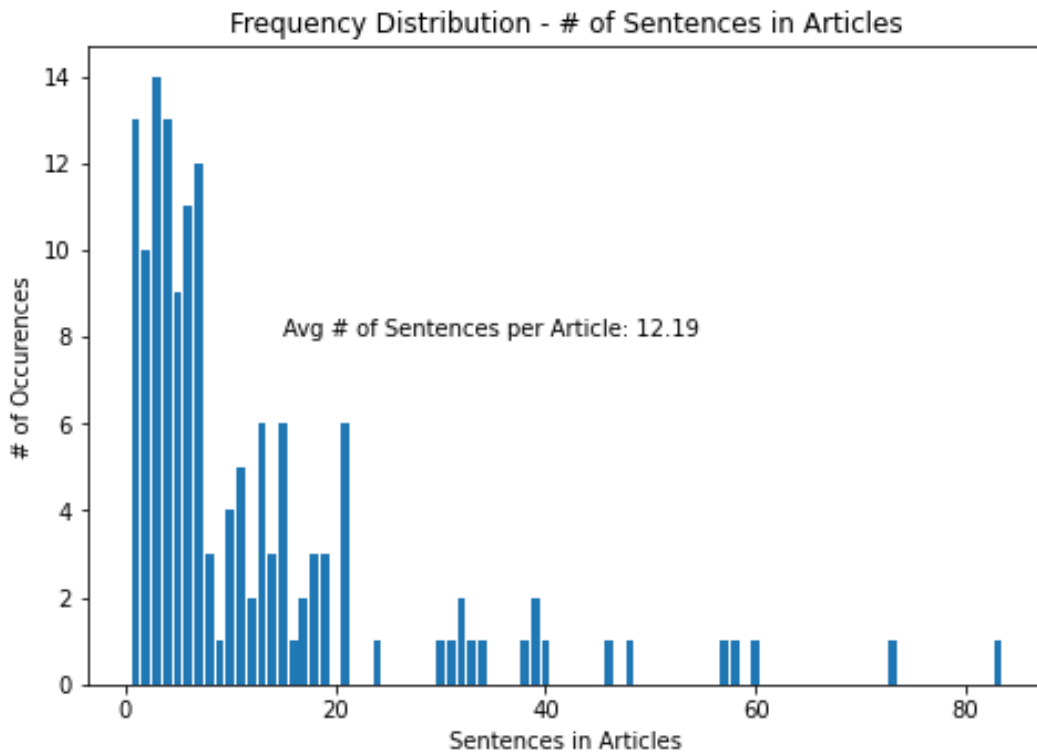


Figure 5.1: Distribution of Sentence Count per Article

Batch Size	X-Small		Small		Medium		Large	
Set	Train	Test	Train	Test	Train	Test	Train	Test
# of Articles	5	140	9	136	23	122	45	100
# of Sentences	49.6	1,718.4	100.6	1,667.4	246.3	1,521.7	505.8	1,262.2
# of Non-Entities	841.2	31,734.8	1,753.1	30,822.9	4,380.9	28,195.1	9,321.7	23,254.3
# of Entities	181.1	6,249.9	406.6	6,024.4	972.2	5,458.8	2,018.7	4,412.3

Table 5.1: Metrics by Training Test Split  
(Averaged over 10 Randomized Replications)

## 5.2 Data Augmentation Method Results

### Linearization

Attempts to linearize entity tags into the sentence prior to summarization resulted in a high degree of incoherency (wherein the model did not know how to interpret the linearized tags in a summary). As a result, linearization was unsuccessful.

### Named Entity Weighting & Stemming

While in theory weighting named entities to the front of the articles should produce a generation output with a higher number of entities due to how models are trained, this only resulted in one viable generated output. As for stemming, nearly all articles contained a sufficient number of entities within each sentence and as a result would have rarely removed sentences. Furthermore, the generated output from a shuffle post-stem often did not vary from the original un-stemmed version if the sentence removed were to be replaced back in the proper location. As a result, this approach was passed on in favour of the simpler shuffling mechanic.

## 5.3 NER Model Settings

Training settings for the BERT-Base Cased model iterated through 3 epochs with a learning rate of  $2e-5$  (0.00002) and a weight decay rate of 0.01. The model was evaluated with Precision, Recall, and F1 metrics.

## 5.4 Model Results

Table 5.2 below illustrates precision, recall, and F1 metrics on the testing set when the model was trained via their respective training sets. Training was randomized and was replicated 10 times on different train-test variations of the same batch-size to provide an accurate confidence interval on the accuracy metrics.

To make results clearer, the dataset with the highest F1 score for each batch size is bolded. The original dataset pertains to a batch of un-augmented data. In the case of the baseline datasets (Rules and Paraphrased), they contain the original un-augmented data and an additional three augmented sentences per original sentence. In the case of the summarization methods (Article, One Tag, and Uni Tag, as seen in Table 4.2), they contain the original un-augmented data and then take a number of generated summaries that is equal to three times the number of sentences in the original article to approximately match the three additions from the baseline (as each generated summary is often one sentence).

The metrics for the Rules dataset are for the highest-performing method/binomial rate combination (LWTR at a rate of 10%), though results across methods within rates were quite similar. As the rate increased, performance dropped slightly across all rule-based methods.

Overall, we see that the Rules based data augmentation method performed the best by a sizeable amount, with paraphrased and article augmentation methods falling slightly behind. One key difference is that the rules-based approach tends to have a significantly higher recall score at low batch sizes than the other two methods, which contributes to the higher F1 score overall as in comparison precision is quite similar albeit also somewhat higher.

Batch Size	Dataset	Precision	Recall	F1-Score
<b>X-Small</b> (S=50)	Original	3.33%±6.53%	0.00%±0.00%	0.01%±0.01%
	<b>Rules*</b>	<b>26.28%± 7.18%</b>	<b>16.48%± 4.79%</b>	<b>20.01%± 5.16%</b>
	Paraphrased	22.38%± 6.27%	8.37%± 3.08	11.68%± 3.56%
	Article	23.46%± 5.98%	8.69%± 3.65%	12.05%± 4.41%
	One Tag	15.97%± 9.63%	0.91%± 0.88%	1.64%± 1.56%
	Uni Tag	16.38%± 9.10%	1.83%± 1.47%	3.06%± 2.33%
<b>Small</b> (S=100)	Original	8.45%±5.56%	0.59%±0.77%	0.99%±1.23%
	<b>Rules*</b>	<b>53.77%± 3.1%</b>	<b>59.14%± 7.34%</b>	<b>56.18%± 4.84%</b>
	Paraphrased	47.27%± 2.58%	47.64%± 5.06	47.20%± 3.75%
	Article	43.49%± 3.83%	40.2%± 5.62%	41.61%± 4.71%
	One Tag	23.18%± 3.07%	12.56%± 3.37%	16.03%± 3.47%
	Uni Tag	25.32%± 2.35%	15.52%± 3.33%	18.91%± 3.12%
<b>Medium</b> (S=250)	Original	41.52%±2.97%	42.10%±3.90%	41.76%±3.32%
	<b>Rules*</b>	<b>70.07%± 1.7%</b>	<b>75.77%± 1.9%</b>	<b>72.8%± 1.54%</b>
	Paraphrased	67.34%± 1.03%	69.51%± 2.02	68.38%± 1.32%
	Article	66.50%± 1.20%	66.34%± 1.66%	66.41%± 1.34%
	One Tag	56.21%± 1.53%	47.97%± 2.70%	51.72%± 2.18%
	Uni Tag	56.70%± 2.34%	49.23%± 3.24%	52.68%± 2.84%
<b>Large</b> (S=500)	Original	65.93%±1.05%	70.49%±1.45%	68.11%±0.83%
	<b>Rules*</b>	<b>75.91%± 1.76%</b>	<b>82.19%± 1.21%</b>	<b>78.92%± 1.34%</b>
	Paraphrased	74.41%± 1.20%	79.14%± 0.85	76.70%± 0.99%
	Article	72.75%± 0.84%	75.92%± 0.47%	74.30%± 0.63%
	One Tag	68.58%± 1.02%	67.45%± 1.25%	68.00%± 0.96%
	Uni Tag	69.96%± 0.88%	69.91%± 1.40%	69.91%± 0.94%

Table 5.2: Model Results for Tested Methods ( $\alpha = 0.05$ ,  $n=10$ )

*\*Label-Wise Token Replacement, 10% Rate*

# Chapter 6: Future Work and Conclusions

## 6.1 Future Work

### Theoretical Optimal Tuning

With all the aforementioned uncertainty in summary generation, the process of fine-tuning hyper-parameters would be incredibly costly and time-consuming. The approach that could be taken is discussed in this section, but was not done in actuality.

There are three main hurdles that would require various configurations to be tested to determine the optimal hyper-parameters for ideal summary generation.

1. **Model Hyper-Parameters:** As discussed in the previous sections, there are a large number of model parameters that can be adjusted to tweak generated outputs. In total, there are seven parameters of interest that would need to be tested, these being: search method (greedy or beam), sampling type (top-k or top-p), temperature, repetition penalty, and minimum length.
2. **Article Variations:** Also discussed prior, the four variations on the original articles that intend to make generation easier for the model would all need to be evaluated.
3. **Shuffle Variations:** Lastly, different shuffling of article sentences would need to be tested for further consistency.

These variations would then have to be examined by human eye, though ROUGE filtering could help cut down manual labour required. Overall, combining the three core factors that would need to be tested rigorously results in a theoretical number of combinations too costly to examine. For instance, examining 30 shuffled variants across the 4 article variations takes upwards of 12 hours on one set of conservative model hyperparameters. As such, configuring all 7 and trying various combinations in a grid-search method would take many days.

### **Higher Augment Multiples**

While the paraphrasing and abstractive summarization results are fairly similar, there is a limit to how many paraphrased variations models can generate. While the sentence "His name was Bill Gates and he founded Microsoft" could viably be paraphrased in three other ways (e.g. "Named Bill Gates, he is the founder of Microsoft, etc."), results are often still quite similar and not varied. With this in mind, it would be worth exploring the impact of (a) a higher augment multiple (i.e. augment 7 new sentences instead of 3 from one original) and (b) how paraphrasing holds up against summarization in these scenarios.

## **6.2 Project Conclusion**

### **The Aim**

The goal of this project was to explore a new method of data augmentation specifically for the natural language processing task of tagging named entities. Current methods of data augmentation for NER tasks include simple algorithmic rules-based approaches, as well as more complicated generational approaches such as back-translation, paraphrasing, and sequential generation. As a new approach, abstractive summarization was explored and tested.

### **Methods Explored**

With the new approach via abstractive summarization, some unique challenges presented themselves. Firstly, generating a summary requires sen-

tences that are related to one another and pertain to the same context. Furthermore, a single sentence cluster (article) will only generate one to two sentences at most, meaning that an article that is 40 sentences in length will only have 2 new augmented sentences. However, human generated summaries tend to be constructed from the first few sentences of an article; as a result, models that are focused on generating summaries (PEGASUS) tend to have inherently learned a higher weighting for the initial sentences in longer articles. With this in mind, we can shuffle the sentences within an article to significantly alter the generated summary output. Through this technique, the article with 40 sentences can now be re-shuffled many times to provide as many augmented sentences as necessary.

Summary generation is not always consistent; in an attempt to increase consistency, methods of replacing entities with their associated tags (somewhat similar to sentence linearization) prior to generation were also explored.

## **Conclusion**

The results of summary generation with and without tag replacement were significantly lower for all batch sizes, with the tag replaced variants performing much worse on lower sample sizes. As for the original article summary method without tags replaced, it performed only slightly worse to the paraphrasing method performed in other research. Though it did outperform on the smallest batch size, this could simply be due to only 10 replications being performed. When compared to the rules-based approach baseline however, both the paraphrasing method and abstractive summarization techniques end up being outperformed by a fairly significant margin at the X-Small, Small, and Medium batch sizes. At the large batch size, all methods provide a slight boost but are much closer in performance.

# Bibliography

- [1] Xiang Dai and Heike Adel. (2020) *An Analysis of Simple Data Augmentation for Named Entity Recognition*. <https://aclanthology.org/2020.coling-main.343.pdf>
- [2] Bosheng Ding, Linlin Liu, Lidong Bing, Canasai Kruengkrai, Thien Hai Nguyen, Shafiq Joty, Luo Si, and Chunyan Miao. (2020) *DAGA: Data Augmentation with a Generation Approach for Low-resource Tagging Tasks*. <https://aclanthology.org/2020.emnlp-main.488.pdf>
- [3] Shuguang Chen, Gustavo Aguilar, Leonardo Neves, and Tamar Solorio. (2021) *Data Augmentation for Cross-Domain Named Entity Recognition*. <https://arxiv.org/pdf/2109.01758.pdf>
- [4] Usama Yaseen and Stefan Langer. (2021) *Data Augmentation for Low-Resource Named Entity Recognition Using Backtranslation*. <https://arxiv.org/pdf/2108.11703v1.pdf>
- [5] Petr Marek, Štěpán Müller, Jakub Konrád, Petr Lorenc, Jan Pichl, and Jan Šedivý. (2021) *Text Summarization of Czech News Articles Using Named Entities*. <https://arxiv.org/pdf/2104.10454.pdf>
- [6] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. (2019) *RoBERTa: A Robustly Optimized BERT Pretraining Approach*. <https://arxiv.org/abs/1907.11692>
- [7] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. (2019) *BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension*. <https://arxiv.org/abs/1910.13461>



- [8] Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. (2018) *Marian: Fast Neural Machine Translation in C++*. <https://arxiv.org/abs/1804.00344>
- [9] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. (2019) *Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer*. <https://arxiv.org/abs/1910.10683>
- [10] OpenAI (2019) GPT-2: 1.5B Release <https://openai.com/blog/tags/gpt-2/>
- [11] Jingqing Zhang, Yao Zhao, Mohammad Saleh and Peter J. Liu (2020) PEGASUS: Pre-training with Extracted Gap-sentences for Abstractive Summarization <https://arxiv.org/pdf/1912.08777.pdf>
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee and Kristina Toutanova (2019) BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding <https://arxiv.org/pdf/1810.04805.pdf>
- [13] Rui Wang and Ricardo Henao (2021) Unsupervised Paraphrasing Consistency Training for Low Resource Named Entity Recognition. <https://aclanthology.org/2021.emnlp-main.430.pdf>
- [14] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi (2020) THE CURIOUS CASE OF NEURAL TEXT DeGENERATION <https://arxiv.org/pdf/1904.09751.pdf>
- [15] Anonymous Authors (2020) NEURAL TEXT DEGENERATION WITH UNLIKELIHOOD TRAINING [https://openreview.net/attachment?id=SJeYe0NtvH&name=original\\_pdf](https://openreview.net/attachment?id=SJeYe0NtvH&name=original_pdf)
- [16] Pranav Ajit Nair and Anil Kumar Singh (2021) On Reducing Repetition in Abstractive Summarization <https://aclanthology.org/2021.ranlp-srw.18.pdf>
- [17] Feng Nan, Ramesh Nallapati, Zhiguo Wang, Cicero Nogueira dos Santos, Henghui Zhu, Dejiao Zhang, Kathleen McKeown and Bing

- Xiang (2021) *Entity-level Factual Consistency of Abstractive Text Summarization* <https://arxiv.org/pdf/2102.09130.pdf>
- [18] Patrick Von Platen (2020, March 18). *How to generate text: using different decoding methods for language generation with Transformers* HuggingFace. <https://huggingface.co/blog/how-to-generate>
- [19] Pritish Uplavikar (2016). *Named Entity Recognition, Wikigold CONLL* <https://github.com/pritishuplavikar/Named-Entity-Recognition/blob/master/wikigold.conll.txt>
- [20] HuggingFace. *How do Transformers Work?* <https://huggingface.co/course/chapter1/4?fw=pt>

# Chapter 8: Appendices

## 8.1 Code Files

**Note:** Files are not fully integrate-able in an automated manner. Functions were run in an IDE and thus utilized pre-declared variables that would otherwise be unknown due to not being fed as a function input (e.g. "df").

### 8.1.1 Main Code Files

#### Summarization/Paraphrase Model

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Thu Aug  4 17:46:46 2022
4
5 @author: Doug
6 """
7
8 #####
9 # LIBRARIES
10 #####
11
12 import pandas as pd
13 import csv
14 import numpy as np
15 import tensorflow as tf
16 import time
17
18 #For Copies
19 import copy
20
21 #For Synonyms
```

```

22 import requests
23 from bs4 import BeautifulSoup
24
25 #For Shuffle
26 import random
27
28 #Hugging Face
29 from transformers import PegasusTokenizer, PegasusForConditionalGeneration
    , T5Tokenizer, T5ForConditionalGeneration, MT5Tokenizer,
    MT5ForConditionalGeneration, AutoModel, AutoTokenizer
30
31 #OS
32 import os.path
33 from os import path as os_path
34
35 #ITERTOOLS
36 import itertools
37
38 #ROUGE METRIC
39 from rouge_score import rouge_scorer
40
41 #FACTORIAL
42 import math
43
44 #NLTK
45 import nltk
46 nltk.download("punkt")
47
48 #MEMORY CLEARING
49 from GPUtil import showUtilization as gpu_usage
50 import torch
51 from numba import cuda
52
53 #####
54 ##### CACHE #####
55 #####
56
57 #Change Cache
58 import os
59 os.environ['TRANSFORMERS_CACHE'] = 'H:/TempHF_Cache/cache/transformers/'
60 os.environ['HF_HOME'] = 'H:/TempHF_Cache/cache/'
61 os.environ['XDG_CACHE_HOME'] = 'H:/TempHF_Cache/cache/'
62
63 #####
64 # CODE

```

```

65 #####
66 #https://huggingface.co/spaces/Wootang01/Paraphraser_two/blob/main/app.py
67
68 #INITIAL SETUP
69 #Set Device
70 torch_device = "cuda" #If throwing CUDA error, restart Python
71
72 #Set Models
73 tokenizer1 = PegasusTokenizer.from_pretrained("google/pegasus-xsum")
74 model1 = PegasusForConditionalGeneration.from_pretrained("google/pegasus-
xsum").to(torch_device)
75
76 tokenizer2 = PegasusTokenizer.from_pretrained("tuner007/pegasus_paraphrase
")
77 model2 = PegasusForConditionalGeneration.from_pretrained("tuner007/
pegasus_paraphrase").to(torch_device)
78
79
80 #####
81 ##### DATA #####
82 #####
83 def gold_dataframe():
84     #Create List of Articles, Tokens
85     df = pd.read_csv("H:/My Files/School/Grad School WLU/MRP/Research/
Files/Data/wikigold.txt",
86                     sep=' ', header=None, doublequote = True, quotechar='
',
87                     skipinitialspace = False, quoting=csv.QUOTE_NONE)
88
89     df.columns = ["Token", "Entity"]
90
91     current_article, current_token, article_list, token_list = [], [], []
92     ,[]
93
94     for index, row in df.iterrows():
95         if df["Token"][index] != "-DOCSTART-":
96             current_article.append(df["Token"][index])
97             current_token.append(df["Entity"][index])
98         else:
99             article_list.append(current_article), token_list.append(
current_token)
100             current_article, current_token = [], []
101
102     for index in range(len(article_list)):
103         article_list[index] = " ".join(article_list[index])

```

```

103     token_list[index] = " ".join(token_list[index])
104
105     #Back to DF
106     temp_dict = {"Article":article_list,"Entity":token_list}
107     df = pd.DataFrame(temp_dict)
108     return df
109
110 df=gold_dataframe()
111
112 def tagged_dataframe():
113     df=gold_dataframe()
114     df["Tagged_All"]=df["Article"]
115     df["Tagged_One"]=df["Article"]
116     df["Tagged_Uni"]=df["Article"]
117
118     #List of Entities
119     Entity_List = ["ORG","LOC","PER","MISC"]
120
121     for index, row in df.iterrows():
122
123         list_of_words = df["Tagged_All"][index].split(" ")
124         list_of_tags = df["Entity"][index].split(" ")
125
126         ###TAGGING ALL
127         word_sentences = []
128         tag_sentences = []
129
130         word_segment = []
131         tag_segment = []
132
133         for i,word in enumerate(list_of_words):
134             if word != ".":
135                 word_segment.append(word)
136                 tag_segment.append(list_of_tags[i])
137             else:
138                 word_segment.append(word)
139                 word_sentences.append(word_segment)
140                 word_segment = []
141
142                 tag_segment.append(list_of_tags[i])
143                 tag_sentences.append(tag_segment)
144                 tag_segment = []
145
146         for i,sentence in enumerate(word_sentences):
147             for j,word in enumerate(sentence):

```

```

148         tag = tag_sentences[i][j]
149         if tag != "0":
150             word_sentences[i][j] = tag[2:]
151
152     for i,sentence in enumerate(word_sentences):
153         word_sentences[i] = " ".join(sentence)
154
155     df["Tagged_All"][index] = " ".join(word_sentences)
156
157
158
159
160
161     #TAGGING SEGMENTS AS ONE
162     previous = ''
163     segment_list = []
164     tag_list = []
165     temp1 = []
166     temp2 = []
167
168     for index2, tag in enumerate(list_of_tags):
169         if tag == previous or previous == '':
170             temp1.append(list_of_words[index2])
171             temp2.append(tag)
172         elif tag != previous:
173             segment_list.append(temp1.copy())
174             tag_list.append(temp2.copy())
175
176             temp1.clear()
177             temp2.clear()
178
179             temp1.append(list_of_words[index2])
180             temp2.append(tag)
181
182             previous = tag
183
184     segment_list.append(temp1)
185     tag_list.append(temp2)
186
187     for index3, group in enumerate(segment_list):
188         segment_list[index3] = " ".join(group)
189         tag_list[index3] = tag_list[index3][0]
190
191     #print(segment_list)
192     #print(tag_list)

```

```

193
194     for index4, thingy in enumerate(segment_list):
195         new_tag = tag_list[index4]
196         if new_tag != "0":
197             segment_list[index4] = new_tag[2:]
198
199
200     df["Tagged_One"][index] = " ".join(segment_list)
201
202
203
204
205
206
207
208     #TAGGING UNIQUE
209     Entity_List_New = Entity_List.copy()
210     list_of_words_tagged = df["Tagged_One"][index].split(" ")
211     #list_of_words = df["Article"][index].split(" ")
212     #If time, make numbering unique i.e. if band shows up twice give
same # for it
213
214     for i, word in enumerate(list_of_words_tagged):
215         if word in Entity_List:
216             Tag_Index = Entity_List.index(word)
217             Original_Tag = Entity_List_New[Tag_Index]
218
219             if Original_Tag[-1].isdigit():
220                 New_Tag = Original_Tag[:-1]+str(int(Original_Tag[-1])
+1)
221             else:
222                 New_Tag = Original_Tag+"1"
223
224             Entity_List_New[Tag_Index] = New_Tag
225             list_of_words_tagged[i] = New_Tag
226
227     df["Tagged_Uni"][index] = " ".join(list_of_words_tagged)
228
229     return df
230
231 df = tagged_dataframe()
232
233 #####
234 ##### SUMMARIZATION #####
235 #####

```



```

236
237 def abs_summary(
238     input_text, num_return_sequences, num_beams, min_length, temperature
    =1.5
239 ):
240     #PEGASUS XSUM
241
242     batch1 = tokenizer1(input_text, truncation=True, padding="longest",
    return_tensors="pt").to(torch_device)
243     translated1 = model1.generate(**batch1, temperature=temperature,
    min_length=min_length,
244                                     num_beams=num_beams,
    num_return_sequences=num_return_sequences)#,do_sample=False,top_k=None
    )
245     Pegasus = tokenizer1.batch_decode(translated1, skip_special_tokens=
    True)
246
247     return Pegasus
248
249 def article_iter(shuffled, num_samples, num_beams, temperature, df_name):
250     for i in range(num_samples):
251         header = "Sample" + str(i+1)
252         df[header] = " "
253
254     for index, row in df.iterrows():
255         print(index)
256
257         if shuffled==False:
258             article_txt = df["Article"][index]
259         else:
260             current = df['Article'][index].split(" .")
261             current = current[:-1] #removes final empty portion
262             random.shuffle(current)
263             article_txt = " .".join(current)
264
265     prop_length = article_txt.count(" ")#/6 #arbitrarily Picked
266
267     results = abs_summary(
268         input_text = article_txt,
269         num_return_sequences=num_samples,
270         num_beams=num_beams,
271         min_length=int(prop_length),
272         temperature=temperature
273     )
274

```

```

275         for i in range(num_samples):
276             header = "Sample" + str(i+1)
277             df[header][index] = results[i]
278
279 df.to_csv("H:/My Files/School/Grad School WLU/MRP/Research/Files/Data/
"+df_name+".csv")
280
281 return df
282
283 def weighted_entity_order(high_first):
284
285     for index, row in df.iterrows(): #Each Article
286         print(index)
287
288         tokens = df['Article'][index].split(" ")
289         entities = df['Entity'][index].split(" ")
290         sentences = df['Article'][index].split(" .")[:-1]
291
292         for i,sentence in enumerate(sentences):
293             sentences[i] = sentence+" ."
294
295         len_sentence = 0
296         num_entities = 0
297         weights = []
298
299         for i, token in enumerate(tokens):
300             if token == ".":
301                 wt_sentence = num_entities/len_sentence
302                 weights.append(wt_sentence)
303                 len_sentence = 0
304                 num_entities = 0
305             elif entities[i]!="0":
306                 len_sentence += 1
307                 num_entities += 1
308             else:
309                 len_sentence += 1
310
311         order = sorted(range(len(weights)), key=lambda k: weights[k],
reverse=high_first)
312         new_order = [""]*len(weights)
313
314         for i,val in enumerate(order):
315             new_order[i] = sentences[val]
316
317         new_order = "".join(new_order)

```

```

318         df["Article"][index] = new_order
319
320     return df
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340 def run_samples_overnight(method_choice):
341     tic = time.perf_counter()
342     df = tagged_dataframe()
343     shuffled = True
344     num_shuffle=50
345     method = method_choice
346
347     for i in range(num_shuffle):
348         header = "Sample" + str(i+1)
349         df[header] = " "
350
351     for index, row in df.iterrows():
352         i=0
353         cnt=0
354         timeout_cnt = 0
355
356         if shuffled==False:
357             article_text = df[method][index]
358
359             results = abs_summary(
360                 input_text = article_text,
361                 num_return_sequences=1,
362                 num_beams=32,

```

```

363         min_length=int(article_text.count(" "),
364         temperature=4
365         )
366
367     if shuffled==True:
368         while i < num_shuffle:
369             current = df[method][index].split(" .")
370             current = current[:-1] #removes final empty portion
371
372             total_permutations = math.factorial(len(current))/1
373             print("Article: "+str(index)+" and method: "+method)
374             print("Permutations: "+str(total_permutations))
375             print("Shuffle "+str(i+1))
376
377             random.shuffle(current)
378             article_text = " ".join(current)
379
380             if i==0:
381                 shuffle_list = []
382                 shuffle_list.append(article_text)
383                 i=i+1
384                 cnt+=1
385
386             elif article_text not in shuffle_list:
387                 shuffle_list.append(article_text)
388                 i=i+1
389                 cnt+=1
390
391             timeout_cnt +=1
392
393             print("Count: "+str(cnt)+"\n")
394
395             if cnt==total_permutations or timeout_cnt == num_shuffle
*3:
396                 print("Permutation Limit Reached\n")
397                 break
398
399         #Shuffling Complete
400         results=[]
401         print("Generating Results!")
402         for i in range(cnt):
403
404             summary_text = abs_summary(
405                 input_text = shuffle_list[i],
406                 num_return_sequences=1,

```

```

407         num_beams=32,
408         min_length=int(shuffle_list[i].count(" ")),
409         temperature=4
410     )
411
412     results.append(summary_text)
413
414     for i in range(cnt):
415         header = "Sample" + str(i+1)
416         df[header][index] = results[i][0]
417
418
419     #USE IF NUM SHUFFLED=0 NOT TRUE FALSE
420     df.to_csv("H:/My Files/School/Grad School WLU/MRP/Research/Files/Data/
dataframes/"+method+".csv", encoding="utf-8")
421     toc = time.perf_counter()
422
423     seconds_taken = toc-tic
424     minutes = seconds_taken/60
425     print("Seconds: %0.4f" % seconds_taken)
426     print("Minutes: %0.4f" % minutes)
427
428     return minutes
429
430
431     #FOR OVERNIGHT SAVING OF FILES
432     #9 Hours to Run with: 4 Variations, 50 Samples
433     minutes_to_run = []
434     for method_chosen in ["Article", "Tagged_All", "Tagged_Uni", "Tagged_One"]:
435         print("Starting Method: "+method_chosen)
436         mins = run_samples_overnight(method_choice=method_chosen)
437         minutes_to_run.append(mins)
438
439         gpu_usage()
440         torch.cuda.empty_cache()
441         torch_device="cuda"
442
443     print(sum(minutes_to_run))
444
445
446
447
448
449
450     #####

```

```

451 ##### PARAPHRASING #####
452 #####
453
454 def paraphraser(
455     input_text, num_return_sequences, num_beams, max_length=60,
456     temperature=1.5
457 ):
458     #PEGASUS XSUM
459
460     batch2 = tokenizer2(input_text, truncation=True, padding="longest",
461         return_tensors="pt").to(torch_device)
462     translated2 = model2.generate(**batch2, temperature=temperature,
463         max_length=max_length,
464         num_beams=num_beams,
465         num_return_sequences=num_return_sequences)#,do_sample=False,top_k=None
466     )
467     Pegasus = tokenizer2.batch_decode(translated2, skip_special_tokens=
468         True)
469
470     return Pegasus
471
472 def run_para_overnight(method_choice):
473     tic = time.perf_counter()
474     df = tagged_dataframe()
475     num_responses=5
476     method = method_choice
477
478     for i in range(num_responses):
479         header = "Sample" + str(i+1)
480         df[header] = " "
481
482     for index, row in df.iterrows():
483         print("Article #" + str(index))
484         new_article = []
485
486         current = df[method][index].split(" .")
487         current = current[:-1] #removes final empty portion
488
489         for entry in current: #Sentence in Article
490             results = paraphraser(
491                 input_text=entry,
492                 num_return_sequences=num_responses,
493                 num_beams=5,
494                 max_length=60,

```

```

490         temperature=1)
491
492         new_article.append(results)
493
494         np_array = np.array(new_article)
495         sentence_list = np_array.T.tolist()
496
497         for i in range(num_responses):
498             header = "Sample" + str(i+1)
499             df[header][index] = " ".join(sentence_list[i])
500
501
502         df.to_csv("H:/My Files/School/Grad School WLU/MRP/Research/Files/Data/
dataframes/Para_"+method+".csv", encoding="utf-8")
503         toc = time.perf_counter()
504
505         seconds_taken = toc-tic
506         minutes = seconds_taken/60
507         print("Seconds: %0.4f" % seconds_taken)
508         print("Minutes: %0.4f" % minutes)
509
510         return minutes, df
511
512 #FOR OVERNIGHT SAVING OF FILES
513 #9 Hours to Run with: 4 Variations, 50 Samples
514 minutes_to_run = []
515 for method_chosen in ["Article", "Tagged_All", "Tagged_Uni", "Tagged_One"]:
516     print("Starting Method: "+method_chosen)
517     mins, df = run_para_overnight(method_choice=method_chosen)
518     minutes_to_run.append(mins)
519
520     gpu_usage()
521     torch.cuda.empty_cache()
522     torch_device="cuda"
523
524 print(sum(minutes_to_run))

```

## NER Model

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Wed Aug 3 08:15:08 2022
4
5 @author: Doug
6 """
7
8 #####
9 ##### NER MODEL #####
10 #####
11
12 #Guide: https://huggingface.co/course/chapter7/2
13 #Fix: https://github.com/huggingface/datasets/issues/4099
14 #Fix 2: https://huggingface.co/datasets/nielsr/XFUN/commit/73
    ba5e026621e05fb756ae0f267eb49971f70ebd
15
16 #####
17 ##### CACHE #####
18 #####
19 import time
20 from datasets import load_dataset
21 from GPUUtil import showUtilization as gpu_usage
22 import torch
23 #df=tagged_dataframe()
24 #Change Cache
25 import os
26 os.environ['TRANSFORMERS_CACHE'] = 'H:/TempHF_Cache/cache/transformers/'
27 os.environ['HF_HOME'] = 'H:/TempHF_Cache/cache/'
28 os.environ['XDG_CACHE_HOME'] = 'H:/TempHF_Cache/cache/'
29
30 from transformers import *
31
32 #####
33 ##### CODE #####
34 #####
35 #For Variants
36 variants = ["Article","Tagged_One","Tagged_Uni"]
37 sizes = [50,100,250,500]
38 repetition = list(range(0,10))
39 #NEED TO DO ARTICLE--250--v9 LATER (not uploaded)
40
41 #For Rule-Based
42 rates = [0.3,0.5,0.7]
```



```

43 variants = ["LWTR","SR","SIS"]
44 sizes = [50,100,250,500]
45 repetition = list(range(0,10))
46
47 #For Paraphrased
48 #variants=["Paraphrased"]
49 #sizes = [50,100,250,500]
50 #repetition = list(range(0,10))
51
52 #Other
53 tic = time.perf_counter()
54 my_epochs = 3
55 time_dict = {}
56
57 for rate in rates: #Remove this and tab backwards everyhting else to do
    variants version (summary versions)
58     for variant in variants:
59         for size in sizes:
60             for j in repetition:
61
62                 #TEMP TESTING
63                 #variant="Paraphrased"
64                 #size=500
65                 #j=0
66
67                 #For VARIANTS Version (Summarization)
68                 #string_path = "H:\My Files\School\Grad School WLU\MRP\
Research\Files\Data\Textfiles\\"+variant+"\\"+str(size)+"\\"+variant+
str(size)+"v"+str(j)+"_wikigold_split.py"
69                 #new_model_path = "H:\\TempHF_Cache\\TrainingArgs\\"+
variant+"_"+str(size)+"v"+str(j)+"_NER_Model_"+str(my_epochs)+"
Epochs_UNAUGMENTED"
70
71                 #For RULES Version (Rules-Based)
72                 string_path = "H:\My Files\School\Grad School WLU\MRP\
Research\Files\Data\Textfiles\Rule_"+str(rate)+"\\"+variant+"\\"+str(
size)+"\\0"+str(rate)[2]+variant+str(size)+"v"+str(j)+"_wikigold_split
.py"
73                 new_model_path = "H:\\TempHF_Cache\\TrainingArgs\\0"+str(
rate)[2]+"_"+variant+"_"+str(size)+"v"+str(j)+"_NER_Model_"+str(
my_epochs)+"Epochs_AUGMENTED"
74
75                 #For Paraphrased Version
76                 #string_path = "H:\My Files\School\Grad School WLU\MRP\
Research\Files\Data\Textfiles\\"+variant+"\\"+str(size)+"\\"+variant+

```

```

77     str(size)+"v"+str(j)+"_wikigold_split.py"
78         #new_model_path = "H:\\TempHF_Cache\\TrainingArgs\\"+
79         variant+"_"+str(size)+"v"+str(j)+"_NER_Model_"+str(my_epochs)+"
80         Epochs_AUGMENTED"
81
82
83         #Continue..
84         cached_path = "H:\\TempHF_Cache\\cache\\datasets\\"
85         raw_datasets=load_dataset(string_path, cache_dir=
86         cached_path)
87
88         print(raw_datasets)
89
90         from transformers import AutoTokenizer
91
92         model_checkpoint = "bert-base-cased"
93         tokenizer = AutoTokenizer.from_pretrained(model_checkpoint
94         , cache_dir="H:\\TempHF_Cache\\Base_Tokenizer")
95
96         tokenizer.is_fast
97
98         ner_feature = raw_datasets["train"].features["ner_tags"]
99         ner_feature
100
101         label_names = ner_feature.feature.names
102         label_names
103
104         #PREPROCESS DATA
105
106         inputs = tokenizer(raw_datasets["train"][0]["tokens"],
107         is_split_into_words=True)
108         inputs.tokens()
109         inputs.word_ids()
110
111
112         def align_labels_with_tokens(labels, word_ids):
113             new_labels = []
114             current_word = None
115             for word_id in word_ids:
116                 if word_id != current_word:
117                     # Start of a new word!
118                     current_word = word_id
119                     label = -100 if word_id is None else labels[
120                     word_id]
121
122                     new_labels.append(label)
123             elif word_id is None:

```

```

115         # Special token
116         new_labels.append(-100)
117     else:
118         # Same word as previous token
119         label = labels[word_id]
120         # If the label is B-XXX we change it to I-XXX
121         if label % 2 == 1:
122             label += 1
123         new_labels.append(label)
124
125     return new_labels
126
127
128     labels = raw_datasets["train"][0]["ner_tags"]
129     word_ids = inputs.word_ids()
130     print(labels)
131     print(align_labels_with_tokens(labels, word_ids))
132
133
134     def tokenize_and_align_labels(examples):
135         tokenized_inputs = tokenizer(
136             examples["tokens"], truncation=True,
137             is_split_into_words=True
138         )
139         all_labels = examples["ner_tags"]
140         new_labels = []
141         for i, labels in enumerate(all_labels):
142             word_ids = tokenized_inputs.word_ids(i)
143             new_labels.append(align_labels_with_tokens(labels,
144                 word_ids))
145
146         tokenized_inputs["labels"] = new_labels
147         return tokenized_inputs
148
149     #Takes a bit...
150     tokenized_datasets = raw_datasets.map(
151         tokenize_and_align_labels,
152         batched=True,
153         remove_columns=raw_datasets["train"].column_names,
154     )
155
156     from transformers import
DataCollatorForTokenClassification

```

```

157         data_collator = DataCollatorForTokenClassification(
tokenizer=tokenizer)
158
159         batch = data_collator([tokenized_datasets["train"][i] for
i in range(2)])
160         batch["labels"]
161
162         for i in range(2):
163             print(tokenized_datasets["train"][i]["labels"])
164
165
166         #EVAL -- required pip install sequeval
167         import evaluate
168         metric = evaluate.load("sequeval")
169
170         labels = raw_datasets["train"][0]["ner_tags"]
171         labels = [label_names[i] for i in labels]
172         labels
173
174         predictions = labels.copy()
175         predictions[2] = "0"
176         metric.compute(predictions=predictions, references=[
labels])
177
178
179
180
181
182         #NOT SURE
183         import numpy as np
184
185         def compute_metrics(eval_preds):
186             logits, labels = eval_preds
187             predictions = np.argmax(logits, axis=-1)
188
189             # Remove ignored index (special tokens) and convert to
labels
190             true_labels = [[label_names[l] for l in label if l !=
-100] for label in labels]
191             true_predictions = [
192                 [label_names[p] for (p, l) in zip(prediction,
label) if l != -100]
193                 for prediction, label in zip(predictions, labels)
194             ]
195             all_metrics = metric.compute(predictions=

```

```

true_predictions, references=true_labels)
196         return {
197             "precision": all_metrics["overall_precision"],
198             "recall": all_metrics["overall_recall"],
199             "f1": all_metrics["overall_f1"],
200             "accuracy": all_metrics["overall_accuracy"],
201         }
202
203     #DEFINING THE MODEL
204     id2label = {str(i): label for i, label in enumerate(
label_names)}
205     label2id = {v: k for k, v in id2label.items()}
206
207
208     from transformers import AutoModelForTokenClassification
209
210     model = AutoModelForTokenClassification.from_pretrained(
211         model_checkpoint,
212         id2label=id2label,
213         label2id=label2id,
214         cache_dir="H:\TempHF_Cache\Base_Model"
215     )
216
217     #Check # of Labels is Correct:
218     model.config.num_labels
219
220     from transformers import TrainingArguments
221
222     args = TrainingArguments(
223         output_dir=new_model_path,
224         evaluation_strategy="epoch",
225         save_strategy="epoch",
226         learning_rate=2e-5,
227         num_train_epochs=my_epochs,
228         weight_decay=0.01,
229         push_to_hub=True,
230         hub_token = "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"
231     )
232
233
234
235     #TUNE
236     from transformers import Trainer
237
238     trainer = Trainer(

```

```

239         model=model,
240         args=args,
241         train_dataset=tokenized_datasets["train"],
242         eval_dataset=tokenized_datasets["test"],
243         data_collator=data_collator,
244         compute_metrics=compute_metrics,
245         tokenizer=tokenizer,
246     )
247
248     #IF ERROR WHEN PUSHING TO HUB USE !git lfs install
249     trainer.train()
250     #trainer.evaluate()
251     fin_results=trainer.evaluate()
252
253     #No Longer Uploading to Hub (Glitchy, Time Waste)
254     #trainer.push_to_hub()
255
256     #Save to TXT Files (FOR PARAPHRASED)
257     # output_save_path = "H:\My Files\School\Grad School WLU\
MRP\Research\Files\Models\Paraphrase_Results\\"
258     # with open(output_save_path+variant+str(size)+"v"+str(j)
+ "_Metrics.txt",'a',encoding='utf-8') as out_file:
259         #     out_file.write("Precision/Recall/F1\n")
260         #     out_file.write(str(fin_results["eval_precision"])+"\
t"+
261
262         #         str(fin_results["eval_recall"])+"\t"+
263         #         str(fin_results["eval_f1"]))
264
265     #Save to TXT Files (FOR RULES BASED)
266     output_save_path = "H:\My Files\School\Grad School WLU\MRP
\Research\Files\Models\Rule_Results\\"
267     with open(output_save_path+str(rate)[0]+str(rate)[2]+"\\"+
variant+"\\"+str(size)+"\\"+"v"+str(j)+"_Metrics.txt",'a',encoding='
utf-8') as out_file:
268         out_file.write("Precision/Recall/F1\n")
269         out_file.write(str(fin_results["eval_precision"])+"\t"
+
270
271         str(fin_results["eval_recall"])+"\t"+
272         str(fin_results["eval_f1"]))
273
274     #RESET
275     #gpu_usage()
276     torch.cuda.empty_cache()

```

```

277         torch_device="cuda"
278
279         toc = time.perf_counter()
280
281         seconds_taken = toc-tic
282         minutes = seconds_taken/60
283         print("Seconds: %0.4f" % seconds_taken)
284         print("Minutes: %0.4f" % minutes)
285         time_dict[variant+str(size)+"v"+str(j)]=minutes
286
287     #Test
288     #from transformers import pipeline
289     #classifier = pipeline("ner", model=model, tokenizer=tokenizer)
290     #classifier("My name is John Smith")

```

## Baseline Augmenting

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Wed Aug  3 16:08:28 2022
4
5 @author: Doug
6 """
7
8 import random
9 import numpy as np
10 import copy
11 import pandas as pd
12 import csv
13 import nltk
14 from nltk.corpus import wordnet
15
16 #####
17 ##### DF #####
18 #####
19
20 def gold_dataframe():
21     #Create List of Articles, Tokens
22     df = pd.read_csv("H:/My Files/School/Grad School WLU/MRP/Research/
23 Files/Data/wikigold.txt",
24                     sep=' ', header=None, doublequote = True, quotechar='
25 ',
26                     skipinitialspace = False, quoting=csv.QUOTE_NONE)
27     df.columns = ["Token","Entity"]
28
29     current_article, current_token, article_list, token_list = [], [], []
30     ,[]
31
32     for index, row in df.iterrows():
33         if df["Token"][index] != "-DOCSTART-":
34             current_article.append(df["Token"][index])
35             current_token.append(df["Entity"][index])
36         else:
37             article_list.append(current_article), token_list.append(
38 current_token)
39             current_article, current_token = [], []
40
41     for index in range(len(article_list)):
42         article_list[index] = " ".join(article_list[index])
43         token_list[index] = " ".join(token_list[index])
```



```

40
41     #Check it worked
42     for index in range(len(article_list)):
43         pass
44         #print(article_list[index],"\n",token_list[index])
45
46     #Back to DF
47     temp_dict = {"Article":article_list,"Entity":token_list}
48     df = pd.DataFrame(temp_dict)
49
50     return df
51
52 df = gold_dataframe()
53
54 #####
55 ##### ARTICLE VERS #####
56 #####
57 import ast
58 sizes = [50,100,250,500]
59
60 id_dict = {}
61 for size in sizes:
62     with open("H:\My Files\School\Grad School WLU\MRP\Research\Files\Data\
        Textfiles\Article\\"+str(size)+"\000_ID_LIST.txt") as indx_lst:
63         for line in indx_lst:
64             version = line[1]
65             ids = ast.literal_eval(line[4:])
66             id_dict[str(size)+"v"+version] = ids
67
68
69
70 #####
71 ##### RULE FXNS #####
72 #####
73 def LWTR_Method(augments):
74     DA_LWTR_additions = []
75
76     for k in range(0,augments):
77         DA_LWTR_word_sentences = copy.deepcopy(word_sentences)
78         for i,sentence in enumerate(word_sentences):
79             for j,word in enumerate(sentence):
80                 tag = tag_sentences[i][j]
81
82                 if np.random.binomial(1, rate, size=None):
83                     substitute = random.choice(tag_groups[tag])

```

```

84         DA_LWTR_word_sentences[i][j] = substitute
85
86         DA_LWTR_additions.append(DA_LWTR_word_sentences[i])
87
88         #print(word_sentences[0])
89         #print(DA_LWTR_word_sentences[0])
90         #print(DA_LWTR_additions[50])
91         #print(tag_sentences[0])
92         #print(len(DA_LWTR_word_sentences))
93         #print(len(DA_LWTR_additions))
94
95         return word_sentences+DA_LWTR_additions, tag_sentences+tag_sentences*
augments
96
97 def SR_Method(augments):
98     DA_SR_additions = []
99
100     SR_Issue_List = ["in", "In", "It", "it", "does", "Does", "IAEA", "have", "Have",
101                      ", "be", "Be", "less", "Less", "He", "he", "Pesos", "Inc", "inc", "acts", "Acts",
102                      "an", "An", "units"]
103
104     for k in range(0, augments):
105         DA_SR_word_sentences = copy.deepcopy(word_sentences)
106
107         for i, sentence in enumerate(word_sentences):
108             for j, word in enumerate(sentence):
109                 tag = tag_sentences[i][j]
110
111                 if np.random.binomial(1, rate, size=None):
112                     try:
113                         substitute = wordnet.synsets(word)[0].lemmas()[0].
name()
114
115                     if word not in SR_Issue_List:
116                         DA_SR_word_sentences[i][j] = substitute
117                     except:
118                         DA_SR_word_sentences[i][j] = word
119
120                 DA_SR_additions.append(DA_SR_word_sentences[i])
121
122         #print(word_sentences[0])
123         #print(DA_SR_word_sentences[0])
124         #print(tag_sentences[0])
125
126         return word_sentences+DA_SR_additions, tag_sentences+tag_sentences*
augments

```

```

124
125 def SIS_Method(augments):
126     DA_SIS_additions = []
127
128     for k in range(0,augments):
129         DA_SIS_word_sentences = copy.deepcopy(word_sentences)
130
131         for i,sentence in enumerate(word_sentences):
132
133             previous = ''
134             segment_list = []
135             temp1 = []
136
137
138             for j,word in enumerate(sentence):
139                 tag = tag_sentences[i][j]
140
141                 if tag == previous or previous == '':
142                     temp1.append(word)
143                 elif tag != previous:
144                     segment_list.append(temp1.copy())
145
146                     temp1.clear()
147
148                     temp1.append(word)
149
150                 previous = tag
151
152             segment_list.append(temp1)
153
154             #print(segment_list)
155
156             for j,entry in enumerate(segment_list):
157                 if np.random.binomial(1, rate, size=None):
158                     random.shuffle(entry)
159
160             DA_SIS_word_sentences[i] = np.hstack(segment_list).tolist()
161             DA_SIS_additions.append(DA_SIS_word_sentences[i])
162
163             #print(word_sentences[0])
164             #print(DA_SIS_word_sentences[0])
165             #print(tag_sentences[0])
166
167     return word_sentences+DA_SIS_additions, tag_sentences+tag_sentences*
augments

```

```

168
169
170
171 #####
172 ###SETUP###
173 #####
174 rates=[0.1,0.3,0.5,0.7]
175 for rate in rates:
176     print(rate)
177     for key in id_dict.keys():
178         print(key)
179         df=gold_dataframe()
180
181         breakdown=key.split("v")
182         size_choice=breakdown[0]
183         version=breakdown[1]
184
185         df = df.iloc[id_dict[key]]
186
187         word_sentences = []
188         tag_sentences = []
189         for index, row in df.iterrows():
190
191             words = df["Article"][index].split(" ")
192             tags = df["Entity"][index].split(" ")
193
194             word_segment = []
195             tag_segment = []
196             for i,word in enumerate(words):
197                 if word != ".":
198                     word_segment.append(word)
199                     tag_segment.append(tags[i])
200                 else:
201                     word_segment.append(word)
202                     word_sentences.append(word_segment)
203                     word_segment = []
204
205                     tag_segment.append(tags[i])
206                     tag_sentences.append(tag_segment)
207                     tag_segment = []
208
209
210
211 word_sentences_flat = np.hstack(word_sentences)
212 tag_sentences_flat = np.hstack(tag_sentences)

```

```

213
214     ### MAP WORDS TO TAGS
215     token_map = pd.DataFrame({"Words":word_sentences_flat, "Tags":
tag_sentences_flat})
216     tag_list = token_map["Tags"].unique().tolist()
217     tag_groups = token_map.groupby("Tags")["Words"].apply(list)
218
219     #####
220     ### RUN FUNCTIONS ###
221     #####
222     LWTR_Results, LWTR_Tags = LWTR_Method(3)
223     SR_Results, SR_Tags = SR_Method(3)
224     SIS_Results, SIS_Tags = SIS_Method(3)
225
226     #####
227     ### TO TEXT FILE ###
228     #####
229     LWTR_Path = "H:\\My Files\\School\\Grad School WLU\\MRP\\Research
\\Files\\Data\\Textfiles\\Rule_"+str(rate)+"\\LWTR\\"+str(size_choice)
+"\\v"+str(version)+"_Augmented.txt"
230     SR_Path = "H:\\My Files\\School\\Grad School WLU\\MRP\\Research\\
Files\\Data\\Textfiles\\Rule_"+str(rate)+"\\SR\\"+str(size_choice)+"\\
v"+str(version)+"_Augmented.txt"
231     SIS_Path = "H:\\My Files\\School\\Grad School WLU\\MRP\\Research\\
Files\\Data\\Textfiles\\Rule_"+str(rate)+"\\SIS\\"+str(size_choice)+"
\\v"+str(version)+"_Augmented.txt"
232
233     with open(LWTR_Path, 'w', encoding="utf-8") as LWTR_File:
234         for h,sentence in enumerate(LWTR_Results): #sentence from list
of sentences
235             for g,word in enumerate(sentence): #word from sentence
236                 LWTR_File.write(LWTR_Results[h][g]+" "+LWTR_Tags[h][g
]+ "\n")
237                 LWTR_File.write("\n")
238
239     #SR VARIANT
240     with open(SR_Path, 'w', encoding="utf-8") as SR_File:
241         for h,sentence in enumerate(SR_Results): #sentence from list
of sentences
242             for g,word in enumerate(sentence): #word from sentence
243                 SR_File.write(SR_Results[h][g]+" "+SR_Tags[h][g]+ "\n")
244                 SR_File.write("\n")
245
246     #SIS VARIANT
247     with open(SIS_Path, 'w', encoding="utf-8") as SIS_File:

```

```

248         for h,sentence in enumerate(SIS_Results): #sentence from list
of sentences
249             for g,word in enumerate(sentence): #word from sentence
250                 SIS_File.write(SIS_Results[h][g]+" "+SIS_Tags[h][g)+"\
n")
251                 SIS_File.write("\n")
252
253
254
255
256
257
258
259
260
261
262
263
264
265 #####
266 ###PARAPHRASE###
267 #####
268 for key in id_dict.keys():
269     df_og = pd.read_csv("H:/My Files/School/Grad School WLU/MRP/Research/
Files/Data/dataframes/Finished/Original_Paraphrase.csv", encoding="utf
-8", index_col=0)
270     df_mp = pd.read_csv("H:/My Files/School/Grad School WLU/MRP/Research/
Files/Data/dataframes/Finished/Mapped_Paraphrase.csv", encoding="utf-8
", index_col=0)
271
272     info = key.split("v")
273     size = info[0]
274     version = info[1]
275
276     df_og = df_og.iloc[id_dict[key]]
277     df_mp = df_mp.iloc[id_dict[key]]
278
279     sentence_list = []
280     tag_list = []
281
282     #Get Sentences and Tags in Simple Form
283     for index,row in df_og.iterrows():
284         for header in ["Article","Sample1","Sample2","Sample3"]:
285             if header=="Article":
286                 tokens = df_og[header][index].split(" ")

```

```

287         tags = df_mp["Entity"][index].split(" ")
288     else:
289         tokens = df_og[header][index].split(" ")
290         tags = df_mp[header][index].split(" ")
291
292     current_sentence=[]
293     current_tag=[]
294
295     for i,token in enumerate(tokens):
296         current_sentence.append(token)
297         current_tag.append(tags[i])
298
299         if token=="." or i==len(tokens)-1:
300             sentence_list.append(" ".join(current_sentence))
301             tag_list.append(" ".join(current_tag))
302             current_sentence,current_tag = [],[]
303
304     #print(sentence_list[0])
305
306     #Write to Txt File
307     txt_file_path = "H:\\My Files\\School\\Grad School WLU\\MRP\\Research
\\Files\\Data\\Textfiles\\Paraphrased\\"+str(size)+"\\v"+str(version)+
+"_Augmented.txt"
308     with open(txt_file_path,'w', encoding="utf-8") as aug_file:
309         for h,sentence in enumerate(sentence_list): #sentence from list of
sentences
310             for g,word in enumerate(sentence.split(" ")): #word from
sentence
311                 aug_file.write(word+" "+tag_list[h].split(" ")[g]+"\\n")
312                 aug_file.write("\\n")

```

## Mapping Functions

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Tue Aug 9 08:29:19 2022
4
5 @author: Doug
6 """
7
8 import pandas as pd
9 import nltk
10 nltk.download("punkt")
11 import csv
12 import random
13 import copy
14 from rouge_score import rouge_scorer
15 import numpy as np
16
17
18 def article_map():
19     # load dataframe
20     df = pd.read_csv("H:/My Files/School/Grad School WLU/MRP/Research/
21 Files/Data/dataframes/Article.csv", encoding="utf-8", index_col=0)
22     #df_taggedsamples = df.copy(deep=True)
23     # 1-gram words to be excluded (i.e. 'and' on its own should never be
24     # assigned a tag)
25     blacklist = ['can', 'of', 'A', 'a', 'for', 'and', 'nor', 'but', 'or', '
26 yet', 'so', 'both', 'and', 'whether', 'or', 'either', 'neither', 'just
27 ', 'the', 'The', 'as', 'if', 'then', 'rather', 'than', 'such', 'that']
28
29
30     # format samples
31     for idx, row in df.iterrows():
32         for column in df.iloc[:, 5:]:
33             empty_sent_list=[]
34             for t in nltk.sent_tokenize(df[column][idx]):
35                 words = nltk.word_tokenize(t)
36                 empty_sent_list.append(" ".join(words))
37
38             df[column][idx] = " ".join(empty_sent_list)
39
40     df_taggedsamples = df.copy(deep=True)
41
42     def getNgrams(article, tags, n):
```



```

40
41     if( n < 1 or n > len(tags) -1 ):
42         raise Exception("n must be between 1 and total number of tags")
43
44     if( len(article) != len(tags)):
45         raise Exception("article length and tag length do not match")
46
47     mapping = {}
48
49     # sliding window of size n
50     for i in range(len(tags) - n + 1):
51         # collect sliding window of tags
52         sequence = tags[i: i + n]
53         # if the tag is real (NOT "0") and all the tags match
54         if(sequence[0] != '0' and len(set(sequence)) <= 1):
55             # add the full sentence to the dictionary with a tag
56             mapping[" ".join(article[i:i+n])] = tags[i]
57
58     return mapping
59
60
61
62
63     # loop row by row
64     for idx, row in df.iterrows():
65         ngrams = 5
66         # loop for number of n-grams you want (currently 5)
67         fullMapping = {}
68         for i in range(1,ngrams):
69             mapping = getNgrams(df['Article'][idx].split(' '), df['Entity'][
idx].split(' '), i)
70             fullMapping = {**fullMapping, **mapping}
71
72         # remove blacklisted entries from dictionary
73         for word in blacklist:
74             try:
75                 del fullMapping[word]
76             except KeyError:
77                 pass
78
79         # apply the mapping to samples
80         for column in df.iloc[:, 5:]:
81             if df[column][idx]=='':
82                 continue
83             # get the sample and fill entity array with no-tag ('0')

```

```

84     sample = df[column][idx].split(' ')
85     entities = ['0'] * len(sample)
86     # iterate through each ngram length
87     for i in range(1, ngrams):
88         # sliding window loop for each ngram length of the full
sample
89         for j in range(len(sample) - i + 1):
90             sequence = sample[j: j + i]
91             # attempt to find an entity tag for the window
92             try:
93                 entity = fullMapping[' '.join(sequence)]
94                 # if no tag is found, do nothing
95             except KeyError:
96                 pass
97             # if a tag is found
98             #replace the list of entities at the indicies of the
sliding window
99             else:
100                 for k in range(i):
101                     entities[j+k] = entity
102             # in the copied dataframe, replace the sample with the entity
tags
103             df_taggedsamples[column][idx] = ' '.join(entities)
104
105     df.to_csv("H:/My Files/School/Grad School WLU/MRP/Research/Files/Data/
dataframes/Finished/Original_Article.csv", encoding="utf-8")
106     df_taggedsamples.to_csv("H:/My Files/School/Grad School WLU/MRP/
Research/Files/Data/dataframes/Finished/Mapped_Article.csv", encoding=
"utf-8")
107
108
109     return
110 article_map()
111
112
113
114 def paraphrase_map():
115     # load dataframe
116     df = pd.read_csv("H:/My Files/School/Grad School WLU/MRP/Research/
Files/Data/dataframes/Para_Article.csv", encoding="utf-8", index_col
=0)
117     df_taggedsamples = df.copy(deep=True)
118     # 1-gram words to be excluded (i.e. 'and' on its own should never be
assigned a tag)
119     blacklist = [',', 'can', 'of', 'A', 'a', 'for', 'and', 'nor', 'but', 'or'

```

```

, 'yet', 'so', 'both', 'and', 'whether', 'or', 'either', 'neither', '
just', 'the', 'The', 'as', 'if', 'then', 'rather', 'than', 'such', '
that']

120
121
122
123 # format samples
124 for idx, row in df.iterrows():
125     for column in df.iloc[:, 5:]:
126         empty_sent_list=[]
127         for t in nltk.sent_tokenize(df[column][idx]):
128             words = nltk.word_tokenize(t)
129             empty_sent_list.append(" ".join(words))
130
131         df[column][idx] = " ".join(empty_sent_list)
132
133 df_taggedsamples = df.copy(deep=True)
134
135 def getNgrams(article, tags, n):
136
137     if( n < 1 or n > len(tags) -1 ):
138         raise Exception("n must be between 1 and total number of tags")
139
140     if( len(article) != len(tags)):
141         raise Exception("article length and tag length do not match")
142
143     mapping = {}
144
145     # sliding window of size n
146     for i in range(len(tags) - n + 1):
147         # collect sliding window of tags
148         sequence = tags[i: i + n]
149         # if the tag is real (NOT "0") and all the tags match
150         if(sequence[0] != '0' and len(set(sequence)) <= 1):
151             # add the full sentence to the dictionary with a tag
152             mapping[" ".join(article[i:i+n])] = tags[i]
153
154     return mapping
155
156
157
158
159 # loop row by row
160 for idx, row in df.iterrows():
161     ngrams = 5

```

```

162     # loop for number of n-grams you want (currently 5)
163     fullMapping = {}
164     for i in range(1,ngrams):
165         mapping = getNgrams(df['Article'][idx].split(' '), df['Entity'][
idx].split(' '), i)
166         fullMapping = {**fullMapping, **mapping}
167
168     # remove blacklisted entries from dictionary
169     for word in blacklist:
170         try:
171             del fullMapping[word]
172         except KeyError:
173             pass
174
175     # apply the mapping to samples
176     for column in df.iloc[:, 5:]:
177         if df[column][idx]=='':
178             continue
179         # get the sample and fill entity array with no-tag ('0')
180         sample = df[column][idx].split(' ')
181         entities = ['0'] * len(sample)
182         # iterate through each ngram length
183         for i in range(1, ngrams):
184             # sliding window loop for each ngram length of the full
sample
185             for j in range(len(sample) - i + 1):
186                 sequence = sample[j: j + i]
187                 # attempt to find an entity tag for the window
188                 try:
189                     entity = fullMapping[' '.join(sequence)]
190                     # if no tag is found, do nothing
191                 except KeyError:
192                     pass
193                 # if a tag is found
194                 #replace the list of entities at the indices of the
sliding window
195                 else:
196                     for k in range(i):
197                         entities[j+k] = entity
198             # in the copied dataframe, replace the sample with the entity
tags
199             df_taggedsamples[column][idx] = ' '.join(entities)
200
201 df.to_csv("H:/My Files/School/Grad School WLU/MRP/Research/Files/Data/
dataframes/Finished/Original_Paraphrase.csv", encoding="utf-8")

```

```

202 df_taggedsamples.to_csv("H:/My Files/School/Grad School WLU/MRP/
    Research/Files/Data/dataframes/Finished/Mapped_Paraphrase.csv",
    encoding="utf-8")
203
204
205     return
206 paraphrase_map()
207
208
209
210 def tagged_one_map():
211     df = pd.read_csv("H:/My Files/School/Grad School WLU/MRP/Research/
    Files/Data/dataframes/Tagged_One.csv", encoding="utf-8", index_col=0)
212
213     headers = []
214     for i in range(1,51):
215         headers.append("Sample"+str(i))
216
217     #SETUP MAP
218     mapping={}
219     for index,row in df.iterrows():
220         list_of_words = df["Article"][index].split(" ")
221         list_of_tags = df["Entity"][index].split(" ")
222
223         previous = ''
224         segment_list = []
225         tag_list = []
226         temp1 = []
227         temp2 = []
228         article_map = {"ORG": [], "PER": [], "LOC": [], "MISC": []}
229
230         for index2, tag in enumerate(list_of_tags):
231             if tag == previous or previous == '':
232                 temp1.append(list_of_words[index2])
233                 temp2.append(tag)
234             elif tag != previous:
235                 segment_list.append(temp1.copy())
236                 tag_list.append(temp2.copy())
237
238                 temp1.clear()
239                 temp2.clear()
240
241                 temp1.append(list_of_words[index2])
242                 temp2.append(tag)
243

```

```

244         previous = tag
245
246     segment_list.append(temp1)
247     tag_list.append(temp2)
248
249     for index3, group in enumerate(segment_list):
250         segment_list[index3] = " ".join(group)
251         tag_list[index3] = tag_list[index3][0]
252
253     for index4, thingy in enumerate(segment_list):
254         new_tag = tag_list[index4]
255         if new_tag != "0":
256             #segment_list[index4] = new_tag[2:]
257             article_map[new_tag[2:]].append(thingy)
258
259     mapping[index] = article_map
260
261
262     #COPY MAP
263     #mapping_reset = copy.deepcopy(mapping)
264     #APPLY MAP
265     # format samples
266     for idx, row in df.iterrows():
267         for column in df.iloc[:, 5:]:
268             empty_sent_list=[]
269             if df[column][idx]==' ':
270                 continue
271
272             for t in nltk.sent_tokenize(df[column][idx]):
273                 words = nltk.word_tokenize(t)
274                 empty_sent_list.append(" ".join(words))
275
276             df[column][idx] = " ".join(empty_sent_list)
277
278
279
280     df_taggedsamples = df.copy(deep=True)
281
282     for index,row in df.iterrows():
283         for i,entry in enumerate(headers):
284             if df[entry][index]==' ':
285                 continue
286             #mapping = copy.deepcopy(mapping_reset)
287             list_of_words = df[entry][index].split(" ")
288             list_of_tags = df_taggedsamples[entry][index].split(" ")

```

```

289
290         if len(list_of_words)!=len(list_of_tags):
291             print("Initial Error")
292
293     #Go Through Each Word
294     for j,word in enumerate(list_of_words):
295         if word in ["ORG","PER","LOC","MISC"]:
296             replacement = random.choice(mapping[index][word])
297             #mapping[index][word].remove(replacement)
298             len_replace = len(replacement.split(" "))
299             list_of_words[j] = replacement
300
301             #print("Initial Tag: "+word)
302             #print("Replacement: "+replacement)
303
304             if len_replace == 1:
305                 list_of_tags[j] = "I-"+word
306                 #print("Tag Replace: "+list_of_tags[j])
307             else:
308                 list_of_tags[j] = str(str("I-"+word+" ")*
len_replace)[:len_replace-1]
309                 #print("Tag Replace: "+list_of_tags[j])
310             else:
311                 list_of_tags[j] = "0"
312
313         if len(list_of_words)!=len(list_of_tags):
314             print("New ERROR")
315
316         #print(list_of_words)
317         #print(list_of_tags)
318
319     df[entry][index] = " ".join(list_of_words)
320     df_taggedsamples[entry][index] = " ".join(list_of_tags)
321
322
323     if len(" ".join(list_of_words).split(" ")) != len(" ".join(
list_of_tags).split(" ")):
324         print("SECOND ERROR")
325
326     df.to_csv("H:/My Files/School/Grad School WLU/MRP/Research/Files/Data/
dataframes/Finished/Original_Tagged_One.csv", encoding="utf-8")
327     df_taggedsamples.to_csv("H:/My Files/School/Grad School WLU/MRP/
Research/Files/Data/dataframes/Finished/Mapped_Tagged_One.csv",
encoding="utf-8")
328

```

```

329     return
330 tagged_one_map()
331
332
333
334 def tagged_uni_map():
335     df = pd.read_csv("H:/My Files/School/Grad School WLU/MRP/Research/
Files/Data/dataframes/Tagged_Uni.csv", encoding="utf-8", index_col=0)
336
337     headers = []
338     for i in range(1,51):
339         headers.append("Sample"+str(i))
340
341     #SETUP MAP
342     mapping={}
343     for index,row in df.iterrows():
344         Entity_List = ["ORG","LOC","PER","MIS"]
345         Entity_List_New = ["ORG","LOC","PER","MIS"]
346         list_of_words = df["Article"][index].split(" ")
347         list_of_tags = df["Entity"][index].split(" ")
348
349         previous = ''
350         segment_list = []
351         tag_list = []
352         temp1 = []
353         temp2 = []
354         article_map = {}
355
356         for index2, tag in enumerate(list_of_tags):
357             if tag == previous or previous == '':
358                 temp1.append(list_of_words[index2])
359                 temp2.append(tag)
360             elif tag != previous:
361                 segment_list.append(temp1.copy())
362                 tag_list.append(temp2.copy())
363
364                 temp1.clear()
365                 temp2.clear()
366
367                 temp1.append(list_of_words[index2])
368                 temp2.append(tag)
369
370             previous = tag
371
372         segment_list.append(temp1)

```



```

373     tag_list.append(temp2)
374
375     for index3, group in enumerate(segment_list):
376         segment_list[index3] = " ".join(group)
377         tag_list[index3] = tag_list[index3][0]
378
379     #Part That Matters
380     for index4, thingy in enumerate(segment_list):
381         new_tag = tag_list[index4]
382         new_tag = new_tag[2:5]
383
384         if new_tag in Entity_List:
385             tag_index = Entity_List.index(new_tag)
386             original_tag = Entity_List_New[tag_index]
387
388             if original_tag[-1].isdigit():
389                 new_tag = original_tag[:-1]+str(int(original_tag[-1])
+1)
390
391                 article_map[new_tag] = thingy
392             else:
393                 new_tag = original_tag+"1"
394                 article_map[new_tag] = thingy
395
396             Entity_List_New[tag_index] = new_tag
397
398
399
400     mapping[index] = article_map
401
402
403     #COPY MAP
404     #mapping_reset = copy.deepcopy(mapping)
405     #APPLY MAP
406     # format samples
407     # format samples
408     for idx, row in df.iterrows():
409         for column in df.iloc[:, 5:]:
410             empty_sent_list=[]
411             if df[column][idx]==' ':
412                 continue
413             for t in nltk.sent_tokenize(df[column][idx]):
414                 words = nltk.word_tokenize(t)
415                 empty_sent_list.append(" ".join(words))
416

```

```

417         df[column][idx] = " ".join(empty_sent_list)
418
419
420     df_taggedsamples = df.copy(deep=True)
421     #df=df.head(1)
422     #headers=["Sample1"]
423
424
425     for index,row in df.iterrows():
426         for i,entry in enumerate(headers):
427             if df[entry][index]==" ":
428                 continue
429             list_of_words = df[entry][index].split(" ")
430             list_of_tags = df_taggedsamples[entry][index].split(" ")
431
432             if len(list_of_words)!=len(list_of_tags):
433                 print("Initial Error")
434
435             #Go Through Each Word
436             for j,word in enumerate(list_of_words):
437                 if word[:3] in ["ORG","PER","LOC","MIS"]:
438                     #CATCHING NICHE CASES WHERE GEN CUTS NUMBER OFF
439                     if word not in mapping[index]:
440                         word = word[:3]+"1"
441
442                     replacement = mapping[index][word]
443                     len_replace = len(replacement.split(" "))
444                     list_of_words[j] = replacement
445
446                     #print("Initial Tag: "+word)
447                     #print("Replacement: "+replacement)
448
449                     if len_replace == 1:
450                         list_of_tags[j] = "I-"+word[:3]
451                         #print("Tag Replace: "+list_of_tags[j])
452                     else:
453                         list_of_tags[j] = str(str("I-"+word[:3]+" ")*
len_replace)[-1]
454                         #print("Tag Replace: "+list_of_tags[j])
455                     else:
456                         list_of_tags[j] = "0"
457
458             if len(list_of_words)!=len(list_of_tags):
459                 print("New ERROR")
460

```

```

461         #print(list_of_words)
462         #print(list_of_tags)
463
464         df[entry][index] = " ".join(list_of_words)
465         df_taggedsamples[entry][index] = " ".join(list_of_tags)
466
467
468         if len(" ".join(list_of_words).split(" ")) != len(" ".join(
469             list_of_tags).split(" ")):
470             print("SECOND ERROR")
471
472         df.to_csv("H:/My Files/School/Grad School WLU/MRP/Research/Files/Data/
473             dataframes/Finished/Original_Tagged_Uni.csv", encoding="utf-8")
474         df_taggedsamples.to_csv("H:/My Files/School/Grad School WLU/MRP/
475             Research/Files/Data/dataframes/Finished/Mapped_Tagged_Uni.csv",
476             encoding="utf-8")
477
478     return
479 tagged_uni_map()
480
481
482
483
484
485
486 ### CONFIRM THAT ALL TOKENS HAVE A MATCHING TAG
487 uh_oh=0
488 headers = []
489 for i in range(1,51):
490     headers.append("Sample"+str(i))
491
492 for value in ["Article","Tagged_One","Tagged_Uni"]:
493     df_og = pd.read_csv("H:/My Files/School/Grad School WLU/MRP/Research/
494         Files/Data/dataframes/Finished/Original_"+value+".csv", encoding="utf
495         -8", index_col=0)
496     df_mp = pd.read_csv("H:/My Files/School/Grad School WLU/MRP/Research/
497         Files/Data/dataframes/Finished/Mapped_"+value+".csv", encoding="utf-8"
498         , index_col=0)
499
500     for index,row in df_og.iterrows():
501         for header in headers:

```

```

498         if df_og[header][index]==" " or pd.isna(df_og[header][index]):
499             continue
500         else:
501             og_len = len(df_og[header][index].split(" "))
502             mp_len = len(df_mp[header][index].split(" "))
503
504             if og_len!=mp_len:
505                 print("VARIANT: "+value+", ARTICLE: "+str(index)+", SAMPLE
: "+header)
506                 uh_oh = 1
507
508         if uh_oh!=1:
509             print("No Issues!")
510
511
512
513
514
515
516 #####
517 ##### SCORING SECTION #####
518 #####
519 def score_dfs():
520     variants = ["Article","Tagged_One","Tagged_Uni"]
521
522     for variant in variants:
523         print(variant)
524         df=pd.read_csv("H:/My Files/School/Grad School WLU/MRP/Research/
Files/Data/dataframes/Finished/Original_"+variant+".csv", encoding="
utf-8", index_col=0)
525
526         #Setup
527         df_scores = df.copy(deep=True)
528         f1_threshold = 0.2
529
530         scorer = rouge_scorer.RougeScorer(['rouge1', 'rouge2', 'rougeL'],
use_stemmer=True)
531
532         #Set Sample Headers to Iterate
533         headers = []
534         for i in range(1,51):
535             headers.append("Sample"+str(i))
536
537         #Iterate and Apply Scores
538         for index,row in df.iterrows():

```

```

539         original = df[variant][index]
540
541         for header_title in headers:
542             if pd.isna(df[header_title][index]) or df[header_title][
index]==" ":
543                 df[header_title][index]="NA"
544                 df_scores[header_title][index]="NA"
545                 continue
546
547             new = df[header_title][index]
548             scores = scorer.score(original, new)
549             f1_score = scores["rouge1"][2]
550
551             if f1_score > f1_threshold:
552                 df_scores[header_title][index] = f1_score
553             elif f1_score < f1_threshold and new != " ":
554                 df_scores[header_title][index] = "Low Score"
555             else:
556                 df_scores[header_title][index] = "NA"
557
558             #df.replace('N/A',pd.NA)
559             df_scores.to_csv("H:/My Files/School/Grad School WLU/MRP/Research/
Files/Data/dataframes/Finished/Scored_"+variant+".csv", encoding="utf
-8")
560
561         return
562 score_dfs()

```

## 8.1.2 Secondary Code Files

### Dataset Statistics

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Sat Aug 6 22:12:28 2022
4
5 @author: Doug
6 """
7
8 #####
9 # LIBRARIES
10 #####
11
12 import pandas as pd
13 import csv
14 import numpy as np
15 import tensorflow as tf
16 import time
17
18 #For Copies
19 import copy
20
21 #For Synonyms
22 import requests
23 from bs4 import BeautifulSoup
24
25 #For Shuffle
26 import random
27
28 #Hugging Face
29 from transformers import PegasusTokenizer, PegasusForConditionalGeneration
30 , T5Tokenizer, T5ForConditionalGeneration, MT5Tokenizer,
31 MT5ForConditionalGeneration, AutoModel, AutoTokenizer
32
33 #OS
34 import os.path
35 from os import path as os_path
36
37 #ITERTOOLS
38 import itertools
39
40 #ROUGE METRIC
41 from rouge_score import rouge_scorer
```

```

40
41 #FACTORIAL
42 import math
43
44 #NLTK
45 import nltk
46 nltk.download("punkt")
47
48 #MEMORY CLEARING
49 from GPUUtil import showUtilization as gpu_usage
50 import torch
51 from numba import cuda
52
53 #####
54 ##### CACHE #####
55 #####
56
57 #Change Cache
58 import os
59 os.environ['TRANSFORMERS_CACHE'] = 'H:/TempHF_Cache/cache/transformers/'
60 os.environ['HF_HOME'] = 'H:/TempHF_Cache/cache/'
61 os.environ['XDG_CACHE_HOME'] = 'H:/TempHF_Cache/cache/'
62
63 #####
64 # CODE
65 #####
66
67 def gold_dataframe():
68     #Create List of Articles, Tokens
69     df = pd.read_csv("H:/My Files/School/Grad School WLU/MRP/Research/
70                     Files/Data/wikigold.txt",
71                     sep=' ', header=None, doublequote = True, quotechar='
72                     ",
73                     skipinitialspace = False, quoting=csv.QUOTE_NONE)
74
75     df.columns = ["Token", "Entity"]
76
77     current_article, current_token, article_list, token_list = [], [], []
78
79     for index, row in df.iterrows():
80         if df["Token"][index] != "-DOCSTART-":
81             current_article.append(df["Token"][index])
82             current_token.append(df["Entity"][index])
83         else:

```

```

82     article_list.append(current_article), token_list.append(
current_token)
83     current_article, current_token = [], []
84
85     for index in range(len(article_list)):
86         article_list[index] = " ".join(article_list[index])
87         token_list[index] = " ".join(token_list[index])
88
89     #Back to DF
90     temp_dict = {"Article":article_list,"Entity":token_list}
91     df = pd.DataFrame(temp_dict)
92     return df
93
94 df=gold_dataframe()
95
96 def tagged_dataframe():
97     df=gold_dataframe()
98     df["Tagged_All"]=df["Article"]
99     df["Tagged_One"]=df["Article"]
100    df["Tagged_Uni"]=df["Article"]
101
102    #List of Entities
103    Entity_List = ["ORG","LOC","PER","MISC"]
104
105    for index, row in df.iterrows():
106
107        list_of_words = df["Tagged_All"][index].split(" ")
108        list_of_tags = df["Entity"][index].split(" ")
109
110        ###TAGGING ALL
111        word_sentences = []
112        tag_sentences = []
113
114        word_segment = []
115        tag_segment = []
116
117        for i,word in enumerate(list_of_words):
118            if word != ".":
119                word_segment.append(word)
120                tag_segment.append(list_of_tags[i])
121            else:
122                word_segment.append(word)
123                word_sentences.append(word_segment)
124                word_segment = []
125

```



```

126         tag_segment.append(list_of_tags[i])
127         tag_sentences.append(tag_segment)
128         tag_segment = []
129
130     for i,sentence in enumerate(word_sentences):
131         for j,word in enumerate(sentence):
132             tag = tag_sentences[i][j]
133             if tag != "0":
134                 word_sentences[i][j] = tag[2:]
135
136     for i,sentence in enumerate(word_sentences):
137         word_sentences[i] = " ".join(sentence)
138
139     df["Tagged_All"][index] = " ".join(word_sentences)
140
141
142
143
144
145     #TAGGING SEGMENTS AS ONE
146     previous = ''
147     segment_list = []
148     tag_list = []
149     temp1 = []
150     temp2 = []
151
152     for index2, tag in enumerate(list_of_tags):
153         if tag == previous or previous == '':
154             temp1.append(list_of_words[index2])
155             temp2.append(tag)
156         elif tag != previous:
157             segment_list.append(temp1.copy())
158             tag_list.append(temp2.copy())
159
160             temp1.clear()
161             temp2.clear()
162
163             temp1.append(list_of_words[index2])
164             temp2.append(tag)
165
166             previous = tag
167
168     segment_list.append(temp1)
169     tag_list.append(temp2)
170

```

```

171     for index3, group in enumerate(segment_list):
172         segment_list[index3] = " ".join(group)
173         tag_list[index3] = tag_list[index3][0]
174
175     #print(segment_list)
176     #print(tag_list)
177
178     for index4, thingy in enumerate(segment_list):
179         new_tag = tag_list[index4]
180         if new_tag != "0":
181             segment_list[index4] = new_tag[2:]
182
183
184     df["Tagged_One"][index] = " ".join(segment_list)
185
186
187
188
189
190
191
192     #TAGGING UNIQUE
193     Entity_List_New = Entity_List.copy()
194     list_of_words_tagged = df["Tagged_One"][index].split(" ")
195     #list_of_words = df["Article"][index].split(" ")
196     #If time, make numbering unique i.e. if band shows up twice give
same # for it
197
198     for i, word in enumerate(list_of_words_tagged):
199         if word in Entity_List:
200             Tag_Index = Entity_List.index(word)
201             Original_Tag = Entity_List_New[Tag_Index]
202
203             if Original_Tag[-1].isdigit():
204                 New_Tag = Original_Tag[:-1]+str(int(Original_Tag[-1])
+1)
205
206             else:
207                 New_Tag = Original_Tag+"1"
208
209             Entity_List_New[Tag_Index] = New_Tag
210             list_of_words_tagged[i] = New_Tag
211
212     df["Tagged_Uni"][index] = " ".join(list_of_words_tagged)
213
214     return df

```

```

214
215
216
217
218
219
220 #####
221 ### STATS ###
222 #####
223 df = tagged_dataframe()
224
225 WPA = []
226 SPA = []
227 EPA = {"O": [], "I-ORG": [], "I-PER": [], "I-LOC": [], "I-MISC": []}
228 CPA = []
229
230 for index, row in df.iterrows():
231
232     list_of_words = df["Article"][index].split(" ")
233     list_of_tags = df["Entity"][index].split(" ")
234
235     list_of_sentence_words = df["Article"][index].split(". ")
236
237     words_per_article = len(list_of_words)
238     sent_per_article = len(list_of_sentence_words)
239
240     WPA.append(words_per_article)
241     SPA.append(sent_per_article)
242
243     for key in EPA.keys():
244         EPA[key].append(list_of_tags.count(key))
245
246
247     cnt = 0
248     for word in list_of_words:
249         if word in ["(", ")", ".", ",", ">", "<", "[", "]", "-", "{", "}", ":", ";", ":
250             cnt+=1
251
252     CPA.append(cnt)
253
254 sum(CPA)
255 sum(EPA["I-ORG"]+EPA["I-LOC"]+EPA["I-PER"]+EPA["I-MISC"])
256
257 SPA_unique = np.unique(SPA).tolist()
258 SPA_counts = [SPA.count(num) for num in SPA_unique]

```

```

259
260 WPA_unique = np.unique(WPA).tolist()
261 WPA_counts = [WPA.count(num) for num in WPA_unique]
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297 #####
298 ##### ARTICLE VERS #####
299 #####
300 import ast
301 sizes = [50,100,250,500]
302
303 id_dict = {}

```

```

304 for size in sizes:
305     with open("H:\My Files\School\Grad School WLU\MRP\Research\Files\Data\
Textfiles\Article\\"+str(size)+"\000_ID_LIST.txt") as indx_lst:
306         for line in indx_lst:
307             version = line[1]
308             ids = ast.literal_eval(line[4:])
309             id_dict[str(size)+"v"+version] = ids
310
311
312 set_vers = "Train"
313 df=gold_dataframe()
314 sent_length_dict = {}
315 art_cnt_dict = {}
316 O_cnt_dict = {}
317 NO_cnt_dict = {}
318 for size in [50,100,250,500]:
319     batch_avg_sent_length=[]
320     batch_tot_sent_length=[]
321     vers_art_cnt = []
322     vers_O_cnt = []
323     vers_NO_cnt = []
324
325     for version in range(0,10):
326         index_ids_to_use = id_dict[str(size)+"v"+str(version)]
327
328         if set_vers == "Test":
329             index_ids_to_use = list(set(list(range(0,145)))-set(
index_ids_to_use))
330
331         new_df=df.iloc[index_ids_to_use]
332
333         sentence_lengths = []
334         article_cnt = 0
335         O_tag_cnt = []
336         NO_tag_cnt = []
337         for index,row in new_df.iterrows():
338             sentence_lengths.append(SPA[index])
339             article_cnt+=1
340
341             tag_splitter = new_df["Entity"][index].split(" ")
342             O_tag_cnt.append(tag_splitter.count("O"))
343             NO_tag_cnt.append(tag_splitter.count("I-PER")+tag_splitter.
count("I-LOC")+tag_splitter.count("I-ORG")+tag_splitter.count("I-MISC"
))
344

```

```

345     vers_art_cnt.append(article_cnt)
346     tot_sentences = sum(sentence_lengths)
347     avg_sentence_length = sum(sentence_lengths)/len(sentence_lengths)
348     batch_avg_sent_length.append(avg_sentence_length)
349     batch_tot_sent_length.append(tot_sentences)
350
351     vers_O_cnt.append(sum(O_tag_cnt))
352     vers_NO_cnt.append(sum(NO_tag_cnt))
353
354     avg_batch_sent_length = sum(batch_avg_sent_length)/len(
batch_avg_sent_length)
355     tot_batch_sent_length = sum(batch_tot_sent_length)/len(
batch_tot_sent_length)
356     tot_batch_art_cnt = sum(vers_art_cnt)/len(vers_art_cnt)
357     avg_O_cnt = sum(vers_O_cnt)/len(vers_O_cnt)
358     avg_NO_cnt = sum(vers_NO_cnt)/len(vers_NO_cnt)
359
360     art_cnt_dict[size]=tot_batch_art_cnt
361     sent_length_dict[size]=tot_batch_sent_length
362     O_cnt_dict[size]=avg_O_cnt
363     NO_cnt_dict[size]=avg_NO_cnt
364
365     print(sent_length_dict)
366     print(art_cnt_dict)
367     print(O_cnt_dict)
368     print(NO_cnt_dict)
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387

```

```
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425 #####
426 #### PLOT ####
427 #####
428 #https://matplotlib.org/stable/gallery/statistics/histogram_multihist.html
429
430 import matplotlib.pyplot as plt
431
432
```

```

433 #SENTENCES
434 fig = plt.figure()
435 ax = fig.add_axes([0,0,1,1])
436 langs = SPA_unique
437 students = SPA_counts
438 ax.bar(langs,students)
439 plt.xlabel("Sentences in Articles")
440 plt.ylabel("# of Occurences")
441 plt.title('Frequency Distribution - # of Sentences in Articles')
442 plt.text(15, 8, 'Avg # of Sentences per Article: '+str(round(sum(SPA)/len(
    SPA),2)))
443 plt.show()
444
445 #WORDS
446 fig = plt.figure()
447 ax = fig.add_axes([0,0,1,1])
448 langs = WPA_unique
449 students = WPA_counts
450 ax.bar(langs,students)
451 plt.xlabel("Sentences in Articles")
452 plt.ylabel("# of Occurences")
453 plt.title('Frequency Distribution - # of Sentences in Articles')
454 plt.text(6, 3, 'Avg # of Sentences per Article: '+str(round(sum(SPA)/len(
    SPA),2)))
455 plt.show()
456
457
458 #ENTITIES
459 colors = ["red","blue","green","black","orange"][1]
460 entity_lab = list(EPA.keys())[1]
461 entity_cnt = list(EPA.values())[1]
462 fig = plt.figure()
463 plt.hist(entity_cnt, density=False, bins=50, stacked=True, color=colors,
    label=entity_lab)
464 plt.legend(prop={'size': 10})
465 plt.ylabel('# of Occurences')
466 plt.xlabel('# of Words per Article');
467 plt.show()

```



## Train-Test Split

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Tue Aug 9 15:21:27 2022
4
5 @author: Doug
6 """
7 #####
8 # LIBRARIES
9 #####
10
11 import pandas as pd
12 import csv
13 import numpy as np
14 import tensorflow as tf
15 import time
16
17 #For Copies
18 import copy
19
20 #For Synonyms
21 import requests
22 from bs4 import BeautifulSoup
23
24 #For Shuffle
25 import random
26
27 #Hugging Face
28 from transformers import PegasusTokenizer, PegasusForConditionalGeneration
29 , T5Tokenizer, T5ForConditionalGeneration, MT5Tokenizer,
30 MT5ForConditionalGeneration, AutoModel, AutoTokenizer
31
32 #OS
33 import os.path
34 from os import path as os_path
35
36 #ITERTOOLS
37 import itertools
38
39 #ROUGE METRIC
40 from rouge_score import rouge_scorer
41
42 #FACTORIAL
43 import math
```

```

42
43 #NLTK
44 import nltk
45 nltk.download("punkt")
46
47 #MEMORY CLEARING
48 from GPUUtil import showUtilization as gpu_usage
49 import torch
50 from numba import cuda
51
52 #####
53 ##### CACHE #####
54 #####
55
56 #Change Cache
57 import os
58 os.environ['TRANSFORMERS_CACHE'] = 'H:/TempHF_Cache/cache/transformers/'
59 os.environ['HF_HOME'] = 'H:/TempHF_Cache/cache/'
60 os.environ['XDG_CACHE_HOME'] = 'H:/TempHF_Cache/cache/'
61
62 #####
63 # CODE
64 #####
65
66 df = pd.read_csv("H:/My Files/School/Grad School WLU/MRP/Research/Files/
    Data/dataframes/Article.csv", encoding="utf-8", index_col=0)
67
68 WPA = []
69 SPA = []
70 EPA = {"O": [], "I-ORG": [], "I-PER": [], "I-LOC": [], "I-MISC": []}
71 CPA = []
72
73 for index, row in df.iterrows():
74
75     list_of_words = df["Article"][index].split(" ")
76     list_of_tags = df["Entity"][index].split(" ")
77
78     list_of_sentence_words = df["Article"][index].split(". ")
79
80     words_per_article = len(list_of_words)
81     sent_per_article = len(list_of_sentence_words)
82
83     WPA.append(words_per_article)
84     SPA.append(sent_per_article)
85

```

```

86     for key in EPA.keys():
87         EPA[key].append(list_of_tags.count(key))
88
89
90     cnt = 0
91     for word in list_of_words:
92         if word in ["(", ")", ".", ",", ">", "<", "[", "]", "-", "{", "}", ":", ";"]:
93             cnt+=1
94
95     CPA.append(cnt)
96
97 sum(CPA)
98 sum(EPA["I-ORG"]+EPA["I-LOC"]+EPA["I-PER"]+EPA["I-MISC"])
99
100 SPA_unique = np.unique(SPA).tolist()
101 SPA_counts = [SPA.count(num) for num in SPA_unique]
102
103 WPA_unique = np.unique(WPA).tolist()
104 WPA_counts = [WPA.count(num) for num in WPA_unique]
105
106
107
108
109
110 def rng_train_set(min_exclusion, max_exclusion, sentence_goal, bounds_rate
111 ):
112     if min_exclusion<0 or max_exclusion<0 or min_exclusion>=max_exclusion
113     or sentence_goal<0 or bounds_rate<0:
114         print("Error in Variable Settings! Exited.")
115         return
116
117     in_bounds = False
118
119     while in_bounds == False:
120         original_article_list = list(range(0,145))
121
122         filter_out_list = []
123         filter_cnt = 0
124         for i,length in enumerate(SPA):
125             if length <= min_exclusion or length >= max_exclusion:
126                 filter_out_list.append(i)
127                 filter_cnt += length
128
129         new_article_list = list(set(original_article_list)-set(
130 filter_out_list))

```

```

128
129     new_cnt = 0
130     for index in new_article_list:
131         new_cnt+=SPA[index]
132
133     #Confirm
134     if filter_cnt+new_cnt!=1768 and len(filter_out_list)+len(
new_article_list)!=145:
135         print("ERROR IN SENTENCE OR ARTICLE COUNT")
136
137     #Get Avg SPA
138     avg_spa = new_cnt / len(new_article_list)
139
140     expected_iterations = math.ceil(sentence_goal/avg_spa)
141     train_articles = random.sample(new_article_list,
expected_iterations)
142     test_articles = list(set(original_article_list)-set(train_articles
))
143
144
145     total_sentences = 0
146     for article_index in train_articles:
147         total_sentences += SPA[article_index]
148
149     avg_sentences = total_sentences/len(train_articles)
150
151     #Check Bounds Acceptance
152     if bounds_rate >= 0 and bounds_rate < 1:
153         if total_sentences >= sentence_goal*(1-bounds_rate) and
total_sentences <= sentence_goal*(1+bounds_rate):
154             in_bounds = True
155         else:
156             if total_sentences >= sentence_goal - bounds_rate and
total_sentences >= sentence_goal + bounds_rate:
157                 in_bounds = True
158
159     return train_articles, test_articles, total_sentences, avg_sentences#,
filter_out_list, new_article_list, avg_spa
160
161 #Run Replications for Sample Sizes
162 replications = 10
163 batch_pools = [50,100,250,500]
164 train_batch_dict = {50:[],100:[],250:[],500:[]}
165 test_batch_dict = {50:[],100:[],250:[],500:[]}
166

```

```

167 for size in batch_pools:
168     print("Batch Size: "+str(size))
169     for i in range(replications):
170         train_list, test_list, cnt_sentences, avg_sentences =
171         rng_train_set(
172             min_exclusion=2
173             , max_exclusion=40
174             , sentence_goal=size
175             , bounds_rate=0.1
176         )
177         train_batch_dict[size].append(train_list)
178         test_batch_dict[size].append(test_list)
179
180
181
182
183
184
185
186
187
188 #Select Scored Samples
189 def final_selection(augmented_sentence_multiplier):
190     variants = ["Article", "Tagged_One", "Tagged_Uni"]
191     augment_multiple=3#augmented_sentence_multiplier
192
193     #Set Sample Headers to Iterate
194     headers = []
195     for i in range(1,51):
196         headers.append("Sample"+str(i))
197
198     nested_dict = {}
199     for variant in variants:
200         nested_dict[variant] = {}
201         for size in train_batch_dict:
202             nested_dict[variant][size] = {}
203             for attempt in range(replications):
204                 nested_dict[variant][size][attempt] = {}
205                 df=pd.read_csv("H:/My Files/School/Grad School WLU/MRP/
Research/Files/Data/dataframes/Finished/Original_"+variant+".csv",
encoding="utf-8", index_col=0)
206                 df_scored = pd.read_csv("H:/My Files/School/Grad School
WLU/MRP/Research/Files/Data/dataframes/Finished/Scored_"+variant+".csv
", encoding="utf-8", index_col=0)

```

```

207         df=df.iloc[train_batch_dict[size][attempt]]
208
209
210         for index,row in df.iterrows():
211             possible_shuffles = math.factorial(SPA[index])
212             NA_samples = list(df_scored.iloc[index][headers]).
count("NA")
213             LowScore_samples = list(df_scored.iloc[index][headers
]).count("Low Score")
214
215             max_possible = min(50-NA_samples-LowScore_samples,
possible_shuffles)
216             samples_to_take = min(SPA[index]*augment_multiple,
max_possible)
217
218             score_list = list(df_scored.iloc[index][headers])
219             for i,score in enumerate(score_list):
220                 if score == "NA" or score == "Low Score":
221                     score_list[i]=0
222                 else:
223                     score_list[i]=float(score_list[i])
224
225             index_order = sorted(range(len(score_list)), key=
lambda k: score_list[k], reverse=True)
226             final_list = index_order[0:samples_to_take]
227             final_list = ["Sample"+str(x+1) for x in final_list]
228
229             cleaned_list=[]
230             for j,sample in enumerate(final_list):
231                 if df[sample][index]!=" " or df[sample][index]!="
NA" or not pd.isnull(df[sample][index]):
232                     #print(variant+"-"+str(size)+"-"+str(attempt)
+ "-Article ID "+str(index)+sample)
233                     cleaned_list.append(sample)
234                     #break
235
236             #nested_dict[variant][size][attempt][index] =
final_list
237             nested_dict[variant][size][attempt][index] =
cleaned_list
238
239
240             #print(nested_dict[variant][size][attempt][index])
241             #print("\n")
242

```

```

243     return nested_dict
244
245 fin_results = final_selection(3)
246
247
248
249
250
251
252
253
254
255 #####
256 ##### CREATE TXTS #####
257 #####
258 full_list = list(range(0,145))
259
260 for variant in list(fin_results.keys()): #Variants (Article, One Tag,
    Unique Tag)
261     #Set DFs
262     df_og=pd.read_csv("H:/My Files/School/Grad School WLU/MRP/Research/
    Files/Data/dataframes/Finished/Original_"+variant+".csv", encoding="
    utf-8", index_col=0)
263     df_mp=pd.read_csv("H:/My Files/School/Grad School WLU/MRP/Research/
    Files/Data/dataframes/Finished/Mapped_"+variant+".csv", encoding="utf
    -8", index_col=0)
264
265     for size in list(fin_results[variant].keys()): #Sizes (50/100/250/500)
266         #Write Article IDs for Reference per Repetition Version
267         with open('H:/My Files/School/Grad School WLU/MRP/Research/Files/
    Data/Textfiles/'+variant+'/' +str(size)+'000_ID_LIST.txt', 'w',
    encoding="utf-8") as f_ids:
268             for repetition in list(fin_results[variant][size].keys()): #
    Repetition (0-->10)
269                 #Write TOKEN/TAGs to 3 text files, test, original, and
    original+augmented
270                 with open('H:/My Files/School/Grad School WLU/MRP/Research
    /Files/Data/Textfiles/'+variant+'/' +str(size)+'v'+str(repetition)+'
    _Augmented.txt', 'w', encoding="utf-8") as f_aug, open('H:/My Files/
    School/Grad School WLU/MRP/Research/Files/Data/Textfiles/'+variant+'/'
    +str(size)+'v'+str(repetition)+'_UnAugmented.txt', 'w', encoding="utf
    -8") as f_org, open('H:/My Files/School/Grad School WLU/MRP/Research/
    Files/Data/Textfiles/'+variant+'/' +str(size)+'v'+str(repetition)+'
    _Testing.txt', 'w', encoding="utf-8") as f_tst:
271                     f_ids.write("v"+str(repetition)+":\t")

```

```

272         f_ids.write(str(list(fin_results[variant][size][
repetition].keys()))))
273         f_ids.write("\n")
274
275         #Iterate through Articles, Then Samples
276         for article_id in list(fin_results[variant][size][
repetition].keys()): #Article IDs
277             #Get the IDs for Articles
278             train_ids = list(fin_results[variant][size][
repetition].keys())
279             test_ids = list(set(full_list)-set(train_ids))
280
281             #token_list_train = []
282             #tag_list_train = []
283
284             df_og_train = df_og.iloc[train_ids]
285             df_mp_train = df_mp.iloc[train_ids]
286             df_og_test = df_og.iloc[test_ids]
287             #df_mp_test = df_mp.iloc[test_ids]
288
289             df_og_train=df_og_train.fillna("NA")
290             df_mp_train=df_mp_train.fillna("NA")
291             df_og_test=df_og_test.fillna("NA")
292             #df_og_train(replace)
293
294             #Get Augmented Samples
295             sample_list = fin_results[variant][size][
repetition][article_id]
296             if sample_list != []:
297                 for sample in sample_list:
298                     #SKIP NAs
299                     if df_og_train[sample][article_id]=="NA":
300                         continue
301
302                     tokens_train = df_og_train[sample][
article_id].split(" ")
303                     tags_train = df_mp_train[sample][
article_id].split(" ")
304
305                     if tokens_train[-1]!=".": tokens_train.
append("."), tags_train.append("0")#, print("Added a period on article
"+str(article_id))
306
307                     #token_list_train.append(tokens_train)
308                     #tag_list_train.append(tags_train)

```



```

309
310         #CONFIRM NO MISMATCH
311         for index,entry in enumerate(tags_train):
312             if len(tags_train)!=len(tokens_train):
313                 print("MISMATCH on ARTICLE "+str(
article_id)+" and "+sample)
314
315         for index,the_token in enumerate(
tokens_train):
316             to_append = the_token+" "+tags_train[
index]
317             f_aug.write(to_append)
318             f_aug.write("\n")
319
320             if the_token==".":
321                 f_aug.write('\n')
322
323             f_aug.write("-DOCSTART- 0\n\n") #SWAP BACK
TO "-DOCSTART- 0\n\n" LATER
324             #f_aug.write("-DOCSTART- 0\n\n") #SWAP BACK
325
326             #GET ORIGINAL UN-AUGMENTED SAMPLES FOR TRAINING
327             #GET ORIGINAL UN-AUGMENTED SAMPLES FOR TRAINING
328             #GET ORIGINAL UN-AUGMENTED SAMPLES FOR TRAINING
329             tokens_train = df_og_train["Article"][article_id].
split(" ")
330             tags_train = df_og_train["Entity"][article_id].
split(" ")
331             if tokens_train[-1]!=".": tokens_train.append(".")
, tags_train.append("0")#, print("Added a period on article"+str(
article_id))
332             #CONFIRM NO MISMATCH
333             for index,entry in enumerate(tags_train):
334                 if len(tags_train)!=len(tokens_train):
335                     print("MISMATCH on ARTICLE "+str(
article_id)+" and "+sample)
336             for index,the_token in enumerate(tokens_train):
337                 to_append = the_token+" "+tags_train[index]
338                 f_aug.write(to_append)
339                 f_aug.write("\n")
340                 f_org.write(to_append)
341                 f_org.write("\n")
342
343             if the_token==".":
344                 f_aug.write('\n')

```

```

345         f_org.write('\n')
346
347         f_aug.write("-DOCSTART- O\n\n") #SWAP BACK TO "-
DOCSTART- O\n\n" LATER
348         f_org.write("-DOCSTART- O\n\n") #SWAP BACK TO "-
DOCSTART- O\n\n" LATER
349
350         #GET ORIGINAL UN-AUGMENTED SAMPLES FOR TESTING
351         #GET ORIGINAL UN-AUGMENTED SAMPLES FOR TESTING
352         #GET ORIGINAL UN-AUGMENTED SAMPLES FOR TESTING
353         for a_index,row in df_og_test.iterrows():
354             tokens_train = df_og_test["Article"][a_index].
split(" ")
355             tags_train = df_og_test["Entity"][a_index].split("
")
356             if tokens_train[-1]!=".": tokens_train.append(".")
, tags_train.append("O")#, print("Added a period on article"+str(
article_id))
357             #CONFIRM NO MISMATCH
358             for index,entry in enumerate(tags_train):
359                 if len(tags_train)!=len(tokens_train):
360                     print("MISMATCH on ARTICLE "+str(a_index)+
" and "+sample)
361             for index,the_token in enumerate(tokens_train):
362                 to_append = the_token+" "+tags_train[index]
363                 f_tst.write(to_append)
364                 f_tst.write("\n")
365
366                 if the_token==".":
367                     f_tst.write('\n')
368
369             f_tst.write("-DOCSTART- O\n\n") #SWAP BACK TO "-
DOCSTART- O\n\n" LATER

```