

# Data Extraction using Semantic Similarity

Arunim Garg, Dhivya Chandrasekaran, Mohiuddin MD Abdul Qudar  
Kazi Zainab Khanam

**Abstract**—Natural Language Processing tasks are finding greater significance by Artificial Intelligence experts over the years and one of the major challenges faced by them is the sparsity of structured text datasets to enable rigorous training of developed models. In this article we propose a data extraction technique that extracts aligned sentence pairs from parallel corpora. Aligned sentences are extracted from the Newsela simplified datasets and simple English Wikipedia where the datasets are used to increase the performance of the auto-encoder sequence-to-sequence text simplification model. The performance is seen to improve with the availability of quality datasets. The major focus of the project is on implementing the semantic similarity algorithm and show the impact of the extracted data on the text simplification algorithm.

## I. INTRODUCTION

With the exponential increase in text data generated over time, Natural Language Processing (NLP) has gained significant attention from Artificial Intelligence (AI) experts. Measuring the semantic similarity between various text components like words, sentences, or documents plays a significant role in a wide range of NLP tasks like information retrieval [2], text summarization [3], text classification [4], essay evaluation [5], machine translation [6], question answering [7], [8], among others. One of the major challenges faced by various NLP tasks is the availability of efficient and sufficient data resources to train models in-order to achieve better performances. This paper implements a technique for data extraction using MASSAligner [1] to improve the performance of a simple sequence-to-sequence auto-encoder model for text simplification task. The remaining paper is organised as follows: Section II provides a brief overview on semantic similarity and text simplification tasks. Section III discusses methodology used in the project. Section IV tabulates the results and Section VI concludes with details of proposed future work and limitations.

## II. BACKGROUND

### A. Semantic Similarity

Text similarity is assessing the similarity between two words, sentences or expressions, based on the likelihood of how similar is their meaning. There are two predominant methods used for calculating the similarity between expressions, corpus-based or distributional semantic models (DSMs), and knowledge-based models. Knowledge-based methods use word senses, parts of speech and taxonomic information to calculate the similarity between expressions, whereas the corpus-based methods determine similarity based on the assumption that similar words occur in similar documents. The major drawback of corpus-based methods is that they ignore the sentence structure and the difference in word meaning based on the context, while the knowledge-based approach is restricted by the availability of human-crafted dictionaries. Recent researchers have worked on hybrid models combining the above methods to generate similar sentence pairs with better performance than the traditional methods [10], [11]. A clear understanding between ‘text-relatedness’ and ‘text-similarity’ is an integral part of text extraction based on similarity. Recent work by [12], [13], uses web search engine results for calculating word relatedness, where words with opposite meanings also have a high similarity score.

### B. Text Simplification

Most of the early work in text simplification involved extractive methods of summarization – extracting the sentences that convey the most meaning. These are easiest to implement and available readily online for free. The simplest of all summarization techniques is TF-IDF (Term Frequency – Inverse Document Frequency). This technique involves generating a frequency table for all the words in the corpus, after removing stop words. The

sentence weight is then calculated based on the word frequencies, and is normalized by sentence length, to calculate the sentences with the most “substance”. We can use a threshold to retain sentences with the most substance or summarize a given text in a given number of sentences by choosing the ones with the highest values. As research in NLP has branched widely over the past decade, primarily due to increased availability of computing resources, the research in simplification has shifted towards abstractive approaches – actual generation of text. Abstractive approaches have mostly focused on lexical or phrasal substitutions for sentence-level simplification. This approach is the simplest to implement as it involves minimal data preprocessing. Truly abstractive simplification approaches involve sentence splitting, text deletion, and addition. One of the approaches used for simplification, and the one we plan to adapt, involves seq2seq modeling, which also used in machine translation. We take aligned complex and simplified sentence pairs, instead of sentences in two languages, convert them to a fixed encoded vector and train an RNN or LSTM to be able to simplify our content. The performance of this approach relies heavily on the dataset being used to train the model.

### III. METHODOLOGY AND EXPERIMENTAL ANALYSIS

This section describes the semantic similarity model, text simplification model, and the two different datasets (Newsela simplified dataset and simple English Wikipedia datasets) used.

#### A. Datasets

We use two different datasets to provide a comparison of how the performance increase with a bigger dataset and the impact of semantic similarity in building the dataset.

- **English Wikipedia and Simple English Wikipedia (SEW):** Simple English Wikipedia is widely used as it is publicly available and because of the popularity of the regular English Wikipedia. SEW includes simplified versions of articles in regular English Wikipedia, and datasets are available with aligning “equivalent” sentences from the two, to allow seq2seq model training.

- **Newsela Corpus:** This dataset, created by Xu et al. contains news articles with 4 simplified versions for each, produced manually by professional editors. Corpus level simplification is available, however, it has to be processed for sentence-level simplification.

From these datasets, we use the below mentioned Semantic similarity model to extract complex and simple sentence pairs that are then fed to the text simplification autoencoder model to analyze the performance.

#### B. Semantic similarity Model:

The model we implement to extract semantically similar aligned sentence pairs was proposed by Paetzold et al., The MASSAligner is available as an open-source Python library. The steps involved can be divided into two sections:

1) **Measuring semantic similarity:** The model calculates the semantic similarity between sentence pairs using the TF-IDF model. The model converts documents to a bag of words, forms word vectors based on the normalized term frequency (frequency of a term / total number of terms) and finally calculates the cosine similarity between word vectors. The model returns a similarity matrix containing the similarity between all the sentences.

2) **Aligning sentences based on similarity:** The model follows a vicinity driven approach to extract the sentences that are similar based on a threshold value provided as a hyperparameter. The aligner then traverses through the similarity matrix to establish an alignment path that searches for the best pair of similar sentences. Given a coordinate  $[i, j]$  in a matrix there are three vicinities taken into consideration  $V1 = [i, j + 1], [i + 1, j], [i + 1, j + 1]$ ,  $V2 = [i + 1, j + 2], [i + 2, j + 1]$ , and  $V3$  as all remaining  $[x, y]$  where  $x > i$  and  $y > j$ . The initial point begins at the coordinate that is closest to the first point  $(0, 0)$  and has a similarity greater than 0.2 (hyperparameter achieved by [14]) (The algorithm searches the three vicinities to find the coordinates having the highest similarity. Thus, the similar pairs are extracted and we write them to two different files.

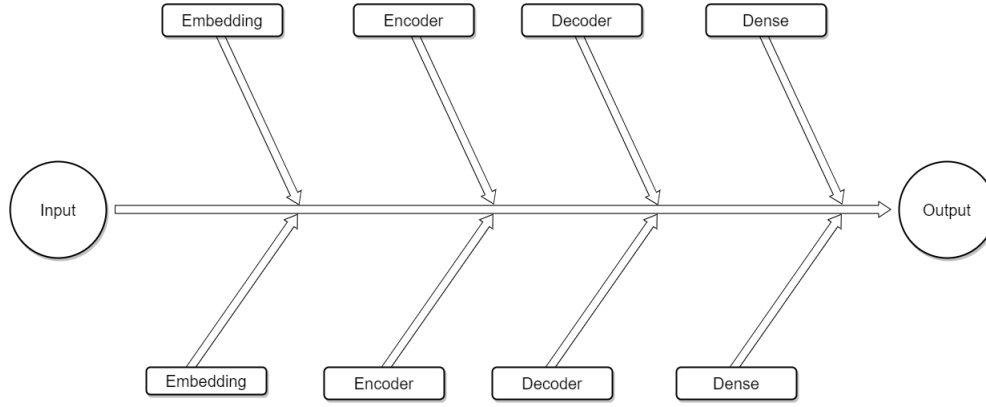


Fig. 1. Architecture of the Auto-encoder Text Simplification Model

### C. Text Simplification Model:

Text-simplification is done using a model inspired by a standard machine-translation model called the seq2seq model. This is an auto-encoder model that predicts the simplified version of a complex sentence. The model is simple and at present, uses a simple long short term memory (LSTM) model for encoder/decoder functionalities. The model consists of an embedding layer where the sentences are converted to vectors based on a comprehensive vocabulary built using all the words present in both the simplified and complex documents. The vectors are then passed through two LSTM layers of which one serves as an encoder; encoding the input sequence and producing internal state vectors which serve as conditioning for the decoder. The other would serve as the decoder that is responsible for predicting the target sequence. Figure 1 shows the architecture of the model.

## IV. RESULTS

We establish a baseline score using the precompiled Wikipedia dataset that has 167,000 aligned normal and simplified sentence pairs. We extract sentences from the Newsela corpus using the MASSAlign semantic similarity algorithm based on five different semantic similarity thresholds. Five datasets generated by using 30%, 40%, 50%, 60% and 70% similarity percentage between the normal and simplified sentence pairs. Table I shows the number of sentence pairs generated at each semantic similarity thresholds. From the Fig 2, we can see that the number of sentences generated is

highest when the similarity threshold is set at 30% and the least number of sentences are generated threshold is set at 70%.

SNo	Percentage Similarity (%)	Number of sentences
1	30	277082
2	40	249200
3	50	217680
4	60	184336
5	70	151588

TABLE I  
SENTENCE COUNT COMPARISON AT VARIOUS SIMILARITY LEVELS.

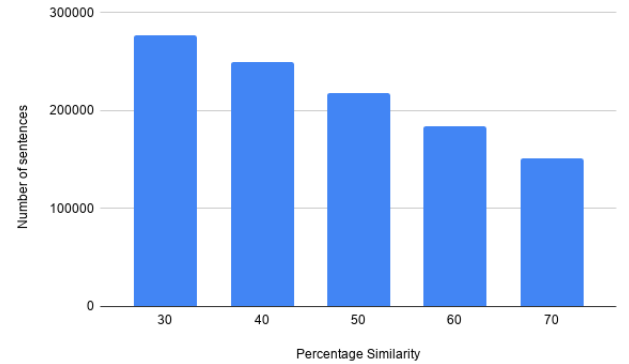


Fig. 2. Comparison of datasets generated using different similarity measures

We then use the Wikipedia based base-line model to perform hyperparameter tuning. We tune the following hyperparameters: sentence length, removing punctuation, removing stop\_words, number of

Hyper-parameters	Option 1	Option 2	Option 3	Best Option
Remove punctuation	True	False	–	<b>True</b>
Remove stopwords	True	False	–	<b>False</b>
Words cutoff	15	20	25	<b>15</b>
Aggr length	15	12	10	<b>15</b>
Hidden size	128	256	512	<b>256</b>
Batch size	128	64	32	<b>64</b>

TABLE II  
HYPERPARAMETER TUNING.

hidden neurons, and batch\_size using the baseline model. Table I shows the hyperparameters that we have tuned our model with and the optimal values of the different hyperparameters. The performance is measured using bilingual evaluation understudy (BLUE) score. The BLUE score is within the range between 0 and 1. The baseline model achieved the best results by removing punctuation, using batch size of 64, hidden size of 512, total aggregated length of 15 words in the simplified sentence and when the word limit(words\_cutoff) was set at 15 words as shown in table II.

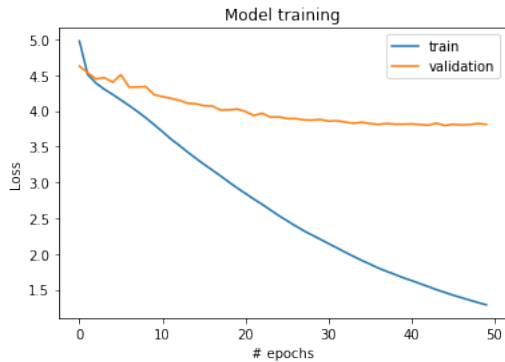


Fig. 3. Training and Validation loss score for our model with 50% sentence similarity with normal and simplified text file

The five different datasets extracted using our semantic similarity algorithm were used to enhance the performance of the base-line model. The performance of the model increases with the use of the Newsela dataset. However, the best results were obtained with the complex and simplified sentence pair dataset extracted at a semantic similarity threshold of 50% as shown in Table III. Fig 3 is a graphical representation of the variation in BLEU score over different datasets. The training and validation graph of the best performance achieved is shown in Fig 4

SNo	Semantic Similarity Threshold	BLEU Score
1	30 % similarity	0.1895
2	40 % similarity	0.2942
3	<b>50 % similarity</b>	<b>0.3134</b>
4	60 % similarity	0.3128
5	70 % similarity	0.2910

TABLE III  
PERFORMANCE OF THE MODEL WITH DIFFERENT DATASETS.

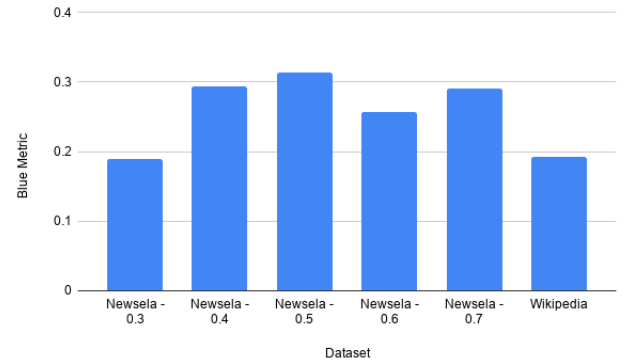


Fig. 4. Comparison of datasets generated using different similarity measures

## V. FUTURE WORK

This project has mainly focused on implementing the semantic similarity algorithm to perform data extraction from parallel corpora, hence the text simplification model used is rather simple. However, with better resources and more time, we plan to implement a more complex text simplification model and enhance the performance further.

## VI. AUTHORS' CONTRIBUTIONS

- Mohiuddin Md Qudar Abdul - Concept development, implementation and designing figures and tables. 25% percent.
- Arunim Garg - Concept development, implementation and proof reading. 25% percent.

- Dhivya Chandrasekaran - Concept development, implementation and report content development . 25% percent.
- Kazi Zainab - Concept development, implementation and log file compilation and designing figures and tables. 25% percent.

## REFERENCES

- [1] Paetzold, Gustavo, Fernando Alva-Manchego, and Lucia Specia. "Massalign: alignment and annotation of comparable documents." In Proceedings of the IJCNLP 2017, System Demonstrations, pp. 1-4. 2017.
- [2] Kim, Sun, Nicolas Fiorini, W. John Wilbur, and Zhiyong Lu. "Bridging the gap: Incorporating a semantic similarity measure for effectively mapping PubMed queries to documents." Journal of biomedical informatics 75 (2017): 122-127.
- [3] Mohamed, Muhidin, and Mourad Oussalah. "SRL-ESA-TextSum: A text summarization approach based on semantic role labeling and explicit semantic analysis." Information Processing & Management 56, no. 4 (2019): 1356-1372.
- [4] Kim, Yoon. "Convolutional neural networks for sentence classification." arXiv preprint arXiv:1408.5882 (2014).
- [5] Janda, Harneet Kaur, Atish Pawar, Shan Du, and Vijay Mago. "Syntactic, Semantic and Sentiment Analysis: The Joint Effect on Automated Essay Evaluation." IEEE Access 7 (2019): 108486-108503.
- [6] Zou, Will Y., Richard Socher, Daniel Cer, and Christopher D. Manning. "Bilingual word embeddings for phrase-based machine translation." In Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, pp. 1393-1398. 2013.
- [7] Bordes, Antoine, Sumit Chopra, and Jason Weston. "Question answering with subgraph embeddings." arXiv preprint arXiv:1406.3676 (2014).
- [8] Lopez-Gazpio, Iñigo, Montse Maritxalar, Aitor Gonzalez-Agirre, German Rigau, Larraitz Uribe, and Eneko Agirre. "Interpretable semantic textual similarity: Finding and explaining differences between sentences." Knowledge-Based Systems 119 (2017): 186-199.
- [9] D. Bollegala, Y. Matsuo, and M. Ishizuka, "Measuring semantic similarity between words using Web search engines," in Proc. WWW, vol. 7, 2007, pp. 757–766.
- [10] D. Bollegala, Y. Matsuo, and M. Ishizuka, "Measuring semantic similarity between words using Web search engines," in Proc. WWW, vol. 7, 2007, pp. 757–766.
- [11] R. L. Cilibrasi and P. M. B. Vitanyi, "The Google similarity distance," IEEE Trans. Knowl. Data Eng., vol. 19, no. 3, pp. 370–383, Mar. 2007.
- [12] Atoum, Issa, and Ahmed Otoom. "Efficient hybrid semantic text similarity using WordNet and a corpus." Int. J. Adv. Comput. Sci. Appl 7, no. 9 (2016): 124-130.
- [13] R. Ferreira, R. D. Lins, F. Freitas, S. J. Simske, and M. Riss, "A New Sentence Similarity Assessment Measure Based on a Three-layer Sentence Representation," in Proceedings of the 2014 ACM Symposium on Document Engineering, 2014, pp. 25–34
- [14] Paetzold, Gustavo Henrique, and Lucia Specia. "Vicinity-driven paragraph and sentence alignment for comparable corpora." arXiv preprint arXiv:1612.04113 (2016).

## APPENDIX

### A. Importing the necessary libraries

```

1 import os
2 from sklearn.utils import shuffle
3 import math
4 from nltk.tokenize import word_tokenize,
   sent_tokenize
5 from massalign.core import *

1 import nltk
2 nltk.download('stopwords')
3 nltk.download('punkt')
4 import nltk
5 from nltk.corpus import stopwords
6 from nltk.tokenize import word_tokenize
7 stop_words = set(stopwords.words('english'))
8 import string
9 import re
10 import datetime
11 import time
12 from numpy import array, argmax, random, take
13 import pandas as pd
14 from keras.models import Sequential
15 from sklearn.model_selection import
   train_test_split
16 from keras.layers import Dense, LSTM,
   Embedding, Bidirectional, RepeatVector,
   TimeDistributed
17 from keras.preprocessing.text import Tokenizer
18 from keras.callbacks import ModelCheckpoint
19 from keras.preprocessing.sequence import
   pad_sequences
20 from keras.models import load_model
21 from keras import optimizers
22 import matplotlib.pyplot as plt
23 % matplotlib inline
24 pd.set_option('display.max_colwidth', 200)
25 from nltk.translate.bleu_score import
   sentence_bleu

```

### B. Uploading the Newsela Extractor.py file

```

1 def files(path):
2     for file in os.listdir(path):
3         if os.path.isfile(os.path.join(path,
4             file)):
5             yield file

```

### C. Formatting the raw Newsela data

```

1 #formatting Newsela raw data into format
   needed by aligner
2 rawDataPath = 'NewselaRaw'
3 formattedDataPath = 'NewselaFormatted'
4 ignoreFiles = []
5 for f in files(rawDataPath):
6     print(f)
7     if f in ignoreFiles:
8         continue
9     fileLines = [line.rstrip('\n') for line in
   open(rawDataPath+ '/' +f)]
10    firstSent = fileLines[0]

```

```

11 if firstSent.split(' ')[0].isupper():
12     fileLines[0] = firstSent[firstSent.
    find("  ")+3:].strip()
13
14     file_content = ''
15     for l in fileLines:
16         if l != '' and (not l.startswith("##")):
17             file_content += l.strip() + ' '
18
19     try:
20         fileSentences = sent_tokenize(
            file_content)
21
22         outFile = open(formattedDataPath + '/' + f, "w")
23         for s in fileSentences:
24             # write line to output file
25             outFile.write(s)
26             outFile.write("\n")
27             outFile.write("\n")
28         outFile.close()
29     except:
30         ignoreFiles.append(f)

```

*D. Extracting the list of articles from the Newsela corpus*

```

1 #get list of all articles in the corpus
2 articles = [line.rstrip('\n') for line in open
    ('articleNames.csv')]
3 artDict = {}
4 for a in articles:
5     tmpA = a.split(',')
6     if tmpA[0] in artDict.keys():
7         artDict[tmpA[0]].append(tmpA[5])
8     else:
9         artDict[tmpA[0]] = [tmpA[5]]

```

*E. Generating the Normal and Simplified text files*

```

1 path = 'alignedSentences-0.5.txt'
2 #output separated files for normal and
    simplified sentences
3 outSimpFile = open('alignedSentencesSimp-0.5.
    txt', "w")
4 outNormFile = open('alignedSentencesNorm-0.5.
    txt', "w")
5 lines = [line.rstrip('\n') for line in open(
    path)]
6 for l in lines:
7     outNorm = l.split('\t')[0][3:-2] + "\n"
8     outSimp = l.split('\t')[1][3:-2] + "\n"
9     outNormFile.write(outNorm)
10    outSimpFile.write(outSimp)
11 outNormFile.close()
12 outSimpFile.close()

```

*F. For each of the sentence similarity percentage, combining the aligned sentence pairs*

```

1 for sim in simVals:
2     #for each similarity value from 0.3 to 0.7

```

```

3 fileOutput = open("alignedSentences-" +
    str(sim) + ".txt", "w")
4 for k in artDict.keys():
5     #combine data into aligned sentence
    pairs
6     print k
7     for i in range(len(artDict[k])-1):
8         for j in range(i+1, len(artDict[k]
    ))):
9             f1 = 'NewselaFormatted/' +
                artDict[k][i]
10             f2 = 'NewselaFormatted/' +
                artDict[k][j]
11
12             try:
13                 aps = getAlignedParagraphs
                    (f1, f2, sim)
14                 for a in aps:
15                     s = str(a[0]) + "\t" +
                        str(a[1]) + "\n"
16                     fileOutput.write(s)
17             except:
18                 continue
19 fileOutput.close()

```

*G. Opening the aligned sentence pairs text file and pre-processing the text files and building our model*

```

1 def read_text(filename):
2     # open the file
3     file = open(filename, mode='rt', encoding=
        'utf-8')
4     # read all text
5     text = file.read()
6     file.close()
7     return text
8
9 # split a text into sentences
10 def to_lines(text):
11     sents = text.strip().split('\n')
12     sents = [i.split('\t') for i in sents]
13     return sents
14
15 # function to build a tokenizer
16 def tokenization(lines):
17     tokenizer = Tokenizer()
18     tokenizer.fit_on_texts(lines)
19     return tokenizer
20
21 # encode and pad sequences
22 def encode_sequences(tokenizer, length, lines)
    :
23     # integer encode sequences
24     seq = tokenizer.texts_to_sequences(lines)
25     # pad sequences with 0 values
26     seq = pad_sequences(seq, maxlen=length,
        padding='post')
27     return seq
28
29 # build NMT model
30 def build_model(in_vocab, out_vocab,
    in_timesteps, out_timesteps, hidden_size):
31     model = Sequential()

```

```

32 model.add(Embedding(in_vocab, hidden_size,
33                     input_length=in_timesteps, mask_zero=True
34                     ))
35 model.add(LSTM(hidden_size))
36 model.add(RepeatVector(out_timesteps))
37 model.add(LSTM(hidden_size,
38               return_sequences=True))
39 model.add(Dense(out_vocab, activation='
40               softmax'))
41 return model
42
43 def get_word(n, tokenizer):
44     for word, index in tokenizer.word_index.
45     items():
46         if index == n:
47             return word
48     return None
49
50 def calc_BLUE(pred_df):
51     blue = 0
52     for i in range(len(pred_df)):
53         reference = [pred_df.iloc[i]['actual'
54         ].split()]
55         candidate = pred_df.iloc[i]['predicted
56         '].split()
57         blue = blue + sentence_bleu(reference,
58         candidate)
59     return blue/len(pred_df)
60
61 def removeSentStopWords(sent):
62     word_tokens = word_tokenize(sent.lower())
63     filtered_sentence = [w for w in word_tokens
64     if not w in stop_words]
65     newSent = ''
66     for s in range(len(filtered_sentence)):
67         if s < len(filtered_sentence) - 1:
68             if filtered_sentence[s+1].strip() in
69             string.punctuation:
70                 newSent = newSent + filtered_sentence[
71                 s]
72             else:
73                 newSent = newSent + filtered_sentence[
74                 s] + ' '
75         else:
76             newSent = newSent + filtered_sentence[s]
77     return newSent
78
79 H. Building our auto-encoder model and initializ-
80 ing method to calculate bleu score
81
82 def RunAutoEncoder(usingNewsela=True,
83 word_cutoff = 15, aggr_length = 12,
84 trainTestSplitPerc = 0.3,
85 hidden_size = 64,
86 remove_punc = True, batchSize=64,
87 remove_stopwords = False,
88 normFilePath='
89 NewselaAlignedNorm-0.4.txt', simpFilePath=
90 'NewselaAlignedSimp-0.4.txt'):
91     if usingNewsela:
92         data_normal = read_text(normFilePath)
93         norm_data = to_lines(data_normal)
94         norm_data = array(norm_data)
95
96         norm_data = [s[0] for s in norm_data]
97
98         data_simple = read_text(simpFilePath)
99         simp_data = to_lines(data_simple)
100         simp_data = array(simp_data)
101         simp_data = [s[0] for s in simp_data]
102     else:
103         data_normal = read_text(normFilePath)
104         norm_data = to_lines(data_normal)
105         norm_data = array(norm_data)
106         norm_data = [s[2] for s in norm_data]
107
108         data_simple = read_text(simpFilePath)
109         simp_data = to_lines(data_simple)
110         simp_data = array(simp_data)
111         simp_data = [s[2] for s in simp_data]
112
113     #Remove stop words
114     if remove_stopwords:
115         print("Removving stop words...")
116         for n in range(len(norm_data)):
117             norm_data[n] = removeSentStopWords(
118             norm_data[n])
119         for s in range(len(simp_data)):
120             simp_data[s] = removeSentStopWords(
121             simp_data[s])
122     allData= []
123     allData_aggr=[]
124     for i in range(len(norm_data)):
125         if (len(norm_data[i].split()) <
126         word_cutoff) and (len(simp_data[i].split()
127         ) < len(norm_data[i].split())):
128             allData.append([norm_data[i],
129             simp_data[i]])
130     allData = array(allData)
131
132     # Remove punctuation
133     if remove_punc:
134         print("Removing punctuations...")
135         allData[:,0] = [s.translate(str.maketrans(
136         ' ', ' ', string.punctuation)) for s in
137         allData[:,0]]
138         allData[:,1] = [s.translate(str.maketrans(
139         ' ', ' ', string.punctuation)) for s in
140         allData[:,1]]
141
142     # convert to lowercase
143     for i in range(len(allData)):
144         allData[i,0] = allData[i,0].lower().
145         strip()
146
147         allData[i,1] = allData[i,1].lower().
148         strip()
149         allData_aggr.append(allData[i,0])
150         allData_aggr.append(allData[i,1])
151
152     # empty lists
153     norm_l = []
154     simp_l = []

```



```

60 # populate the lists with sentence lengths
61 for i in allData[:,0]:
62     norm_l.append(len(i.split()))
63
64 for i in allData[:,1]:
65     simp_l.append(len(i.split()))
66
67 # prepare aggregate tokenizer
68 aggr_tokenizer = tokenization(allData_aggr)
69 aggr_vocab_size = len(aggr_tokenizer.
70     word_index) + 1
71 # prepare normal tokenizer
72 norm_tokenizer = tokenization(allData[:, 0])
73 norm_vocab_size = len(norm_tokenizer.
74     word_index) + 1
75 # prepare simple tokenizer
76 simp_tokenizer = tokenization(allData[:, 1])
77 simp_vocab_size = len(simp_tokenizer.
78     word_index) + 1
79
80 train, test = train_test_split(allData,
81     test_size=trainTestSplitPerc, random_state
82     = 12)
83
84 # prepare training data
85 trainX = encode_sequences(aggr_tokenizer,
86     aggr_length, train[:, 0])
87 trainY = encode_sequences(aggr_tokenizer,
88     aggr_length, train[:, 1])
89
90 # prepare validation data
91 testX = encode_sequences(aggr_tokenizer,
92     aggr_length, test[:, 0])
93 testY = encode_sequences(aggr_tokenizer,
94     aggr_length, test[:, 1])
95
96 model = build_model(aggr_vocab_size,
97     aggr_vocab_size, aggr_length, aggr_length,
98     hidden_size)
99 rms = optimizers.RMSprop(lr=0.001)
100 model.compile(optimizer=rms, loss='
101     sparse_categorical_crossentropy')
102
103 filename = 'NLPFinalModel'
104 checkpoint = ModelCheckpoint(filename,
105     monitor='val_loss', verbose=1,
106     save_best_only=True, mode='min')
107 start_time = time.time()
108 history = model.fit(trainX, trainY.reshape(
109     trainY.shape[0], trainY.shape[1], 1),
110     epochs=50, batch_size=batchSize,
111     validation_split = 0.2,
112     callbacks=[checkpoint], verbose=1)
113 print("--- %s seconds ---" % (time.time() -
114     start_time))
115
116 model.summary()
117 plt.plot(history.history['loss'])
118 plt.plot(history.history['val_loss'])
119 plt.legend(['train', 'validation'])
120 plt.xlabel('# epochs')
121 plt.ylabel('Loss')
122 plt.title('Model training')

```

```

107 plt.show()
108
109 model = load_model('NLPFinalModel')
110 preds = model.predict_classes(testX.reshape
111     ((testX.shape[0], testX.shape[1])))
112
113 # convert predictions into text (English)
114 preds_text = []
115 for i in preds:
116     temp = []
117     for j in range(len(i)):
118         t = get_word(i[j], aggr_tokenizer)
119         if j > 0:
120             if (t == get_word(i[j-1],
121                 aggr_tokenizer)) or (t == None):
122                 temp.append('')
123             else:
124                 temp.append(t)
125         else:
126             if (t == None):
127                 temp.append('')
128             else:
129                 temp.append(t)
130
131     preds_text.append(' '.join(temp))
132
133 print(calc_BLUE(pd.DataFrame({'actual' :
134     test[:,0], 'predicted' : preds_text})))

```

## I. Calculating the bleu score for the normal and simplified text file

```

1 #user variables
2 #simple wikipedia data
3 #normFilePath = "alignedSentencesNorm-0.3.txt"
4 #simpFilePath = "alignedSentencesSimp-0.3.txt"
5
6 #newsela corpus
7 normFilePath = "alignedSentencesNorm-0.3.txt"
8 simpFilePath = "alignedSentencesSimp-0.3.txt"
9
10 usingNewsela = True
11 remove_punc = True
12 remove_stopwords = False
13 word_cutoff = 15
14 aggr_length = 15
15 hidden_size = 256
16 batchSize = 64
17 trainTestSplitPerc= 0.3
18 RunAutoEncoder(usingNewsela, word_cutoff,
19     aggr_length, trainTestSplitPerc,
20     hidden_size, remove_punc, batchSize,
21     remove_stopwords, normFilePath,
22     simpFilePath)

```