

# Word Representation in Biomedical Domain

## Report of Group 3

## 1 Methodology and Underlying Principles

### 1.1 Parse the Data

Initially, we performed data extraction. Given the large scale of the COVID-19 dataset, we focused on extracting only the titles and abstracts of the papers to form our dataset. The specific extraction process is outlined in the code below.

```
1 import csv
2
3 # Open and read the metadata.csv file
4 with open(r"D:\DOWNLOAD\wetransfer_winter-school-nlp-2025_2025-02-03_1458\all_sources-metadata_2020-03-13.csv", encoding='utf-8') as f_in:
5     reader = csv.DictReader(f_in)
6
7     # Open a txt file for writing
8     with open('titles_abstracts.txt', 'w', encoding='utf-8') as f_out:
9         for row in reader:
10             title = row['title'].strip() if row['title'] else "-"
11             abstract = row['abstract'].strip() if row['abstract'] else "-"
12
13             # Combine title and abstract
14             text = f"{title}\n{abstract}\n\n"
15
16             # Write to file
17             f_out.write(text)
```

### 1.2 Tokenization

Our tokenization process is implemented in four methods, as described below:

1. Use split()
2. Use NLTK or SciSpaCy
3. Use Byte-Pair Encoding (BPE)
4. Build custom Byte-Pair Encoding (BPE)

#### 1.2.1 Use split()

In this approach, we perform tokenization by using the simple split() method together with regex, which splits the text based on spaces and newline characters. The implementation is as follows:

```
1 import re
2 import random
3 from collections import Counter
4
5 def Tokenization_1(content):
6     filename = './titles_abstracts.txt'
7     with open(filename, 'r', encoding='utf-8') as file:
8         content = file.read()
9         words = []
```

```

10     lines = content.split("\n")
11     for line in lines:
12         words.extend(re.split(r"\W+", line)) # Use regular
            expressions to split text
13     # Filter out empty strings to ensure only valid words are
        extracted
14     words = [word for word in words if word]
15
16     # Use the Counter class to count word occurrences and perform
        weighted sampling
17     word_counts = Counter(words)
18     weighted_sample = random.choices(
19         list(word_counts.keys()), weights=word_counts.values(), k=20
20     )
21     print("Weighted-sample:", "-".join(weighted_sample))
22
23     result = "-".join(words)
24     return result

```

### 1.2.2 Use NLTK

In this section, we use NLTK [1] for tokenization. NLTK is a library optimized for processing scientific texts. Its tokenization approach is tailored to handle the specific challenges posed by scientific and technical language. Below are the principles underlying NLTK's tokenization process:

- **Rule-Based Tokenization:** Utilizes predefined rules to handle punctuation, abbreviations, and special symbols common in scientific writing.
- **Regular Expressions:** Implements regex-based patterns to split text efficiently while preserving meaningful units.
- **Treebank Tokenizer:** Uses patterns derived from the Penn Treebank dataset to ensure consistency in breaking down complex sentences.
- **Word and Sentence Segmentation:** Differentiates between sentence boundaries and word boundaries, ensuring accurate parsing of multi-sentence documents.
- **Handling Special Cases:** Accounts for numbers, dates, chemical formulas, and technical notation, preventing incorrect segmentation.

The details of the implementation are as follows.

```

1  import nltk
2  from nltk import word_tokenize, sent_tokenize
3
4  def Tokenization_2(content):
5      # Sentence segmentation
6      sents = sent_tokenize(content)
7      interpuations = [",", ".", ":", ";", "?", "(", ")", "[", "]",
            ", &", "!", "*", "@", "#", "$", "%", "/", "\\", " "]
8
9      words = []
10     for sent in sents:
11         sent_t = word_tokenize(sent) # Tokenization
12         sent_t = [word.lower() for word in sent_t if word not in
            interpuations] # Remove punctuation marks
13         words.append(sent_t)
14
15     # Save tokenization results
16     result = ["-".join(sent) for sent in words]

```

```

17     with open("./nlkt_tokenized.txt", "w", encoding="utf-8") as f:
18         f.write("\n".join(result))
19
20     return result

```

### 1.2.3 Use Byte-Pair Encoding (BPE)

In this section we use Byte-Pair Encoding to tokenize. Byte-Pair Encoding (BPE) is a tokenization method that iteratively merges the most frequent adjacent character or subword pairs in a text until a predefined vocabulary size is reached. It starts by splitting words into individual characters and counting the frequency of adjacent pairs. The most common pair is merged into a new subword, and this process repeats, gradually forming longer tokens. By doing this, BPE efficiently reduces the size of the vocabulary while preserving the ability to handle rare words by breaking them into familiar subunits. This makes it especially useful for natural language processing, as it balances flexibility and efficiency in text representation.

The implementation is as follows:

```

1  from transformers import AutoTokenizer
2
3  def Tokenization_3(content):
4      bert_tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
5
6      # Split text into paragraphs to avoid excessively long inputs
7      max_length = 512
8      result = []
9
10     for i in range(0, len(content), max_length):
11         chunk = content[i : i + max_length]
12         tokens = bert_tokenizer.tokenize(chunk)
13         result.append(" ".join(tokens)) # Join tokens with spaces
14
15     # Save tokenization results
16     with open("bert_tokenized.txt", "w", encoding="utf-8") as f:
17         f.write("\n".join(result))
18
19     return result

```

### 1.2.4 Build custom Byte-Pair Encoding (BPE)

In this section we use build custom Byte-Pair Encoding (BPE). Custom Byte-Pair Encoding (BPE) follows the same iterative merging process as standard BPE but allows modifications for specific tasks or datasets. It starts with character-level tokenization, repeatedly merges the most frequent adjacent pairs into subwords, and continues until a target vocabulary size is reached. Custom implementations can adjust merge rules, handle special characters differently, or optimize tokenization for domain-specific applications, enhancing adaptability while preserving the efficiency of BPE in reducing vocabulary size and handling rare words.

The implementation is as follows:

```

1  from tokenizers import Tokenizer, models, trainers, pre_tokenizers
2
3  def Tokenization_4(content):
4      # Train the BPE model
5      print("Training BPE tokenizer...")
6      bpe_tokenizer = Tokenizer(models.BPE())
7      bpe_tokenizer.pre_tokenizer = pre_tokenizers.Whitespace()
8

```

```

9     trainer = trainers.BpeTrainer(special_tokens=["<unk>", "<pad>",
10         "<s>", "</s>"])
11     bpe_tokenizer.train(files=["titles_abstracts.txt"], trainer=
12         trainer)
13
14     # Save the trained BPE model
15     bpe_tokenizer.save("biomedical_bpe.json")
16
17     # Reload the BPE tokenizer and perform tokenization
18     bpe_tokenizer = Tokenizer.from_file("biomedical_bpe.json")
19     bpe_tokens = bpe_tokenizer.encode(content)
20
21     # Save BPE tokenization results
22     with open("bpe_tokenized.txt", "w", encoding="utf-8") as f:
23         f.write("-".join(bpe_tokens.tokens)) # Join tokens with
24         spaces
25
26     return bpe_tokens

```

## 1.3 Build Word Representations

### 1.3.1 Use N-gram Language Modeling

N-gram language modeling [2] represents words using statistical cooccurrence. Given a sequence  $w_1, w_2, \dots, w_t$ , the probability of a word based on the previous  $n - 1$  words is:

$$P(w_t | w_{t-(n-1)}, \dots, w_{t-1}) = \frac{C(w_{t-(n-1)}, \dots, w_t)}{C(w_{t-(n-1)}, \dots, w_{t-1})}$$

Embeddings are learned by constructing a co-occurrence matrix  $M$ , where:

$$M_{ij} = P(w_j | w_i)$$

Embeddings are obtained by singular value decomposition (SVD) of  $M$ :

$$M = USV^T$$

where  $U, V$  are low-dimensional word embeddings that capture semantic relationships.  $U$  captures the embedding of words by focusing on their co-occurrence patterns across contexts (rows of  $M$ ), while  $V$  captures the embedding of contexts by focusing on their co-occurrence patterns with words (columns of  $M$ ).

The details of the implementation are as follows.

```

1  import collections
2  from collections import defaultdict, Counter
3  from scipy.sparse import csr_matrix
4  from scipy.sparse.linalg import svds
5  import numpy as np
6
7  class NGramLanguageModel:
8      def __init__(self, n, vocab_size=10000, embedding_dim=100):
9          self.n = n
10         self.vocab_size = vocab_size
11         self.embedding_dim = embedding_dim
12
13         # Vocabulary
14         self.vocab = {'<UNK>', '<START>', '<END>'}
15         self.word_to_id = {}
16         self.id_to_word = {}
17

```

```

18     # N-gram statistics
19     self.ngram_counts = defaultdict(lambda: defaultdict(float))
20     self.context_counts = defaultdict(float)
21
22     # Word embeddings
23     self.word_embeddings = None
24     self.context_embeddings = None
25
26     # Smoothing parameter
27     self.alpha = 0.1
28
29 def build_vocab(self, text):
30     """Build vocabulary from text"""
31     words = text.split()
32     word_counts = Counter(words)
33     top_words = sorted(word_counts.items(), key=lambda x: x[1],
34                        reverse=True)
35     top_words = top_words[:self.vocab_size - 3] # Reserve space
36                                                  # for special tokens
37
38     self.vocab.update(word for word, _ in top_words)
39     self.word_to_id = {word: idx for idx, word in enumerate(self
40                      .vocab)}
41     self.id_to_word = {idx: word for word, idx in self.
42                      word_to_id.items()}
43
44 def build_cooccurrence_matrix(self, text):
45     """Build co-occurrence matrix"""
46     words = ['<START>'] + text.split() + ['<END>']
47     vocab_size = len(self.vocab)
48
49     # Initialize sparse matrix components
50     rows, cols, data = [], [], []
51
52     # Compute co-occurrence
53     for i in range(len(words) - self.n + 1):
54         context_words = words[i:i + self.n - 1]
55         target_word = words[i + self.n - 1]
56
57         # Get word IDs
58         context_ids = [self.word_to_id.get(w, self.word_to_id['<
59                        UNK>'])
60                        for w in context_words]
61         target_id = self.word_to_id.get(target_word, self.
62                        word_to_id['<UNK>'])
63
64         # Update n-gram statistics
65         context = tuple(context_ids)
66         self.ngram_counts[context][target_id] += 1
67         self.context_counts[context] += 1
68
69         # Update co-occurrence counts
70         for context_id in context_ids:
71             rows.append(context_id)
72             cols.append(target_id)
73             data.append(1.0)
74
75     # Construct sparse matrix

```

```

70     M = csr_matrix((data, (rows, cols)), shape=(vocab_size,
71                                     vocab_size))
72
73     # Compute conditional probability P(wj|wi)
74     row_sums = M.sum(axis=1).A.flatten()
75     row_sums[row_sums == 0] = 1 # Avoid division by zero
76     M = csr_matrix(M / row_sums[:, np.newaxis])
77
78     return M
79
80 def factorize_matrix(self, M):
81     """Perform SVD on the co-occurrence matrix"""
82     U, S, Vt = svds(M, k=self.embedding_dim)
83
84     # Scale singular values
85     S_sqrt = np.sqrt(S)
86     self.word_embeddings = U * S_sqrt # Scale each column of U
87     self.context_embeddings = Vt.T * S_sqrt # Scale each row of
88     Vt.T
89
90 def train(self, text):
91     """Train the model"""
92     self.build_vocab(text)
93     M = self.build_cooccurrence_matrix(text)
94     self.factorize_matrix(M)
95
96 def get_word_vector(self, word):
97     """Retrieve word vector"""
98     word_id = self.word_to_id.get(word, self.word_to_id['<UNK>'])
99     return self.word_embeddings[word_id]
100
101 def get_context_vector(self, context):
102     """Retrieve context vector"""
103     context_words = context.split()[-self.n + 1:]
104     context_ids = [self.word_to_id.get(w, self.word_to_id['<UNK>'])
105                    for w in context_words]
106     return np.mean([self.context_embeddings[i] for i in
107                     context_ids], axis=0)
108
109 def predict_next(self, context):
110     """Predict the next word"""
111     context_vec = self.get_context_vector(context)
112
113     # Compute similarity with all words
114     similarities = np.dot(self.word_embeddings, context_vec)
115
116     # Find the most similar word
117     best_id = np.argmax(similarities)
118     best_word = self.id_to_word[best_id]
119
120     # Compute probability (using softmax)
121     exp_similarities = np.exp(similarities - np.max(similarities))
122     probability = exp_similarities[best_id] / exp_similarities.sum()

```

```

121         return best_word, float(probability)
122
123     def generate(self, seed, length=10):
124         """Generate text"""
125         current = ['<START>'] + seed.split()
126         result = current.copy()
127
128         for _ in range(length):
129             context = '-'.join(current[-(self.n - 1):])
130             next_word, _ = self.predict_next(context)
131
132             if next_word == '<END>':
133                 break
134
135             result.append(next_word)
136             current = result[-self.n + 1:]
137
138         return '-'.join(result[1:]) # Remove <START> token
139
140     def get_similar_words(self, word, k=5):
141         """Find the top-k similar words"""
142         word_vec = self.get_word_vector(word)
143         similarities = np.dot(self.word_embeddings, word_vec)
144         top_k = np.argsort(similarities)[-k-1:-1][::-1] # Exclude
145                 the word itself
146         return [(self.id_to_word[i], similarities[i]) for i in top_k]
147
148     def Representations_1(result):
149         # Tokenized text input
150         tokenized_text = result
151
152         model = NGramLanguageModel(3)
153         model.train(tokenized_text)
154
155         # Output word vectors
156         for word, vector in model.get_word_vector().items():
157             print(f"Word: {word}, Vector: {vector[:5]}..." ) # Display
158                                     first 5 dimensions
159
160         return model

```

### 1.3.2 Use Skip-gram with Negative Sampling

Skip-gram with Negative Sampling (SGNS) [3] is an efficient method of learning word embeddings by training a neural network to predict context words given a target word, while also distinguishing real context words from randomly sampled negative words.

Given a word sequence  $w_1, w_2, \dots, w_T$ , the skip-gram model maximizes the probability of context words  $w_c$  within a window around a target word  $w_t$ :

$$P(w_c|w_t) = \sigma(v_c^T v_t)$$

where  $v_t$  and  $v_c$  are the embeddings of the target and context words, and  $\sigma(x) = \frac{1}{1+e^{-x}}$  is the sigmoid function.

Negative sampling replaces the full softmax with a binary classification task, optimizing:

$$\log \sigma(v_c^T v_t) + \sum_{n=1}^k \mathbb{E}_{w_n \sim P_n} \log \sigma(-v_n^T v_t)$$

where  $k$  negative samples  $w_n$  are drawn from a noise distribution. This approach significantly reduces computational cost while preserving meaningful semantic structure in the learned embeddings.

The details of our implementation are as follows:

```
1 import numpy as np
2 from gensim.models import Word2Vec
3 from gensim.utils import simple_preprocess
4
5 class SkipGramWithNegativeSampling:
6     def __init__(self, vocab_size, vector_dim, window_size,
7         negative_samples):
8         self.vocab_size = vocab_size
9         self.vector_dim = vector_dim
10        self.window_size = window_size
11        self.negative_samples = negative_samples
12        # Initialize word vectors randomly
13        self.word_vectors = np.random.uniform(-0.5, 0.5, (vocab_size
14            , vector_dim))
15        self.context_vectors = np.random.uniform(-0.5, 0.5, (
16            vocab_size, vector_dim))
17
18    def _sigmoid(self, x):
19        # Activation function
20        return 1 / (1 + np.exp(-x))
21
22    def _loss(self, central_vec, context_vec, negative_vecs):
23        # Loss function
24        positive_score = np.dot(context_vec, central_vec)
25        positive_loss = -np.log(self._sigmoid(positive_score))
26        negative_scores = np.dot(negative_vecs, central_vec)
27        negative_loss = -np.sum(np.log(self._sigmoid(-
28            negative_scores)))
29        return positive_loss + negative_loss
30
31    def _gradient_update(self, central_vec, context_vec,
32        negative_vecs, learning_rate):
33        # Positive sample gradient
34        positive_score = self._sigmoid(np.dot(context_vec,
35            central_vec))
36        positive_grad = (positive_score - 1) * context_vec
37        # Negative sample gradient
38        negative_scores = self._sigmoid(np.dot(negative_vecs,
39            central_vec))
40        negative_grad = np.sum(negative_scores[:, np.newaxis] *
41            negative_vecs, axis=0)
42        # Update vectors
43        central_vec -= learning_rate * (positive_grad +
44            negative_grad)
45        context_vec -= learning_rate * (positive_score - 1) *
46            central_vec
47        for neg_vec in negative_vecs:
48            neg_vec -= learning_rate * negative_scores * central_vec
49
50    def train(self, dataset, epochs, learning_rate):
51        # Train model using gradient descent
52        for epoch in range(epochs):
53            total_loss = 0
54            for sentence in dataset:
55                for i, central_word in enumerate(sentence):
56                    # Set window
```



```

47         start = max(0, i - self.window_size)
48         end = min(len(sentence), i + self.window_size +
49                     1)
49         # Get sentence
50         context_words = [sentence[j] for j in range(
51             start, end) if j != i]
51         for context_word in context_words:
52             # Positive sample
53             central_vec = self.word_vectors[central_word
54                                     ]
54             context_vec = self.context_vectors[
55                 context_word]
55             # Negative samples
56             negative_indices = random.sample(
57                 [idx for idx in range(self.vocab_size)
58                  if idx not in sentence],
59                 self.negative_samples,
60             )
60             negative_vecs = self.context_vectors[
61                 negative_indices]
61             # Calculate loss function
62             loss = self._loss(central_vec, context_vec,
63                               negative_vecs)
63             total_loss += loss
64             # Update gradients
65             self._gradient_update(central_vec,
66                                   context_vec, negative_vecs, learning_rate
67                                 )
66         print(f"Epoch-{epoch+1}/{epochs}, Loss: {total_loss}")
67
68     def get_word_vectors(self):
69         # Get word vectors
70         return self.word_vectors
71
72     def Representations_2(result):
73         # Preprocess sentences
74         tokens = [sent.split(' ') for sent in result]
75         # Set model parameters
76         model = Word2Vec(
77             sentences=tokens[:20000], # Input preprocessed sentences
78             vector_size=256,          # Embedding vector dimension
79             window=5,                 # Window size (consider up to 5
80                                     words before and after)
81             sg=1,                     # Skip-gram mode (sg=1 for Skip-
82                                     gram; sg=0 for CBOW)
83             negative=10,              # Number of negative samples (
84                                     number of negative samples per positive sample)
85             min_count=2,              # Ignore words that appear less
86                                     than 2 times
87             workers=4,                # Number of CPU cores to use
88             epochs=10                 # Number of training iterations
89         )
90         # Output word vectors
91         vocab = model.wv.index_to_key
92         for word in vocab:
93             vector = model.wv[word]
94             print(f"Word: {word}, Vector: {vector[:5]}..." ) # Only
95             display first 5 dimensions

```

### 1.3.3 Use Contextualised Word Representation by Masked Language Model (MLM)

Masked Language Model (MLM) [4] is a technique used to learn \*\*contextualized word representations\*\* by training a model to predict randomly masked words in a sentence based on their surrounding context. This approach enables embeddings to capture dynamic word meanings depending on context, unlike static word embeddings (e.g. Word2Vec).

Given an input sequence  $w_1, w_2, \dots, w_T$ , some tokens are randomly masked (e.g.  $w_m$ ), and the model learns to predict them:

$$P(w_m | w_1, \dots, w_{m-1}, [\text{MASK}], w_{m+1}, \dots, w_T)$$

A transformer-based model (e.g., BERT) processes the full sequence bidirectionally, generating deep contextualized word embeddings. The loss function is typically the \*\*cross-entropy loss\*\* over the masked tokens:

$$\mathcal{L} = - \sum_m \log P(w_m | \mathbf{h}_m)$$

where  $\mathbf{h}_m$  is the hidden representation of the masked token. This method enables embeddings to capture word sense variations and syntactic dependencies based on the given context.

In our implementation, we utilize BERT for embedding generation. Given the large number of parameters in BERT and the high computational cost of training, we employ LoRA [5] fine-tuning to optimize the training process. The details of our implementation are as follows:

```

1  from transformers import BertTokenizer, BertForMaskedLM, Trainer,
    TrainingArguments
2  from datasets import Dataset
3  import torch
4  from transformers import DataCollatorForLanguageModeling
5  from peft import get_peft_model, LoraConfig, TaskType
6
7  def Representations_3():
8      # Load BERT (Masked Language Model version)
9      tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
10     model = BertForMaskedLM.from_pretrained("bert-base-uncased")
11     # Configure LoRA
12     lora_config = LoraConfig(
13         r=4, # Low-rank adaptation dimension
14         lora_alpha=32, # LoRA alpha parameter
15         lora_dropout=0.1, # Dropout probability
16         target_modules=["query", "value"] # Adapt transformer's Q/V
            matrices
17     )
18     model = get_peft_model(model, lora_config)
19     total_params = sum(p.numel() for p in model.parameters()) # All
        parameters
20     trainable_params = sum(p.numel() for p in model.parameters() if
        p.requires_grad) # Only trainable parameters
21     print(f"Total-model-parameters:-{total_params:,}")
22     print(f"Trainable-parameters:-{trainable_params:,}")
23     print(f"Frozen-parameters:-{total_params-trainable_params:,}")
24     # Read text data
25     with open("titles_abstracts.txt", "r", encoding="utf-8") as f:
26         lines = [line.strip() for line in f.readlines() if line.
            strip()]
27     # Create Hugging Face dataset
28     dataset = Dataset.from_dict({"text": lines})

```

```

29 # Split dataset
30 dataset = dataset.train_test_split(test_size=0.1)
31 train_dataset = dataset["train"]
32 eval_dataset = dataset["test"]
33 # Preprocessing: Tokenization + MLM masking
34 def tokenize_function(examples):
35     return tokenizer(examples["text"], truncation=True, padding=
        "max_length", max_length=256)
36 train_dataset = train_dataset.map(tokenize_function, batched=
    True)
37 eval_dataset = eval_dataset.map(tokenize_function, batched=True)
38 # Data collator
39 data_collator = DataCollatorForLanguageModeling(
40     tokenizer=tokenizer,
41     mlm=True, # Perform MLM training
42     mlm_probability=0.15 # 15% of tokens are masked
43 )
44 # Training arguments
45 training_args = TrainingArguments(
46     output_dir="./bert_lora_finetuned",
47     per_device_train_batch_size=16,
48     per_device_eval_batch_size=16,
49     num_train_epochs=3,
50     save_strategy="epoch",
51     evaluation_strategy="epoch",
52     logging_dir="./logs",
53     logging_steps=500,
54 )
55 # Trainer
56 trainer = Trainer(
57     model=model,
58     args=training_args,
59     train_dataset=train_dataset,
60     eval_dataset=eval_dataset,
61     data_collator=data_collator,
62 )
63 # Start training
64 trainer.train()
65 # Set to evaluation mode
66 model.eval()
67 # Read text and perform tokenization
68 with open("titles_abstracts.txt", "r", encoding="utf-8") as f:
69     lines = f.readlines()[:500]
70 # Process line by line
71 all_embeddings = []
72 for line in lines:
73     inputs = tokenizer(line.strip(), return_tensors="pt",
        truncation=True, padding="max_length", max_length=256)
74     with torch.no_grad():
75         outputs = model(inputs) # Direct BERT call
76         embeddings = outputs.logits # Extract token predictions
77         all_embeddings.append(embeddings.squeeze(0))
78 # Print information
79 print(f"Total embeddings: {len(all_embeddings)}")
80 print(f"Shape of single embedding: {all_embeddings[0].shape}")
81 # Save token embeddings
82 torch.save(all_embeddings, "bert_lora_finetuned_embeddings.pt")
83 print("LoRA fine-tuned token embeddings have been saved!")

```

## 1.4 Explore the Word Representations

### 1.4.1 Visualise the word representations by t-SNE

In this section, we employ t-SNE to transform high-dimensional embeddings into a 2D visual representation, enabling effective embedding visualization.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.manifold import TSNE
4 from random import sample
5
6 kimi_BME={
7     "viruses": ["virus", "coronavirus", "covid-19", "influenza", "
8         2019-ncov", "sars-cov-2", "h5n1", "mers-cov", "norovirus", "
9         hiv-1", "h1n1"],
10    "diseases": ["infection", "disease", "pneumonia", "syndrome", "
11        mortality", "cancer", "pathogen", "illness"],
12    "proteins": ["protein", "proteins", "antibodies", "spike", "
13        glycoprotein", "capsid", "membrane", "antigen", "antigens", "
14        epitopes"],
15    "genes": ["gene", "genes", "genome", "rna", "dna", "receptor", "
16        genetic", "mutation", "mutations"],
17    "cells": ["cells", "cell", "host", "cellular", "blood", "serum", "
18        tissue", "antibody", "antibodies"],
19    "treatments": ["treatment", "therapy", "vaccine", "vaccines", "
20        antiviral", "care", "diagnosis", "diagnostic", "drug", "drugs
21        "],
22    "epidemiology": ["outbreak", "epidemic", "pandemic", "
23        transmission", "quarantine", "prevention", "public", "
24        surveillance", "population", "spread"],
25    "research": ["study", "analysis", "data", "research", "model", "
26        models", "sequencing", "clinical", "review", "findings"],
27    "symptoms": ["symptoms", "fever", "cough", "dyspnea", "fatigue", "
28        headache", "sore-throat", "shortness-of-breath"],
29    "methods": ["methods", "assay", "pcr", "rt-pcr", "sequencing", "
30        imaging", "tools", "test", "tests", "diagnostic"],
31    "public_health": ["health", "public", "hospital", "hospitals", "
32        healthcare", "prevention", "quarantine", "epidemiology", "
33        surveillance"],
34    "miscellaneous": ["cases", "patients", "human", "novel", "
35        background", "severity", "risk", "control", "management", "
36        response"]
37 }
38
39 def Exploration_1(model, BME=kimi_BME, max_words=1000, seed=1022):
40     random.seed(seed)
41     # Extract vocabulary and corresponding word vectors
42     BMEs = [word for words in list(BME.values()) for word in words]
43     vocab = model.wv.index_to_key
44     vocab_1 = list(filter(lambda x: x in BMEs, vocab))
45     vocab_2 = sample(vocab, max_words-len(vocab_1))
46
47     vocab = vocab_1 + vocab_2 # Extract vocabulary
48     word_vectors = model.wv[vocab] # Extract word vectors

```

```

32 # Limit the maximum number of words for visualization
33 if len(vocab) > max_words:
34     np.random.seed(seed)
35     indices = np.random.choice(len(vocab), max_words, replace=
        False)
36     vocab = [vocab[i] for i in indices]
37     word_vectors = word_vectors[indices]
38
39 # Use t-SNE for dimensionality reduction to 2D
40 tsne = TSNE(n_components=2, random_state=seed, perplexity=75,
        init='pca', n_iter=1000)
41 reduced_embeddings = tsne.fit_transform(word_vectors)
42
43 # Plot t-SNE visualization
44 plt.figure(figsize=(14, 10))
45 for i, word in enumerate(vocab):
46     x, y = reduced_embeddings[i]
47     plt.scatter(x, y, color='blue', s=10)
48     plt.text(x + 0.2, y + 0.2, word, fontsize=8, color='gray')
49 plt.title("t-SNE Visualization of Word Embeddings")
50 plt.show()
51
52 return None

```

#### 1.4.2 Visualise the Word Representations of Biomedical Entities by t-SNE

In this section, we selected the biomedical entities from the tokenized data, assigned colors based on their categories, and visualized them using t-SNE to observe their distribution in the high-dimensional space.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.manifold import TSNE
4 from random import sample
5
6 kimi_BME={
7     "viruses": ["virus", "coronavirus", "covid-19", "influenza", "
        2019-ncov", "sars-cov-2", "h5n1", "mers-cov", "norovirus", "
        hiv-1", "h1n1"],
8     "diseases": ["infection", "disease", "pneumonia", "syndrome", "
        mortality", "cancer", "pathogen", "illness"],
9     "proteins": ["protein", "proteins", "antibodies", "spike", "
        glycoprotein", "capsid", "membrane", "antigen", "antigens", "
        epitopes"],
10    "genes": ["gene", "genes", "genome", "rna", "dna", "receptor", "
        genetic", "mutation", "mutations"],
11    "cells": ["cells", "cell", "host", "cellular", "blood", "serum",
        "tissue", "antibody", "antibodies"],
12    "treatments": ["treatment", "therapy", "vaccine", "vaccines", "
        antiviral", "care", "diagnosis", "diagnostic", "drug", "drugs
        "],
13    "epidemiology": ["outbreak", "epidemic", "pandemic", "
        transmission", "quarantine", "prevention", "public", "
        surveillance", "population", "spread"],
14    "research": ["study", "analysis", "data", "research", "model", "
        models", "sequencing", "clinical", "review", "findings"],
15    "symptoms": ["symptoms", "fever", "cough", "dyspnea", "fatigue",
        "headache", "sore-throat", "shortness-of-breath"],

```

```

16     "methods": ["methods", "assay", "pcr", "rt-pcr", "sequencing", "
        imaging", "tools", "test", "tests", "diagnostic"],
17     "public_health": ["health", "public", "hospital", "hospitals", "
        healthcare", "prevention", "quarantine", "epidemiology", "
        surveillance"],
18     "miscellaneous": ["cases", "patients", "human", "novel", "
        background", "severity", "risk", "control", "management", "
        response"]
19 }
20
21 def Exploration_2(model, BME=kimi_BME, max_words=1000, seed=1022):
22     random.seed(seed)
23
24     # Extract vocabulary and corresponding word vectors
25     BMEs = [word for words in list(BME.values()) for word in words]
26     BME_classes = [c for c in list(BME.keys())]
27     vocab = model.wv.index_to_key
28     vocab_1 = list(filter(lambda x: x in BMEs, vocab))
29     vocab_2 = sample(vocab, max_words - len(vocab_1))
30
31     vocab = vocab_1 + vocab_2 # Extract vocabulary
32     word_vectors = model.wv[vocab] # Extract word vectors
33
34     # Limit the number of words to visualize
35     if len(vocab) > max_words:
36         np.random.seed(seed)
37         indices = np.random.choice(len(vocab), max_words, replace=
            False)
38         vocab = [vocab[i] for i in indices]
39         word_vectors = word_vectors[indices]
40
41     # Use t-SNE to reduce dimensionality to 2D
42     tsne = TSNE(n_components=2, random_state=seed, perplexity=75,
        init='pca', n_iter=1000)
43     reduced_embeddings = tsne.fit_transform(word_vectors)
44
45     # Plot the t-SNE visualization
46     colors = ['pink', 'blue', 'green', 'yellow', 'purple', 'magenta',
        'red', 'orange', 'cyan', 'lawngreen', 'olive', 'lightblue']
47     colors = dict(zip(BME_classes, colors[:len(BME_classes)]))
48
49     plt.figure(figsize=(14, 10))
50     l = {} # Store data for the legend
51
52     for i, word in enumerate(vocab):
53         x, y = reduced_embeddings[i]
54         for c in BME_classes:
55             if word in BME[c]:
56                 color = colors[c]
57                 size = 30
58                 l[c] = (x, y, color)
59                 break
60             else:
61                 color = 'gray'
62                 size = 5
63         plt.scatter(x, y, color=color, s=size)
64         plt.text(x + 0.2, y + 0.2, word, fontsize=0.3 * size, color=
            color)

```

```

65
66     for key, value in l.items():
67         plt.scatter(value[0], value[1], color=value[2], s=size,
68                     label=key)
69
70     plt.title("t-SNE Visualization of Word Embeddings")
71     plt.legend()
72     plt.show()
73     return None

```

### 1.4.3 Co-occurrence

In this section, we demonstrate the biomedical entities which frequently cooccur with "coronavirus". The words are selected in order from highest to lowest according to the statistical frequency of co-occurrence with the target words in the corpus.

```

1  import pandas as pd
2
3  def Exploration_3(model, target_word='coronavirus', k=10,
4                  output_file='\\cooccurrence-words.xlsx'):
5      if target_word in model.vocab:
6          # Get similar words
7          cooccurrence_words = model.get_similar_words(target_word, k)
8          print("Top-10 similar words:")
9          for word, similarity in cooccurrence_words:
10             print('{0}: {1:.5f}'.format(word, similarity))
11
12         # Save to Excel
13         df = pd.DataFrame(cooccurrence_words, columns=['Word', 'Similarity'])
14         df.to_excel(output_file, index=False)
15         print(f"Co-occurrence words saved to {output_file}")
16         return True
17     else:
18         print(f"{target_word} not in vocabulary.")
19         return False

```

### 1.4.4 Semantic Similarity

In this section, we compare the cosine distances of the embeddings obtained from word2vec in the high-dimensional vector space to identify the tokenized results most semantically similar to "coronavirus".

```

1  from sklearn.metrics.pairwise import cosine_similarity
2  import pandas as pd
3
4  def Exploration_4(model, target_word = "coronavirus", output_file=
5                  '\\semantic-similar-words.xlsx'):
6      if target_word.lower() in model.wv:
7          target_vector = model.wv[target_word.lower()]
8          all_words = model.wv.index_to_key
9          similarities = []
10
11         for word in all_words:
12             if word != target_word.lower():
13                 similarity = cosine_similarity([target_vector], [
14                     model.wv[word]])[0][0]

```

```

13         similarities.append((word, similarity))
14
15     sorted_similarities = sorted(similarities, key=lambda x: x
16                                  [1], reverse=True)
17     top_similar_words = sorted_similarities[:10]
18
19     # Save to Excel
20     df = pd.DataFrame(top_similar_words, columns=['Word', '
21                  Similarity'])
22     df.to_excel(output_file, index=False)
23     print(f"Top-10 similar words saved to {output_file}")
24
25     print("Top-10 similar words:")
26     for word, similarity in top_similar_words:
27         print(' {0}: {1:.5f}'.format(word, similarity))
28     return True
29
30 else:
31     print(f"{target_word.lower()} not in vocabulary.")
32
33     return False

```

## 2 Results and Analysis

### 2.1 Tokenization

#### 2.1.1 Use split()

---

**Algorithm 1** Use split() for Tokenization

---

```

1: Input: Text file
2: Output: List of sampled words
3: Initialize empty list words
4: for each line in the file do
5:     Tokenize line using regex r"\W+"
6:     Append tokens to words
7: end for
8: Remove empty strings from words
9: Compute word frequencies using Counter(words)
10: Set sample_size to 20
11: Perform weighted sampling using word frequencies
12: Return: List of sampled words

```

---

- Tokenization and weighted sampling procedure:
  - Reads file line-by-line, tokenizing using regex `r"\W+"`
  - Filters out empty strings from **words** list
- Frequency-based weighted sampling:
  - Counts word occurrences with `Counter(words)`
  - Uses `random.choices()` for weighted random sampling based on frequency
  - `weights=word_counts.values()` ensures higher frequency words have greater probability of selection
- Random selection:
  - Selects 20 words using weighted probability distribution



- Resulting words are printed as output

source	for	acetate	number	rate
enzymes	of	respiratory	peptide	fruit
disease	virus	MERS	the	with
investigation	that	method	in	activation

Table 1: Weighted Random Sample of 20 Words

### 2.1.2 Use NLTK

---

#### Algorithm 2 Tokenization with NLTK

---

```

1: Input: content of file
2: Output: words (List of tokenized sentences)
3: Split content into sentences using sent_tokenize()
4: Define list of punctuation marks interpunctuations
5: Initialize empty list words
6: for each sentence in sentences do
7:   Tokenize sentence into words using word_tokenize()
8:   Remove punctuation words from tokenized sentence
9:   Convert all words to lowercase
10:  Append processed words to words
11: end for
12: Return: words

```

---

- Sentence segmentation:
  - Use `sent_tokenize()` to split input text into sentences.
  - Each sentence is processed separately.
- Tokenization:
  - For each sentence, `word_tokenize()` is used to split the sentence into words.
  - Punctuation marks are filtered out from tokenized words.
  - All words are converted to lowercase for uniformity.
- Output:
  - Returns a list of sentences, each containing a list of tokenized words (with punctuation removed).

covid-19	sara-cov-2	2019-nCoV	infection	pneumonia
transmission	symptoms	outbreak	wuhan	china
case	virus	fever	cough	diagnosis
treatment	epidemic	prevention	control	clinical

Table 2: Top 20 Representative High-Frequency Words of Tokenized Result of NLTK Tokenizer

### 2.1.1.3 Use Byte-Pair Encoding (BPE)

---

**Algorithm 3** Tokenization with usage of Byte-Pair Encoding (BPE)

---

```
1: Input: content of file
2: Output: tokens_list (List of tokenized text chunks)
3: Load pre-trained BERT tokenizer AutoTokenizer.from_pretrained("bert-base-uncased")
4: Set max_length to 512
5: Initialize empty list tokens_list
6: for each chunk of text with length max_length do
7:   Tokenize chunk using BERT tokenizer bert_tokenizer.tokenize(chunk)
8:   Append tokenized chunk to tokens_list
9: end for
10: Write tokens_list to file "bert_tokenized.txt"
11: Return: tokens_list
```

---

- BERT Tokenizer:
  - Load the pre-trained BERT tokenizer `AutoTokenizer.from_pretrained("bert-base-uncased")`.
  - Tokenizer converts raw text into tokens suitable for BERT-based models.
- Text Chunking:
  - Split the input text into chunks of size `max_length` (512 tokens), ensuring that the input size does not exceed BERT's maximum token limit.
- Tokenization:
  - For each chunk, tokenize using BERT tokenizer.
  - Tokens are joined into strings using spaces and appended to `tokens_list`.
- File Output:
  - Save the tokenized text in a file "`bert_tokenized.txt`" for later use.
- Output:
  - Return the list `tokens_list` containing tokenized text chunks.

COVID-19	SARS-CoV-2	2019-nCoV	virus	infection
pneumonia	transmission	symptoms	patients	diagnosis
treatment	research	outbreak	cases	China
Wuhan	genome	epidemiology	clinical	prevention

Table 3: Top 20 Representative High-Frequency Words of Tokenized Result of BERT Tokenizer

### 2.1.1.4 Build new Byte-Pair Encoding (BPE)

---

**Algorithm 4** Tokenization using BPE (Byte Pair Encoding)

---

```
1: Input: content of file
2: Output: result of tokenization
3: Initialize BPE model
4: Set up the trainer for the model
5: Training BPE model with content
6: Save the trained BPE model
7: Tokenize content using trained BPE model (as Track 2.3)
8: Write result to file
9: Return: result of tokenization
```

---

- Initialization of the BPE tokenizer:
  - Load a BPE tokenizer using `Tokenizer(models.BPE())`.
  - Pre-tokenization is set to whitespace-based, using `pre_tokenizers.Whitespace()`.
- BPE Model Training:
  - Define special tokens like "<unk>", "<pad>", "<s>", "</s>" for unknown, padding, start, and end tokens.
  - Train the tokenizer using the `BpeTrainer` with input data from "titles\_abstracts.txt".
- Model Saving:
  - Save the trained BPE tokenizer to a file named "biomedical\_bpe.json" for future use.
- Tokenization:
  - Reload the saved BPE tokenizer and use it to tokenize the input `content`.
  - Output tokens are joined into a space-separated string.
- Return:
  - The tokenized result is written to "bpe\_tokenized.txt" and returned as a list of tokens.

corona	virus	infection	COVID	pneumonia
epidemic	transmission	cases	patients	symptoms
treatment	control	quarantine	diagnosis	outbreak
SARS	MERS	vaccine	mortality	epidemiology

Table 4: Top 20 Representative High-Frequency Words of Tokenized Result of BERT Tokenizer

## 2.2 Explore the Word Representations

### 2.2.1 Visualise the Word Representations by t-SNE

#### Code Analysis

---

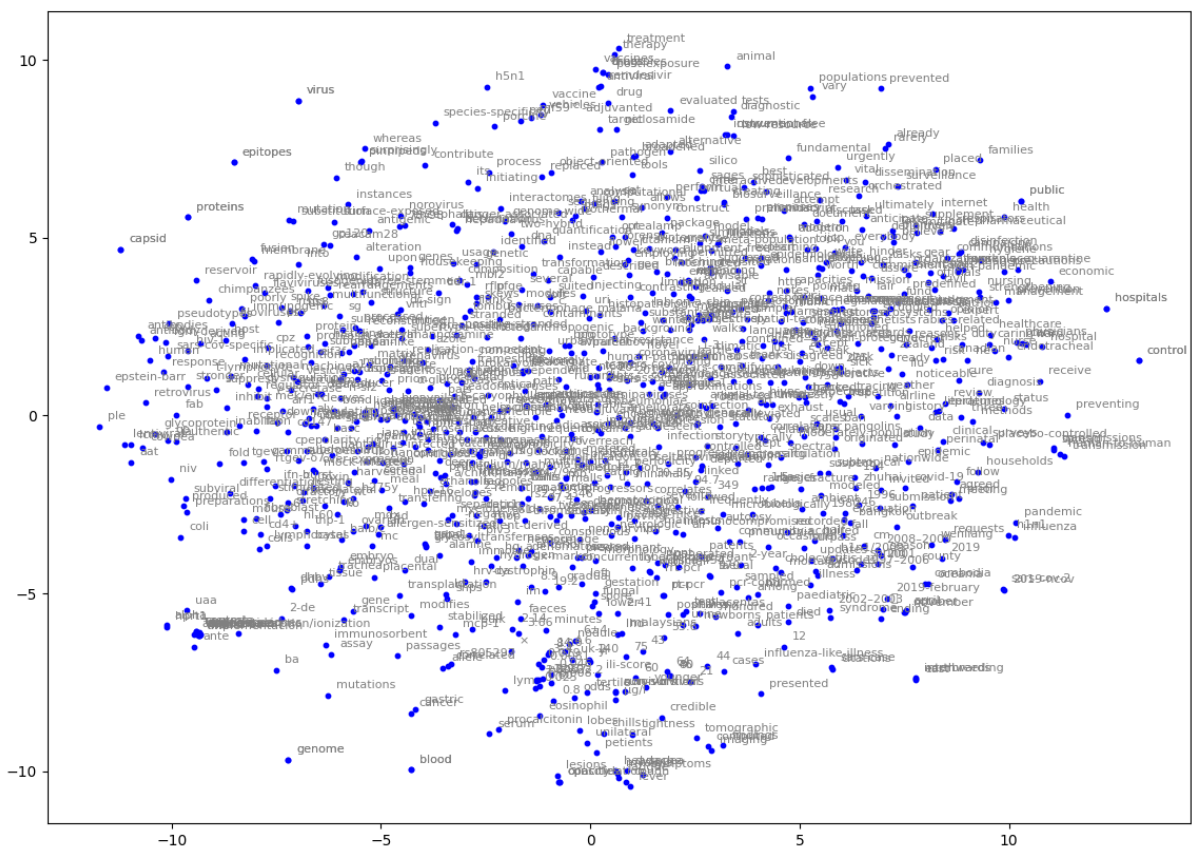
#### Algorithm 5 t-SNE Word Embedding Visualization

---

- 1: **Input:** Word embedding model *model*, Maximum words *max\_words*
  - 2: **Output:** 2D t-SNE visualization of word embeddings
  - 3: Extract vocabulary *vocab* and word vectors *word\_vectors* from *model*
  - 4: **if**  $|vocab| > max\_words$  **then**
  - 5:     Randomly sample *max\_words* words and their vectors
  - 6: **end if**
  - 7: Apply t-SNE to reduce *word\_vectors* to 2D
  - 8: Initialize plot for visualization
  - 9: **for all** word in *vocab* **do**
  - 10:     Plot word at its reduced 2D coordinates
  - 11: **end for**
  - 12: Display the visualization
- 

- Random sampling is used to control visualization density. Only `max_words=1000` are randomly selected from the vocabulary to ensure the plot remains interpretable without overcrowding.
- t-SNE parameters are carefully chosen to optimize for the preservation of local structure:
  - `perplexity=75` strikes a balance between capturing local relationships (similar words) and global relationships (overall structure).

- The algorithm starts with random initialization, which encourages a broad exploration of the word embedding space before settling into a meaningful visualization.
- t-SNE runs for 1000 iterations to ensure sufficient convergence for accurate 2D representation of the word vectors.
- The visualization maintains a uniform and consistent style to ensure clarity:
  - Words are represented by blue dots with a size of `s=10`, making the points distinct without overlapping too much.
  - Text labels are offset by (0.2, 0.2) from the corresponding points and colored gray to avoid visual clutter and overlap.
  - The figure size is fixed at 14x10 to maintain an appropriate aspect ratio and visualization scale for better readability.



## Spatial Topology

### Embedding Dynamics (Enabled by Code Parameters)

### Semantic Clustering Insights

- **Proximity between words** suggests conceptual similarity
- **Scattered blue dots** reveal nuanced semantic networks
- Demonstrates how **word2vec** models capture intricate linguistic contextual relationships

## 2.2.2 Visualise the Word Representations of Biomedical Entities by t-SNE

### Code Analysis

---

**Algorithm 6** t-SNE Word Embedding Visualization

---

```
1: Input: model (word embedding model), vocabulary, max_words
2: Output: 2D scatter plot of word embeddings
3: words ← random.sample(vocabulary, max_words) if len(vocabulary) > max_words else
   vocabulary
4: vectors ← [model[word] for word in words]
5: embedded_words ← TSNE(n_components=2, perplexity=30, random_state=42).fit_transform(vectors)
6: Initialize 2D scatter plot
7: for all word, (x, y) in zip(words, embedded_words) do
8:     Plot point at (x, y)
9:     Annotate point with corresponding word
10: end for
11: Show the scatter plot
```

---

- Hybrid vocabulary selection: - BME disease entities (vocab.1) - Random 300-word sample (vocab.2)
- Visual differentiation: - Disease terms: Larger points (30px), category-specific colors - General terms: Small gray points (5px) - Font size proportional to point size ( $0.3\times$ )
- Color mapping: - Distinct colors for each disease class - Gray for non-classified terms

This implementation enables clear visualization of disease term relationships while maintaining context within the broader vocabulary space.

### Visualization Results

#### Semantic Proximity Insights

- **Viral research terms** densely interconnected
- **Close spatial relationships** between epidemiological and clinical concepts
- **Clear differentiation** between molecular, clinical, and systemic research domains

#### Contextual Categorizations

- **Central region:** Generalized research methodologies
- **Peripheral regions:** Specialized technical terminology

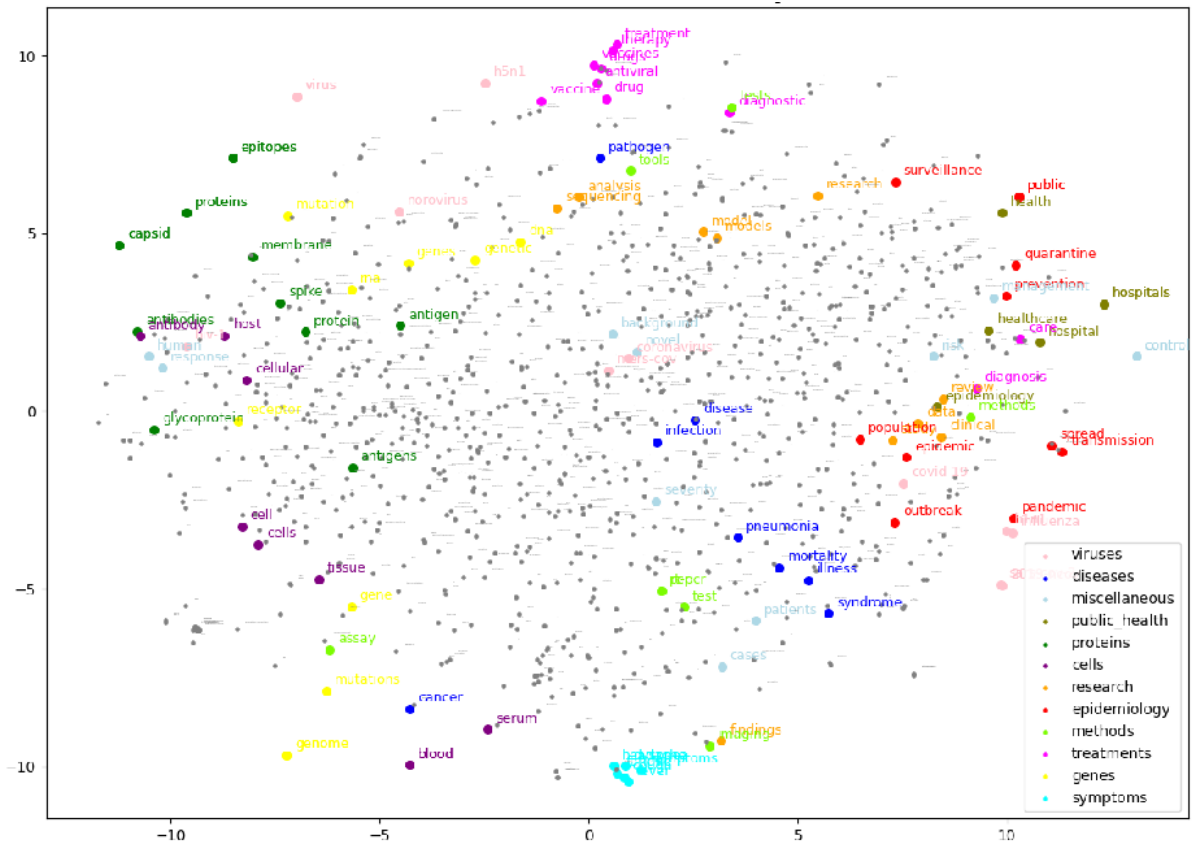


Figure 2: t-SNE Visualization of the Word Representations of Biomedical Entities

viruses	diseases
proteins	genes
cells	treatments
epidemiology	research
symptoms	methods
public health	miscellaneous

### 2.2.3 Co-occurrence

---

#### Algorithm 7 Co-occurrence Word Similarity Analysis

---

```

1: function EXPLORATION_3(model, target_word, k, output_file)
2:   if target_word in model.vocabulary then
3:     cooccurrence_words ← model.get_similar_words(target_word, k)
4:     Print "Top 10 similar words:" (word, similarity) in cooccurrence_words
5:     Print word, similarity
6:
7:     df ← pd.DataFrame(cooccurrence_words, columns=['Word', 'Similarity'])
8:     df.to_excel(output_file, index=False)
9:     Print "Co-occurrence words saved to " + output_file
10:    return True
11:   end if
12:   Print target_word + " not in vocabulary."
13:   return False
14: end function

```

---

Cross-Domain Relevance:

- The list contains terms from virology, epidemiology, and broader medical contexts, reflecting interdisciplinary discussions surrounding "coronavirus".

#### Contextual Associations:

- Many words relate to either the direct effects of COVID-19 (e.g. "persistent," "obstructive") or comparative studies involving other pathogens (e.g., "flaviviral," "SFTSV").

#### Broad Impact of COVID-19:

- The inclusion of terms like "Alzheimer" demonstrates how COVID-19 research spans beyond respiratory symptoms, touching neurological and systemic effects.

Co-occurring Word	Frequency
Covid	0.009722719
multiplicity	0.006220484
flaviviral	0.005976294
Lyme	0.005395208
productive	0.005248906
obstructive	0.005074524
mouth	0.004548549
persistent	0.004354533
Alzheimer	0.004330223
SFTSV	0.004047467

Table 5: 10 biomedical entities with the highest frequency of co-occurrence with coronavirus

## 2.2.4 Semantic Similarity

---

### Algorithm 8 Semantic Similarity Analysis

---

```

1: function EXPLORATION_4(model, target_word, output_file)
2:   if target_word.lower() in model.vocabulary then
3:     target_vector ← model.wv[target_word.lower()]
4:     similarities ← [(w, cosine_similarity([target_vector], [model.wv[w]])[0][0])
5:                   for w in model.wv.index_to_key if w ≠ target_word.lower()]
6:     top_similar_words ← sorted(similarities, key=lambda x: x[1], reverse=True)[:10]
7:     pd.DataFrame(top_similar_words, columns=['Word', 'Similarity']).to_excel(output_file, index=False)
8:     Print "Top similar words saved to " + output_file
9:     Print top_similar_words
10:    return True
11:  end if
12:  Print "Target word not in vocabulary."
13:  return False
14: end function

```

---

#### Model Strengths:

- The skip-gram model successfully captures semantic relationships, synonyms, and abbreviations (e.g., "2019-novel," "ncov-2019").
- Contextual and domain-specific terms like "provisionally" and "abstract" are included, showcasing the model's sensitivity to real-world usage.

### Error Handling:

- Recognition of typographical variants like "coronavirus" reflects the model's robustness against minor text inconsistencies.

### Temporal and Specific References:

- Terms like "2019" and "2019-ncov" highlight the temporal specificity of the dataset and the focus on SARS-CoV-2-related contexts.

Semantic Similar Word	Similarity
novel	0.6320
2019-novel	0.6185
coronavirus	0.6100
abstract	0.5961
provisionally	0.5929
2019-ncov	0.5895
2019	0.5876
ncov-2019	0.5767
cov	0.5766
2019-novel	0.6185

Table 6: 10 biomedical entities with the highest semantic similarity with coronavirus

## 3 Reflection

### 3.1 Strengths

#### 1. Comprehensive Data Handling

- The data preprocessing section is well-designed, efficiently handling JSON parsing and extraction.
- Alternative strategies are considered for handling different data availability scenarios, improving the robustness of the approach.

#### 2. Utilization of Standard Libraries

- The notebook makes effective use of widely adopted NLP libraries, ensuring compatibility, scalability, and performance.
- The structured approach to data extraction enhances readability and maintainability.

#### 3. Reproducibility and Environment Setup

- Clear instructions are provided for setting up the environment, including system requirements and data acquisition procedures.
- The use of standardized commands facilitates easy dataset retrieval and preparation.

### 3.2 Weaknesses

#### 1. Lack of Robust Exception Handling

- Error handling mechanisms are not consistently applied, which may lead to unexpected failures during file operations or data processing.
- Improved exception handling would enhance fault tolerance and user experience.



## 2. Potential Inefficiencies in Large-Scale Processing

- The notebook processes a large dataset, but there is no explicit mention of optimization techniques such as parallel processing or batch operations.
- Leveraging more efficient computation strategies could significantly improve performance.

## 3. Limited Data Visualization

- While the focus is on text processing, the absence of exploratory data analysis (EDA) limits insights into data distribution and characteristics.
- Integrating visualization techniques could provide a more comprehensive understanding of the dataset.

## 3.3 Overall Evaluation

- **Strengths:** Well-structured framework, effective data preprocessing, appropriate library usage, and clear setup instructions.
- **Weaknesses:** Limited portability, lack of robust error handling, inefficiencies in large-scale processing, and insufficient data visualization.

## References

- [1] S. Bird, “Nltk: The natural language toolkit,” in *Annual Meeting of the Association for Computational Linguistics*, 2006. [Online]. Available: <https://api.semanticscholar.org/CorpusID:1438450>
- [2] P. F. Brown, P. V. deSouza, R. L. Mercer, V. J. D. Pietra, and J. C. Lai, “Class-based n-gram models of natural language,” *Comput. Linguist.*, vol. 18, no. 4, p. 467–479, Dec. 1992.
- [3] A. Fonarev, O. Hrinchuk, G. Gusev, P. Serdyukov, and I. Oseledets, “Riemannian optimization for skip-gram negative sampling,” *ArXiv*, vol. abs/1704.08059, 2017.
- [4] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” in *North American Chapter of the Association for Computational Linguistics*, 2019.
- [5] J. E. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, and W. Chen, “Lora: Low-rank adaptation of large language models,” *ArXiv*, vol. abs/2106.09685, 2021.