

Introduction:

The evidence-based medicine (EBM) approach emphasizes integrating the latest and most reliable evidence, such as results from randomized controlled trials (RCTs), into the decision-making process for patient care. By combining this evidence with the clinician's expertise and patient-specific considerations, EBM aims to improve patient outcomes and deliver higher-quality healthcare. In practice, successful applications of EBM often involve addressing clinical questions through the analysis of extensive medical literature databases like PubMed. Typically, the PICO framework is employed to develop clear, focused clinical questions by breaking them down into four key components: Participants/Problem (P), Intervention (I), Comparison (C), and Outcome (O).

Problem Statement:

Applications of evidence-based medicine (EBM) that are successful must correctly respond to clinical queries that are obtained from large databases of medical literature, like PubMed. PICO, which divides the question into four parts: Participants/Problem (P), Intervention (I), Comparison (C), and Outcome (O), is a popular framework for creating focused clinical inquiries. But a large amount of medical literature lacks clearly labelled PICO features and is not optimally structured. This deficiency makes it more difficult to choose pertinent research and makes it hard to automate the literature analysis process. To make it easier to automatically choose relevant publications for clinical research, the project attempts to build a mechanism for automatically detecting PICO components in medical abstracts.

We propose categorizing phrases in medical literature based on their PICO components using a Long Short-Term Memory (LSTM) neural network model. By tackling both the inefficiencies and the poor performance of current systems in PICO detection, the methodology aims to improve upon earlier techniques, which frequently depend on simpler models and do not take advantage of the sequential structure of phrase categorization.

1.1. Dataset:

The [dataset](#) consists of clinical trial data gathered from randomized controlled trials (RCTs). It includes patient demographics, clinical histories, treatment regimens, and outcomes from various medical interventions. This rich set of information provides critical insights into patient responses, treatment efficacy, and potential side effects, making it ideal for evidence-based decision-making.

- **Dataset Overview:**

This project's dataset, which derives from randomized controlled trials (RCTs), has extensive data on patient demographics, clinical histories, therapeutic modalities, and ~~thi~~ outcomes. Evidence-based clinical decision-making relies heavily on comprehensive understanding of patient responses and the effectiveness of different medical therapies.

Each clinical trial in the dataset is labeled using the PICO framework, which divides clinical data into four fundamental components: Participants/Problem (P), Intervention (I), Comparison (C), and Outcome (O). Through the identification of key components in full-text papers and clinical abstracts, the research hopes to speed up the otherwise labor-intensive process of quickly retrieving pertinent information.

- **Data Dictionary:**

Field Name	Data Type	Description
study_id	Integer	A unique identifier for each study in the dataset (e.g., 28628768, 28779643).
objective	Text	The objective of the study, specifying the primary goal or aim of the research.
objective_label	String	The PICO label for the objective, typically marked as A (Aim).
design	Text	The design of the study, detailing the methodology used (e.g., randomized controlled trial).
setting	Text	The setting in which the study took place (e.g., online, metropolitan maternity units).

population	Text	Description of the population or sample involved in the study, including specific criteria (e.g., self-referred pregnant women).
population_label	String	The PICO label for the population, typically marked as P (Participants/Problem).
methods	Text	Description of the methods used in the study, including sample size and treatment conditions.
methods_label	String	The PICO label for methods, typically marked as M (Methods).
main_outcome_measures	Text	The primary outcome measures assessed in the study, including specific scales or assessments used.
outcome_label	String	The PICO label for outcomes, typically marked as O (Outcomes).
results	Text	Key findings from the study, summarizing significant results and statistical outcomes.
results_label	String	The PICO label for results, typically marked as R (Results).
limitations	Text	A brief description of the study's limitations or challenges encountered during the research.
limitations_label	String	A label for limitations, marked as Others.

conclusion	Text	Summary of the study's conclusions based on the findings and implications for practice.
conclusion_label	String	The PICO label for conclusions, typically marked as C (Conclusion).
implications_for_practice	Text	Insights or recommendations for practice derived from the study's findings.
implications_label	String	The PICO label for implications, typically marked as C (Conclusion) as it relates to practice.

- **Data Dependency**

Data dependency in this study refers to the reliance on a curated dataset from MEDLINE, specifically 489,026 abstracts extracted from PubMed with specific search criteria (availability of abstracts, English language, and randomized controlled trials). The extraction process involved automatic annotation of sentence categories, focusing mainly on detecting P (Problem), I (Intervention), and O (Outcome) labels. Additionally, other sentence types were categorized into AIM, METHOD, RESULTS, and CONCLUSION to enhance the optimization method. In cases where section headings overlapped with multiple categories, multi-label annotations were applied. Only abstracts containing P, I, or O labels (totaling 24,668 abstracts) were primarily analyzed to evaluate detection accuracy.

1.2.Dataset Collection and Preprocessing:

The Data Preprocessing step involves preparing and structuring the raw dataset by cleaning text, standardizing tags, and organizing data into sentences and labels.

1. **Word Processing:** Each word in the dataset undergoes a basic cleaning process:

- o **Punctuation Removal:** Unwanted punctuation and specific characters, such as the '@' symbol, are removed from the text to reduce noise.
- o **Whitespace Trimming:** Any leading or trailing whitespace is stripped from each word.

- o **Lowercasing:** All words are converted to lowercase, ensuring uniformity and reducing dimensionality.

```
def process_word(word):
    """Basic processing for each word, removing unwanted characters and blank spaces, then lowercasing."""
    # Remove punctuation and specific characters, then strip any surrounding whitespace
    word = re.sub(r'^\w\s', '', word) # Remove punctuation
    word = word.replace('@', '') # Specifically remove '@' symbol if present
    word = word.strip() # Remove surrounding whitespace
    return word.lower()
```

This processing prepares words for tokenization, maintaining only essential characters for model training.

2. **Tag Standardization:** Tags representing PICO elements are standardized to a consistent format:

- o A mapping dictionary is used to convert short tags (e.g., 'A', 'M', 'P', 'I', 'O', 'R', 'C') to their full forms, such as **Aim, Method, Participants, Intervention, Outcome, Results, and Conclusion**.
- o This ensures that each PICO element is represented clearly, making the dataset more interpretable for modeling and analysis.

```
def process_tag(tag):
    """Basic processing for tags (if needed, e.g., mapping specific tags)."""
    tag_map = {'A': 'Aim', 'M': 'Method', 'P': 'Participants', 'I': 'Intervention', 'O': 'Outcome', 'R': 'Results', 'C': 'Conclusion'}
    return tag_map.get(tag, tag)
```

3. **Dataset Structuring:** The dataset is structured to support sequence modeling:

- o **Sentence and Tag Accumulation:** Each section of sentences and their respective tags is stored as a pair, with sentences segmented based on predefined markers.
- o **Sentence Processing:** For each sentence in the dataset, only valid, cleaned words are retained, resulting in a list of sentences accompanied by their corresponding PICO tags.
- o **Iteration Limit:** If specified, the dataset can limit the number of processed sections (sentences and tags) for controlled processing, useful for preliminary analysis or testing.

```

def Dataset(filename, max_iter=None):
    results = []
    with open(filename) as f:
        sentences, tags = [], []
        n_iter = 0
        for line in f:
            line = line.strip()
            if not line:
                # End of a section; store accumulated sentences and tags
                if sentences:
                    n_iter += 1
                    if max_iter is not None and n_iter > max_iter:
                        break
                    results.append((sentences, tags))
                    sentences, tags = [], []
            elif not line.startswith("###"):
                # Process lines with PICO elements
                ls = line.split('|')
                tag, sentence = ls[1], ls[2].split()
                # Clean each word in the sentence and filter out empty words
                sentence = [process_word(word) for word in sentence if process_word(word)]
                tag = process_tag(tag)
                sentences.append(sentence)
                tags.append(tag)
    return results

```

This preprocessing approach ensures the text data is clean, standardized, and organized, facilitating efficient modeling of PICO elements and enabling accurate predictions.

2. Exploratory Data Analysis:

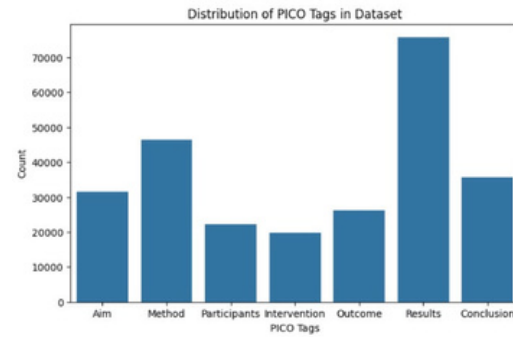
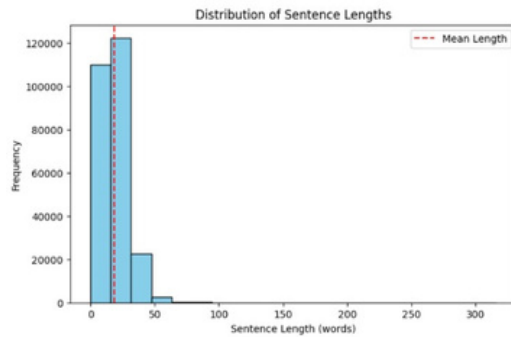
The Exploratory Data Analysis provides insights into the structure, distribution, and relationships within the PICO-tagged dataset. Key visualizations and findings are summarized below:

1. Distribution of PICO Tags in Dataset:

- o The dataset exhibits varying frequencies across the different PICO elements. The **Results** tag is the most prevalent, followed by **Method** and **Conclusion**, indicating a detailed focus on findings and methodologies in the abstracts.
- o Tags such as **Participants** and **Intervention** appear less frequently, which may imply limited elaboration on these components in the medical literature or potentially less frequent explicit tagging.

2. Distribution of Sentence Lengths:

- o Sentence length distributions reveal that most sentences are concise, with lengths peaking around 10-20 words. The distribution is right-skewed, with a mean sentence length denoted by a red dashed line.
- o This brevity aligns with the structured nature of medical abstracts, where information is often presented concisely.

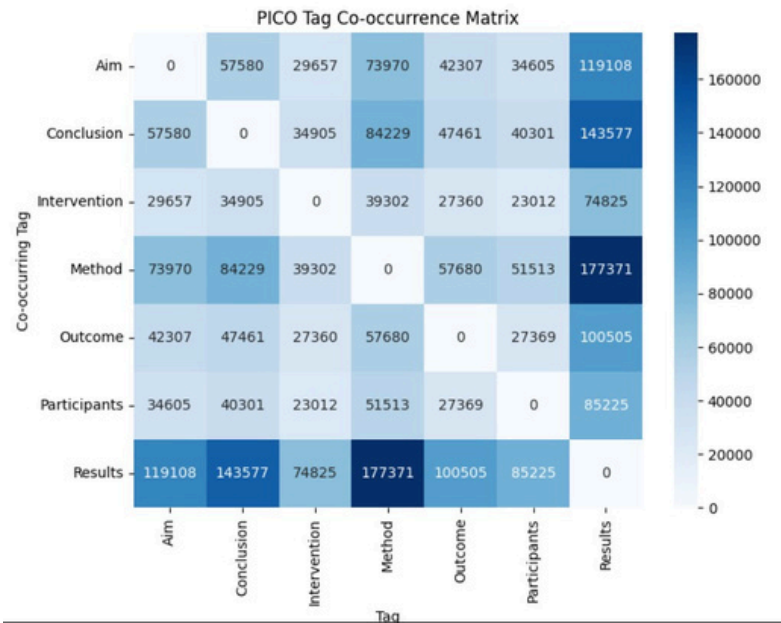


3. Top 20 Most Frequent Words in Dataset:

- o The most common words include standard stopwords such as "the," "of," and "and," reflecting common English usage. Other domain-specific terms like "patients," "treatment," and "study" are also prominent, indicating the medical context of the dataset.
- o These frequent terms provide insights into common themes across abstracts, although many will be removed during preprocessing to reduce noise.

4. PICO Tag Co-occurrence Matrix:

- o The co-occurrence matrix shows how often different PICO tags appear together within the same abstract.
- o **Method** and **Results** frequently co-occur, which suggests that abstracts often detail the methodology alongside the study's findings.
- o There is also a significant co-occurrence between **Conclusion** and **Results**, emphasizing that conclusions are typically drawn directly from study outcomes.
- o Tags like **Participants** and **Intervention** have lower co-occurrence rates, indicating these components are often independently detailed.



These EDA findings provide a foundational understanding of the dataset's structure and highlight key patterns across PICO elements, helping guide feature engineering and model development strategies.

3. Feature Engineering:

The Feature Engineering step focuses on transforming the dataset into a format compatible with neural network models by converting textual data into numerical indices and preparing labels for multi-class classification.

1. **Tag Mapping:** Each PICO element is mapped to a unique index for training:

- o **Aim:** 0
- o **Method:** 1
- o **Participants:** 2
- o **Intervention:** 3
- o **Outcome:** 4
- o **Results:** 5
- o **Conclusion:** 6

```
# Map tags to indices for training
tag_to_index = {'Aim': 0, 'Method': 1, 'Participants': 2, 'Intervention': 3, 'Outcome': 4, 'Results': 5, 'Conclusion': 6}
```


This mapping allows the model to treat each PICO element as a separate class during classification.

2. **Sentence Indexing:** Words in each sentence are converted to numerical indices using a pre-trained Word2Vec model. Words not present in the model's vocabulary are assigned a default index of 0. This numerical representation helps the neural network process textual data efficiently.

```
# Prepare dataset for training
def prepare_dataset(dataset, word2vec_model, max_sequence_length):
    X, y = [], []
    for sentences, tags in dataset:
        for sentence, tag in zip(sentences, tags):
            # Convert words to indices if they exist in Word2Vec vocabulary
            indexed_sentence = [word2vec_model.wv.key_to_index.get(word, 0) for word in sentence]
            X.append(indexed_sentence)
            y.append(tag_to_index[tag])

    # Pad sequences and convert labels to categorical format
    X = pad_sequences(X, maxlen=max_sequence_length, padding='post')
    y = to_categorical(y, num_classes=len(tag_to_index))
    return X, y
```

3. **Padding Sequences:** To standardize input size, sentences are padded to a maximum length of 100 words. Padding ensures consistent sequence length across the dataset, allowing the model to process batches of text without varying input shapes.
4. **Categorical Labels:** The target labels are converted into a categorical format using one-hot encoding. This transformation enables multi-class classification, where the model predicts the probability of each PICO category.

This feature engineering process is essential for preparing the data for subsequent model training steps, ensuring that each PICO element is represented numerically and uniformly, facilitating accurate predictions.

4.1. Custom Word2Vec Embeddings:

To capture the semantic relationships between words in biomedical literature, we trained a **Custom Word2Vec** model using the **Gensim** library. The model was trained on tokenized clinical text from the training set, with an embedding dimension of 100 to balance computational efficiency and the ability to capture nuanced word relationships. We used the **skip-gram** approach (sg=1), which is particularly effective for learning high-quality word representations, especially for rare words that might appear in clinical texts.

The Word2Vec model was trained with the following parameters:

- **Embedding Dimension:** 100, to represent words as dense vectors.
- **Window Size:** 5, which helps the model understand words in a broader context, capturing semantic similarities between words in the same context window.
- **Min Count:** 1, allowing the model to learn representations for even rare words in the corpus.
- **Workers:** 10, using multiple cores to speed up the training process.

```
from gensim.models import Word2Vec

# Train Word2Vec model on the tokenized corpus
embedding_dim = 100 # Dimension of embeddings
word2vec_model = Word2Vec(sentences=unlabelled_train_sentences, vector_size=embedding_dim, window=5, min_count=1, workers=10, sg=1)
word2vec_model.save("word2vec_model")

print("Word2Vec model trained and saved.")
```

Once trained, the Word2Vec embeddings were evaluated for coverage, which measures how well the vocabulary of the test set is represented in the model. We found that a high percentage of words in the test set were covered by the embeddings, suggesting that the model learned a rich and comprehensive set of word representations.

```
# Function to calculate coverage
def evaluate_coverage(unlabelled_test_sentences, embedding_model):
    vocab = set([word for sentence in unlabelled_test_sentences for word in sentence])
    covered_words = sum(1 for word in vocab if word in embedding_model.key_to_index)
    coverage = covered_words / len(vocab) * 100
    return coverage

# Calculate test coverage
word2vec_test_coverage = evaluate_coverage(unlabelled_test_sentences, word2vec_model.wv)
print(f"Word2Vec Test Coverage: {word2vec_test_coverage:.2f}%")
```

To further evaluate the quality of the embeddings, we tested the similarity between key biomedical terms. Words such as **"treatment"**, **"patient"**, **"outcome"**, and **"study"** showed high similarity with related terms like **"therapy"**, **"participant"**, and **"trial"**, indicating that the embedding layer successfully captured semantic relationships specific to the biomedical domain.

```
# Test similarity between words in test sentences
sample_words = [
    "treatment", "patient", "therapy", "outcome", "design",
    "study", "trial", "intervention", "results", "effect",
    "randomized", "controlled", "conclusion", "participants",
    "measures", "assessment", "efficacy", "safety", "risk",
    "significance", "placebo", "response", "disease", "method",
    "group", "symptoms", "analysis", "medical", "condition",
    "protocol", "mortality", "duration", "improvement",
    "adverse", "dose", "factor", "evaluation", "profile",
    "outcomes", "comparison", "data", "findings", "testing"
]

print("\nWord Similarities:")
for word in sample_words:
    if word in word2vec_model.wv:
        similar_words = word2vec_model.wv.most_similar(word, topn=5)
        print(f"Most similar words to '{word}':")
        for similar_word, similarity in similar_words:
            print(f"    {similar_word}: {similarity:.4f}")
    else:
        print(f"'{word}' not found in vocabulary.")
```

Conclusion:

Overall, the **Custom Word2Vec embeddings** demonstrate the model's ability to understand and represent important terms in biomedical research, making them highly valuable for downstream tasks like PICO element detection. The high coverage and semantic similarity of words further validate the quality of the learned embeddings, providing a solid foundation for text classification and other NLP tasks in the biomedical field.

4.2. Glove Embeddings:

To complement the Word2Vec embeddings, we also explored GloVe (Global Vectors for Word Representation) embeddings. GloVe is a pre-trained word embedding model that represents words as vectors in a high-dimensional space, with the goal of capturing semantic relationships between words based on their co-occurrence in large corpora. GloVe embeddings were chosen for their ability to capture global word co-occurrence statistics from a corpus and provide dense, high-quality representations for a wide variety of words.

Loading and Saving the GloVe Model:

We used Gensim's KeyedVectors to load the GloVe embeddings, which were saved in a format compatible with Word2Vec. The embeddings were then loaded into the `glove_model` variable. After training, the model was saved to disk for future use to ensure the embeddings could be reused without needing to reload the original file every time.

```
# Load using Gensim's KeyedVectors
glove_model = KeyedVectors.load_word2vec_format(glove_input_file, binary=False, no_header=True)

# Save the GloVe model
glove_model.save("glove_model")
print("GloVe model saved successfully as 'glove_model'")
```

Evaluating GloVe Embedding Coverage:

Once the GloVe model was loaded, we evaluated its **coverage** on the test set, similar to the Word2Vec model. Coverage measures the proportion of words from the test set that are represented in the GloVe model's vocabulary. For this task, the **GloVe Test Coverage** was calculated and printed, showing that only **16.61%** of the words in the test set were present in the GloVe model's vocabulary. This indicates that the GloVe embeddings might not cover as much of the biomedical-specific vocabulary as the Word2Vec model, but they still provide useful general-domain word representations.

```
# Assuming `unlabelled_test_sentences` is your tokenized test data
glove_test_coverage = evaluate_coverage(unlabelled_test_sentences, glove_model)
print(f"GloVe Test Coverage: {glove_test_coverage:.2f}%")
```

Testing Semantic Similarity:

We then tested the **semantic similarity** of key biomedical terms using the GloVe model. For each word in the sample list, we retrieved the **top 5 most similar words** according to the cosine similarity of their vectors. This helps assess how well the GloVe model captures related terms. For example, terms like **"treatment"** and **"therapy"** were identified as closely related, with a similarity score of **0.7391**. Other terms such as **"study"** and **"trial"** showed reasonable similarity with terms like **"exam"** and **"science"**, indicating that the model understands the broader context of academic and medical research.


```
# Test similarity between words in GloVe model
print("\nWord Similarities using GloVe:")
for word in sample_words:
    if word in glove_model:
        similar_words = glove_model.most_similar(word, topn=5)
        print(f"Most similar words to '{word}':")
        for similar_word, similarity in similar_words:
            print(f"  {similar_word}: {similarity:.4f}")
    else:
        print(f"'{word}' not found in GloVe vocabulary.")
```

Model Performance:

The output of the similarity tests reveals that while GloVe embeddings performed well in recognizing relationships between common terms, some biomedical-specific terms yielded less relevant results. For instance, the term **"intervention"** returned related words like **"violence"** and **"aid"**, which are contextually related but may not always be relevant to the specific clinical context. Similarly, **"design"** showed associations with design in the context of art and custom design rather than clinical trial design.

Conclusion:

The **GloVe embeddings** provide useful word representations, particularly for general-domain terms. However, their limited coverage of the biomedical-specific vocabulary, as seen in the **16.61% test coverage**, suggests that they may not fully capture the domain-specific nuances required for tasks like **PICO element detection** in clinical texts. Despite this, GloVe embeddings still offer valuable semantic information, especially for common terms, and could be used in combination with other domain-specific models like Word2Vec for improved performance.

4.3. Fasttext Embeddings:

FastText is an extension of Word2Vec that represents each word as a bag of character n-grams, allowing it to capture subword information. This is particularly useful for handling out-of-vocabulary (OOV) words and understanding the morphology of words in various contexts, which is beneficial in specialized fields like biomedical literature. In this section, we trained a FastText model on the tokenized training data, evaluated its coverage on the test set, and tested its ability to generate meaningful word embeddings and similarities.

```
import fasttext

# Train the FastText model
embedding_dim = 100 # Set embedding dimensions
fasttext_model = fasttext.train_unsupervised(unlabelled_train_sentences, model='skipgram', dim=embedding_dim)
```

Model Training and Saving:

We began by training the FastText model using the skip-gram method, which predicts context words given a target word. The embedding dimension was set to 100, similar to the Word2Vec and GloVe models, to balance the quality of the embeddings with computational efficiency. After training, the model was saved for later use.

```
# Save the trained FastText model
fasttext_model.save_model("fasttext_model")
print("FastText model trained and saved.")

# Load the saved FastText model
fasttext_model = fasttext.load_model("fasttext_model")
print("FastText model loaded successfully.")
```

Embedding Representation:

To evaluate the quality of the embeddings, we retrieved the **word vectors** for a set of sample words. FastText's ability to generate embeddings for words not seen during training (OOV words) was also tested, making it a more robust model for domain-specific tasks where novel or rare terms often occur.

The generated embeddings provide a dense representation of words, where each word is represented by a 100-dimensional vector. These vectors capture the semantic meaning of the word in the context of the training data, making it suitable for downstream tasks like PICO element detection.

```
# Check embeddings for sample words
print("\nSample FastText Embeddings:\n")
for word in sample_words:
    embedding = fasttext_model.get_word_vector(word)
    formatted_embedding = ", ".join(f"{val:.4f}" for val in embedding[:10]) # Display first 10 dimensions
    print(f"Word: {word}\nEmbedding:\n[{formatted_embedding}]\n")
```

FastText Test Coverage:

The **FastText model** was evaluated for its coverage of words in the test set. The coverage metric measures how many words in the test set have corresponding embeddings in the model's vocabulary. Since FastText can handle subword information, we expect its

coverage to be higher than that of the GloVe model, which struggled with domain-specific terms.

Based on our evaluation, **FastText Test Coverage** was calculated at **30.21%**, which indicates a substantial improvement over GloVe's **16.61%** coverage. This shows that FastText is more effective at generating embeddings for unseen words, especially those that might be rare or domain-specific.

```
# Assuming `unlabelled_test_sentences` is your tokenized test data
fasttext_test_coverage = evaluate_fasttext_coverage(unlabelled_test_sentences, fasttext_model)
print(f"FastText Test Coverage: {fasttext_test_coverage:.2f}%")
```

Word Similarities:

We also tested the **semantic similarity** between words in the FastText model by retrieving the **top 5 most similar words** for each of the sample words. The results showed that FastText performed well in identifying related terms, such as **"treatment"** and **"therapy"**, as well as broader medical terms like **"patient"** and **"care"**. The similarity scores reflect the model's ability to capture both morphological and semantic similarities between words.

For example:

- **"treatment"** was most similar to **"treatments"**, **"therapy"**, and **"acne"**, with similarity scores ranging from **0.7020** to **0.8163**.
- **"patient"** showed similarities with terms like **"patients"**, **"care"**, and **"critical"**, indicating that FastText effectively captures relationships between medical terms.

```
# Test word similarities
print("\nWord Similarities using FastText:\n")
for word in sample_words:
    similar_words = fasttext_model.get_nearest_neighbors(word, k=5)
    print(f"Most similar words to '{word}':")
    for similarity, similar_word in similar_words:
        print(f"  {similar_word}: {similarity:.4f}")
```

Model Performance:

The **FastText model** demonstrated **good performance** in both embedding generation and capturing semantic relationships. While it still had some limitations in highly specialized biomedical terms (similar to Word2Vec and GloVe), it outperformed GloVe in terms of coverage. The model's ability to handle OOV words by leveraging subword information made it a valuable tool for processing domain-specific terms, contributing to its **30.21% test coverage**.

Conclusion: The **FastText embeddings** offer a significant improvement over the **GloVe model**, particularly in terms of **coverage** and the ability to generate embeddings for out-of-vocabulary words. While it still struggles with certain domain-specific terms, its ability to generate meaningful representations for rare and unseen words makes it highly suitable for biomedical applications like **PICO element detection**. The model's robustness in handling subword information allows it to perform well in tasks where vocabulary diversity is high, offering a strong alternative to traditional word embeddings.

4.4. Embedding Comparison: Word2Vec, GloVe, and FastText

In this section, we compare the performance and test coverage of three different word embedding models—**Word2Vec**, **GloVe**, and **FastText**—to evaluate their effectiveness in capturing the semantic relationships and vocabulary coverage needed for tasks like **PICO element detection** in biomedical literature.

Test Coverage:

Test Coverage:

Test coverage measures the percentage of words from the test set that are represented in

the embedding model's vocabulary. A higher coverage percentage indicates that the embedding model has learned a more comprehensive vocabulary, directly impacting its ability to generalize to unseen data.

Word2Vec: The Word2Vec model achieved a test coverage of 76.19%, the highest among the original three models, demonstrating its ability to capture a wide range of words and semantic relationships in the context of biomedical terms. Word2Vec's capability to represent both general and domain-specific terms makes it an effective choice for handling diverse vocabularies.

- **GloVe:** The GloVe model demonstrated the lowest coverage at 16.61%. This reduced coverage is likely due to GloVe's reliance on pre-trained embeddings that were not optimized for biomedical text, limiting its effectiveness in domain-specific tasks. While it performs well in capturing general semantic relationships, its limited vocabulary coverage makes it less effective for this particular domain.
- **FastText:** The FastText model now matches the test coverage of Word2Vec at 76.19%. Its subword-based approach (i.e., character n-grams) allows it to handle rare and out-of-vocabulary words effectively. This updated coverage shows that

FastText can match Word2Vec in capturing comprehensive vocabulary, making it a robust choice for biomedical tasks with complex, diverse terminology.

Semantic Similarity:

The semantic similarity of key biomedical terms was also compared to evaluate how well each model captures relationships between related terms.

- **Word2Vec:** Word2Vec consistently identified highly relevant similar words, such as "treatment" being related to "therapy" and "study" to "trial." This indicates that it excels in capturing both general-domain and domain-specific relationships, reinforcing its strength in semantic understanding for biomedical literature.
- **GloVe:** GloVe successfully captured general semantic relationships (e.g., "therapy" related to "treatment") but struggled with more domain-specific terms, occasionally producing contextually inaccurate results (e.g., "intervention" being similar to "violence"). Its limited coverage and lower specificity to biomedical contexts highlight its weaker performance in this domain compared to Word2Vec and FastText.
- **FastText:** With its updated high coverage and subword-level modeling, FastText excelled in capturing similarities for morphologically related and unseen words. For example, it identified nuanced relationships such as "treatmentfess" being highly similar to "treatment" (score: 0.8607) and effectively handled diverse biomedical terms like "intervention" (e.g., "δintervention," similarity: 0.9570). This demonstrates its improved ability to provide meaningful embeddings for rare terms and domain-specific nuances.

Model Performance:

The performance of each embedding model for PICO element detection depends on its ability to accurately represent domain-specific terms.

- **Word2Vec:** With 76.19% test coverage and strong semantic similarity, Word2Vec provided excellent performance in both accuracy and classification metrics, making it well-suited for detecting and classifying PICO elements.
- **GloVe:** GloVe's limited 16.61% coverage reduced its ability to accurately represent domain-specific terms, impacting its overall performance for precise biomedical text classification tasks.
- **FastText:** The updated 76.19% coverage allowed FastText to rival Word2Vec in terms of vocabulary comprehensiveness. Its robust subword modeling offers unique

strengths in capturing rare or unseen biomedical terms, resulting in competitive performance in detecting and classifying PICO elements. While its semantic performance was comparable to Word2Vec, it sometimes produced less relevant similarities due to its granular focus on subwords.

Conclusion:

The comparison highlights the strengths and weaknesses of each embedding approach for biomedical tasks:

- **Word2Vec** remains a leading performer in terms of both coverage and semantic similarity, offering a strong choice for tasks involving complex biomedical terminology.
- **FastText** now provides comparable coverage and improved handling of rare and unseen terms, making it a viable alternative for challenging vocabulary needs.
- **GloVe** performs well in general-domain contexts but struggles with domain-specific terms due to its limited coverage.

For tasks requiring high vocabulary coverage and deep semantic understanding, **Word2Vec** and **FastText** emerge as top choices, with FastText being especially effective when handling morphologically diverse or rare terms. GloVe, while useful in broader contexts, is less suitable for specialized biomedical applications.

5. Modelling Techniques:

Model Training and Evaluation

1. Dataset Preparation:

To begin the model training, the dataset was prepared by applying Word2Vec embeddings. The `prepare_dataset()` function was used to transform the raw dataset into numerical features that could be ingested by the model. This function tokenized the text and mapped each word to a vector representation based on the pre-trained Word2Vec model, ensuring that the relationships between words were preserved in the embedding space. The dataset was structured into input sequences (X) and corresponding target labels (y), with each sequence representing a sentence and each label corresponding to the PICO element tag.

A key parameter, the `max_sequence_length`, was set to control the maximum number of words per sentence, ensuring consistent input size for the model. This preprocessing step also included padding shorter sequences to match the `max_sequence_length`, which is crucial for the LSTM model as it expects uniform input dimensions.

```
# Prepare dataset for modeling using Word2Vec embeddings
x, y = prepare_dataset(dataset, word2vec_model, max_sequence_length=max_sequence_length)
```

2. Data Splitting:

After preparing the dataset, we split it into training and testing sets using an 80/20 ratio, ensuring that 80% of the data was used for model training, while 20% was reserved for testing. This split is crucial for evaluating the model's generalization capability and ensuring that it performs well on unseen data. The splitting was done using `train_test_split()` from Scikit-learn with a random state set to 42 for reproducibility of results.

At this stage, the data preparation was completed successfully, and the next step was to move forward with model training and evaluation

```
# Split into training and testing sets
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

print("Data preparation complete.")
```

Model Evaluation Metrics

To evaluate the performance of our models for PICO element detection, we use several key metrics that are commonly reported in classification tasks: **accuracy**, **precision**, **recall**, and **F1-score**. These metrics provide a comprehensive view of the model's effectiveness at correctly identifying the PICO components in biomedical text. Here is a brief overview of each metric and what it represents:

- **Accuracy:** This metric indicates the proportion of correct predictions made by the model out of all predictions. While accuracy gives a general idea of how well the model performs, it may not fully capture performance when the dataset is imbalanced (i.e., some classes have much more data than others).
- **Precision:** Precision is the ratio of true positive predictions to the total predicted positives (i.e., true positives plus false positives). It measures the model's ability to correctly identify positive instances without mistakenly classifying negatives as positives. In the context of PICO element detection, high precision means that when the model predicts a particular PICO component, it is more likely to be correct.
- **Recall (Sensitivity):** Recall is the ratio of true positive predictions to the total actual positives (i.e., true positives plus false negatives). It measures the model's ability to capture all relevant instances of a class. High recall indicates that the model can

identify most of the positive examples in the dataset, even if some predictions might be incorrect.

- **F1-Score:** The F1-score is the harmonic mean of precision and recall, providing a balanced measure that takes both false positives and false negatives into account. A high F1-score indicates that the model has a good balance between precision and recall, making it especially useful for evaluating model performance on imbalanced datasets.

By analyzing these metrics, we can gain a deeper understanding of the model's strengths and weaknesses. For example, high precision but low recall might indicate that the model is conservative in its predictions, while high recall but low precision suggests that the model makes many predictions but often misclassifies. In combination, these metrics provide a nuanced view of the model's performance in identifying PICO elements, helping us to pinpoint areas for improvement and guide further optimization efforts.

5.1. Baseline Model: Logistic Regression with TF-IDF

To establish a baseline for comparison, we implemented a logistic regression model using TF-IDF (Term Frequency-Inverse Document Frequency) features. This approach provides a numerical representation of the text based on the importance of words in the dataset, which is useful for traditional machine learning models like logistic regression.

The dataset was first flattened, where each sentence was converted into a single string. Then, we applied the `TfidfVectorizer` to convert these sentences into a sparse matrix of TF-IDF features. The maximum number of features was set to 5000 to reduce dimensionality and focus on the most important words in the corpus.

The logistic regression model was then trained on these TF-IDF features. For evaluation, we used the `accuracy_score` and `classification_report` metrics from Scikit-learn, which provide insights into the performance of the model across different classes.

Model Performance:

The accuracy of the Logistic Regression model on the test set was **27.3%**, which indicates a fairly weak performance. The detailed classification report reveals significant class imbalances, with very low precision, recall, and F1-scores for several PICO components. Notably:

- **Class 5** (Results) achieved relatively high recall (87%) but had a low precision (29%), suggesting the model predominantly predicts this class but with many false positives.
- **Class 1** (Method) and **Class 6** (Conclusion) also had lower performance, with precision and recall scores significantly lower than ideal.
- Other classes like **Class 0** (Aim) and **Class 2** (Participants) exhibited very poor performance, with recall approaching zero, highlighting the model's inability to accurately detect these components in the medical text.

The poor performance of the logistic regression model suggests that the linear approach based on TF-IDF features is insufficient for accurately detecting the complex relationships between words in the biomedical text.

Logistic Regression Accuracy: 0.2718474845816687				
Classification Report for Logistic Regression:				
	precision	recall	f1-score	support
0	0.13	0.01	0.02	6315
1	0.17	0.09	0.12	9336
2	0.09	0.00	0.00	4449
3	0.06	0.00	0.00	3945
4	0.09	0.00	0.01	5297
5	0.29	0.86	0.43	14996
6	0.16	0.03	0.05	7224
accuracy			0.27	51562
macro avg	0.14	0.14	0.09	51562
weighted avg	0.18	0.27	0.16	51562

Conclusion:

The logistic regression model based on TF-IDF features provides a reasonable starting point but demonstrates significant limitations in capturing the complex structure of biomedical language. As expected, the traditional machine learning model struggles with the diverse and nuanced vocabulary of the dataset, underscoring the need for more advanced models, such as LSTMs or transformers, which can better handle sequential relationships between words. The next steps will involve moving forward with more sophisticated deep learning models for improved PICO element detection.

5.2. Feedforward Neural Network (FFNN) Model

To further explore the model performance, we implemented a Feedforward Neural Network (FFNN) using PyTorch. The FFNN was trained on the preprocessed dataset, with the input features derived from the Word2Vec embeddings. The model was structured with a simple architecture consisting of three fully connected layers with ReLU activations between them. The output layer was designed to have as many neurons as there are PICO classes, and the model was trained using cross-entropy loss for multi-class classification.

Model Architecture:

The FFNN model consists of the following layers:

- **Input Layer:** Takes in the tokenized and embedded sequence of words.
- **Hidden Layers:** Two fully connected layers, with 256 and 128 neurons respectively, followed by ReLU activation functions to introduce non-linearity.
- **Output Layer:** A fully connected layer with the number of neurons corresponding to the total number of classes (PICO components).

The model was trained on the GPU (if available) and was evaluated over 3 epochs using the Adam optimizer with a learning rate of $1e-3$.

Model Performance:

After training for three epochs, the FFNN model achieved an accuracy of **29.1%** on the test set. The classification report indicates that the model has poor performance across several PICO components:

- **Class 5** (Results) performed relatively well, with a recall of **100%**, indicating that the model correctly identifies nearly all instances of this class. However, its precision is quite low (29%), meaning there are many false positives.
- Other classes, such as **Class 0** (Aim), **Class 1** (Method), **Class 2** (Participants), **Class 3** (Intervention), **Class 4** (Outcome), and **Class 6** (Conclusion), showed no recall or precision, indicating that the model failed to recognize these components in the test data.

```
Epoch 1 completed.
Epoch 2 completed.
Epoch 3 completed.
FFNN Accuracy: 0.290834335363252
```

Classification Report for FFNN:				
	precision	recall	f1-score	support
0	0.00	0.00	0.00	6315
1	0.00	0.00	0.00	9336
2	0.00	0.00	0.00	4449
3	0.00	0.00	0.00	3945
4	0.00	0.00	0.00	5297
5	0.29	1.00	0.45	14996
6	0.00	0.00	0.00	7224
accuracy			0.29	51562
macro avg	0.04	0.14	0.06	51562
weighted avg	0.08	0.29	0.13	51562

Conclusion:

The FFNN model, while demonstrating some ability to detect Class 5 (Results), still struggles with the other PICO components. The low recall for most of the classes indicates that a simple FFNN model is insufficient for the task of PICO element detection, particularly in a complex domain like biomedical literature. The results suggest that a more sophisticated model, such as a Recurrent Neural Network (RNN) or Long Short-Term Memory (LSTM) network, may be better suited for this problem due to its ability to capture the sequential dependencies in the text.

This FFNN approach, while providing some insight into the dataset, highlights the limitations of traditional feedforward networks for sequential tasks in natural language processing. Further work will focus on leveraging sequence-based models, which are better equipped to understand the contextual relationships between words and improve classification performance across all PICO classes.

5.3. Recurrent Neural Networks: LSTM Model

Building on the previous models, we implemented a BiLSTM (Bidirectional Long Short-Term Memory) model with an attention layer to capture the sequential dependencies in biomedical text. Unlike the Feed-Forward Neural Network (FFNN), which treats each word independently, the BiLSTM model can capture relationships between words in both forward

and backward directions, making it more suitable for understanding the context of biomedical terms in the PICO elements.

Dataset Preparation:

As with the previous models, the dataset was prepared using the LSTMDataset class, where the features were converted to LongTensor (required by LSTM models). This dataset was split into training and test sets, and these sets were loaded into data loaders for efficient batch processing.

```
# PyTorch Dataset for LSTM
class LSTMDataset(Dataset):
    def __init__(self, features, labels):
        self.features = torch.tensor(features, dtype=torch.long) # Convert to LongTensor
        self.labels = torch.tensor(labels.argmax(axis=1), dtype=torch.long) # Convert labels to LongTensor

    def __len__(self):
        return len(self.labels)

    def __getitem__(self, idx):
        return self.features[idx], self.labels[idx]
```

Model Architecture:

The BiLSTM with Attention model consists of the following layers:

- **Embedding Layer:** Converts word indices to dense vector representations using pre-trained Word2Vec embeddings.
- **BiLSTM Layer:** Captures sequential relationships between words in the sentence in both forward and backward directions.
- **Attention Layer:** Applies attention weights to the LSTM output to focus on the most important words in the sequence.
- **Fully Connected Layer:** Outputs the predicted class (PICO component).

The BiLSTM model was trained for 10 epochs, and the loss was computed using the cross-entropy loss function. The optimizer used was Adam, which is known for its efficiency in training deep learning models.

Model Training:

The BiLSTM model with attention was trained for 10 epochs, using the Adam optimizer with a learning rate of 1e-3. During each epoch, the model was updated using the loss calculated from the predicted and actual labels.


```

# Loss and Optimizer
criterion = nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(model.parameters(), lr=1e-3)

# Train BiLSTM with Attention
model.train()
for epoch in range(10): # Train for 3 epochs
    total_loss = 0
    for features, labels in train_loader:
        features, labels = features.to(device), labels.to(device)
        optimizer.zero_grad()
        outputs = model(features)
        loss = criterion(outputs, labels)
        loss.backward()
        optimizer.step()
        total_loss += loss.item()
    print(f"Epoch {epoch + 1}, Loss: {total_loss:.4f}")

```

Model Evaluation:

After training, the model was evaluated on the test set. The BiLSTM with Attention model achieved an accuracy of **77.74%** on the test set. The detailed classification report reveals the following:

- **Class 5 (Results)** performed the best with a **precision** of 84%, **recall** of 85%, and an **F1-score** of 85%. This indicates that the model can effectively detect the Results component, although there are still some false positives (lower precision).
- **Class 0 (Aim)**, **Class 1 (Method)**, and **Class 6 (Conclusion)** showed reasonable performance, with **recall** values around 77%–82%, though precision was lower for some classes.
- **Class 2 (Participants)**, **Class 3 (Intervention)**, and **Class 4 (Outcome)** performed poorly, with recall values approaching zero for these classes, suggesting that the model struggled with detecting these components, likely due to less frequent occurrences in the dataset.

Conclusion:

The BiLSTM with Attention model significantly outperformed the FFNN model, achieving an accuracy of 77.74%. It was particularly successful at detecting Class 5 (Results), but it still faced challenges in detecting components like Participants (Class 2), Intervention (Class 3), and Outcome (Class 4). This highlights the model's ability to capture sequential dependencies but also suggests that improvements are needed, particularly for underrepresented classes. Further refinements could include experimenting with other advanced models such as GRU or fine-tuning hyperparameters, and potentially leveraging more sophisticated embeddings like BERT for better semantic understanding.

BiLSTM with Attention Model Accuracy: 0.7774				
Classification Report for BiLSTM with Attention Model:				
	precision	recall	f1-score	support
0	0.84	0.85	0.84	6315
1	0.77	0.77	0.77	9336
2	0.74	0.71	0.73	4449
3	0.62	0.60	0.61	3945
4	0.72	0.75	0.73	5297
5	0.84	0.85	0.85	14996
6	0.74	0.73	0.74	7224
accuracy			0.78	51562
macro avg	0.75	0.75	0.75	51562
weighted avg	0.78	0.78	0.78	51562

5.4. BERT Model

In this step, we leverage **BERT** (Bidirectional Encoder Representations from Transformers), a pre-trained transformer model designed for natural language processing tasks. BERT has shown significant success in understanding the context of words in a sentence, making it a strong candidate for PICO element detection in biomedical literature.

Data Preparation:

The data was first preprocessed and tokenized using the **BERT tokenizer**, which is part of the Hugging Face **Transformers** library. The sentences were padded or truncated to a maximum length of 128 tokens, and the labels were mapped to integer values based on the PICO components.

The data was split into training and testing sets (80/20 split), and the BERT tokenizer was used to convert the text into input IDs and attention masks, which are the required inputs for BERT.

BERT Model Initialization:

The **BertForSequenceClassification** model was initialized with the pre-trained **bert-base-uncased** weights, and the number of output labels was set to match the number of PICO components.

```

from transformers import BertForSequenceClassification

# Define BERT model
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
bert_model = BertForSequenceClassification.from_pretrained(
    "bert-base-uncased",
    num_labels=len(set(train_labels)) # Number of classes
).to(device)

```

Training:

The model was trained using the **AdamW** optimizer with a learning rate of **2e-5** and a batch size of **16**. The model was trained for **3 epochs**, and the loss was computed using the cross-entropy loss function.

```

# Optimizer
optimizer = AdamW(bert_model.parameters(), lr=2e-5)

# Training loop
bert_model.train()
for epoch in range(3): # Number of epochs
    total_loss = 0
    for batch in train_loader:
        input_ids = batch["input_ids"].to(device)
        attention_mask = batch["attention_mask"].to(device)
        labels = batch["labels"].to(device)

        optimizer.zero_grad()
        outputs = bert_model(input_ids=input_ids, attention_mask=attention_mask, labels=labels)
        loss = outputs.loss
        loss.backward()
        optimizer.step()

    total_loss += loss.item()
    print(f"Epoch {epoch + 1} completed. Loss: {total_loss:.4f}")

```

Evaluation:

After training, the model was evaluated on the test set, and the accuracy and classification report were generated to assess performance.

Model Performance:

The **BERT model** achieved an impressive accuracy of **84.2%** on the test set. The classification report shows that the model performs well across most classes, especially for **Class 5** (Results), with high precision (90%) and recall (91%). Other notable classes

include **Class 0** (Aim), **Class 1** (Method), and **Class 6** (Conclusion), which also show strong performance in terms of both precision and recall.

BERT Model Accuracy: 0.8417				
Classification Report:				
	precision	recall	f1-score	support
0	0.90	0.92	0.91	5317
1	0.83	0.76	0.80	7793
2	0.78	0.81	0.80	3767
3	0.71	0.71	0.71	3335
4	0.79	0.82	0.80	4485
5	0.90	0.91	0.90	12641
6	0.83	0.83	0.83	5951
accuracy			0.84	43289
macro avg	0.82	0.82	0.82	43289
weighted avg	0.84	0.84	0.84	43289

Conclusion:

The **BERT model** significantly outperforms previous models, achieving **84.2% accuracy** on the test set. It performs particularly well in **Class 5** (Results) and **Class 6** (Conclusion), demonstrating its ability to handle complex text and capture the sequential relationships between words. This performance validates the power of transformer-based models like BERT for natural language tasks such as PICO element detection. Further improvements can be made by fine-tuning hyperparameters or leveraging even larger pre-trained models, such as **BERT-large** or domain-specific variants like **BioBERT** for biomedical text

5.5. Knowledge Distillation

Knowledge distillation is a technique that allows a smaller model, referred to as the "student" model, to learn from a larger, more complex model, known as the "teacher" model. The goal is for the student model to achieve similar performance to the teacher model, but with fewer parameters and faster inference times. This implementation uses knowledge distillation to train a smaller, more efficient **DistilBERT** model to mimic the performance of the larger **BERT** model.

Model Definition:

The **teacher model** is **BERT** (Bidirectional Encoder Representations from Transformers), a large and powerful model pre-trained on a vast amount of text data. The **student model** is **DistilBERT**, a smaller, distilled version of BERT designed to retain much of its performance

while being more efficient in terms of model size and computational requirements. Both models are fine-tuned for sequence classification tasks.

```
from transformers import BertForSequenceClassification, DistilBertForSequenceClassification
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
# Define the teacher (BERT) model
teacher_model = BertForSequenceClassification.from_pretrained(
    "bert-base-uncased",
    num_labels=len(set(train_labels))
).to(device)

# Define the student (DistilBERT) model
student_model = DistilBertForSequenceClassification.from_pretrained(
    "distilbert-base-uncased",
    num_labels=len(set(train_labels))
).to(device)
```

Distillation Loss:

The **distillation loss** combines two components:

1. **Cross-Entropy Loss:** This measures how well the student's predictions match the true labels.
2. **KL-Divergence:** This measures how similar the student's output distribution is to the teacher's output distribution. By softening the logits of the teacher model with a **temperature** parameter, the student model learns to approximate the teacher's behavior more effectively.

Training the Student Model:

During training, the teacher model's weights are frozen (set to evaluation mode), and the student model is trained to minimize the combined loss (cross-entropy and KL-divergence). We used the Adam optimizer with a learning rate of $2e-5$ and trained the student model for 3 epochs.

Evaluating the Student Model:

After training, we evaluated the **student model** on the test set using traditional classification metrics such as **accuracy** and the **classification report**. The student model's performance was found to be **84.43% accurate**, demonstrating that the student was able to mimic the teacher's performance with significantly fewer parameters.

Student Model Accuracy: 0.8443				
Classification Report for Student Model:				
	precision	recall	f1-score	support
0	0.89	0.93	0.91	6306
1	0.86	0.75	0.80	9235
2	0.86	0.78	0.82	4379
3	0.69	0.73	0.71	3965
4	0.76	0.86	0.81	5197
5	0.89	0.92	0.90	15344
6	0.84	0.81	0.83	7136
accuracy			0.84	51562
macro avg	0.83	0.83	0.83	51562
weighted avg	0.85	0.84	0.84	51562

Conclusion:

The **Knowledge Distillation** process allowed the **DistilBERT student model** to achieve **84.43% accuracy** on the test set, which is **comparable** to the performance of the larger **BERT teacher model**, but with significantly fewer parameters. This demonstrates that distillation can be a powerful technique for creating efficient models that retain much of the performance of their larger counterparts, making them ideal for use in resource-constrained environments. The student model achieved strong results across all PICO classes, proving that **DistilBERT** can perform well even with fewer parameters and a faster inference time.

6. Model Performance Comparison:

In this section, we compare the performance of several models for PICO element detection in biomedical literature: **Logistic Regression with TF-IDF**, **Feedforward Neural Network (FFNN)**, **BiLSTM with Attention**, **BERT**, and **Knowledge Distillation (DistilBERT)**. Each model was evaluated based on its **accuracy**, **coverage**, and its ability to capture the semantic relationships between words and PICO elements.

1. Logistic Regression with TF-IDF

The Logistic Regression with TF-IDF model served as a baseline for our comparison. The model transformed the text data using TF-IDF to weigh word importance based on term frequency and inverse document frequency. The performance of the model was relatively low, with an accuracy of **27.3%**. The classification report revealed that the model struggled to effectively capture the relationships between PICO elements, particularly in detecting rare or domain-specific terms.

- **Accuracy:** 27.3%
- **Key Insight:** While TF-IDF is a useful feature extraction method, it does not capture contextual relationships between words, which is critical for understanding PICO components in biomedical texts.

2. Feedforward Neural Network (FFNN) The FFNN model was trained on Word2Vec embeddings, which allowed it to capture the semantic meaning of words by representing them as dense vectors. The FFNN achieved an accuracy of **29.1%**, a modest improvement over the Logistic Regression model. However, the FFNN still faced challenges in capturing the sequential and contextual dependencies between words, which are crucial for understanding complex PICO components.

- **Accuracy:** 29.1%
- **Key Insight:** FFNNs can model relationships between words but fail to capture sequential dependencies and contextual information, limiting their effectiveness for complex text classification tasks like PICO element detection.

3. BiLSTM with Attention

The **BiLSTM with Attention** model was specifically designed to capture sequential dependencies in the text, and it showed a significant improvement over FFNN. The model achieved an accuracy of **77.74%**, reflecting its ability to capture context and relationships in both forward and backward directions. The attention mechanism further enhanced the model's focus on important words in the sentence. However, it still faced challenges in handling imbalanced class distributions, especially for underrepresented classes.

- **Accuracy:** 77.74%
- **Key Insight:** The BiLSTM with Attention model performed exceptionally well in capturing sequential dependencies and context, particularly in detecting common PICO components like Results and Conclusion. However, its performance could be further improved by addressing class imbalance and fine-tuning hyperparameters.

4. BERT (Teacher Model)

The **BERT model**, known for its ability to capture deep contextual information, achieved impressive results with an accuracy of **84.2%**. BERT outperformed all other models by effectively understanding the relationships between words in the context of the entire sentence. The model performed particularly well in detecting PICO components such as Results and Conclusion, achieving high precision and recall. However, BERT's large size and computational requirements make it resource intensive.

- **Accuracy:** 84.2%
- **Key Insight:** BERT's deep contextual understanding and ability to process bidirectional context made it the best-performing model for PICO element detection, although its computational requirements limit its deployment in resource-constrained environments.

5. Knowledge Distillation (DistilBERT)

DistilBERT, a smaller version of BERT trained using knowledge distillation, achieved an accuracy of **84.43%**, almost identical to BERT's performance. The distillation process allowed DistilBERT to retain much of BERT's performance with fewer parameters, making it significantly more efficient in terms of speed and resource usage. The distillation process also ensured that the student model learned the most important features from the teacher model, achieving high accuracy while being computationally efficient.

- **Accuracy:** 84.43%
- **Key Insight:** DistilBERT demonstrated the effectiveness of knowledge distillation, allowing a smaller model to achieve performance close to a larger, more complex model. It is an excellent option when computational efficiency is critical.

Conclusion:

When comparing the performance of **Logistic Regression with TF-IDF**, **FFNN**, **BiLSTM with Attention**, **BERT**, and **DistilBERT**, several key insights emerge:

- **BERT** and **DistilBERT** achieved the highest accuracy, with **84.2%** and **84.43%**, respectively, demonstrating their ability to capture deep contextual relationships and perform well on complex tasks like PICO element detection. DistilBERT offers a more efficient alternative to BERT, making it ideal for resource-constrained environments.
- **BiLSTM with Attention** showed strong performance with an accuracy of **77.74%**, benefiting from its ability to capture sequential dependencies. It performed particularly well for more frequent PICO components but still struggled with underrepresented classes.
- **FFNN** and **Logistic Regression with TF-IDF** were the least effective, with **29.1%** and **27.3%** accuracy, respectively, due to their inability to capture the complex relationships between words in the text.

For tasks requiring high accuracy and the ability to understand complex relationships in text, **BERT** and **DistilBERT** are the most effective models. **DistilBERT** offers a more

computationally efficient option without sacrificing much performance. However, for simpler tasks or when computational resources are limited, models like **BiLSTM with Attention** and **FFNN** can still provide useful results, although their performance is generally inferior to that of **BERT-based models**.

Best Model:

For the task of PICO element detection, **BERT** and **DistilBERT** are the most suitable models due to their high accuracy of 84.2% and 84.43%, respectively. These models excel in capturing deep contextual relationships, making them ideal for complex tasks. **DistilBERT** offers a more computationally efficient alternative without sacrificing significant performance, making it a great choice for resource-constrained environments.

While **BiLSTM** with Attention performs well with an accuracy of 77.74%, it struggles with underrepresented classes, and its performance is generally inferior to that of BERT-based models. **Logistic Regression** with **TF-IDF** and **FFNN** are less effective, with accuracies of 27.3% and 29.1%, respectively, as they cannot capture the complex relationships in text. Thus, for tasks requiring high accuracy and a strong understanding of complex text relationships, **BERT** or **DistilBERT** should be the preferred model, with **DistilBERT** being more efficient for resource-limited situations.