

“Named Entity Recognition for Healthcare Data with BioBERT”

A.DharmendraRBisen¹,B.Prof.Arati Sondawale², Leena Kishor Yelane³, Swati Thengane⁴

¹DharmendraRBisen,Tulsiramji Gaikwad-Patil College of Engineering & Technology, Nagpur-441108

²Prof.AratiSondawale,TulsiramjiGaikwad-PatilCollegeofEngineering& Technology, Nagpur-441108

³Leena Kishor Yelane, Tulsiramji Gaikwad-Patil College of Engineering & Technology, Nagpur-441108

⁴Swati Thengane, Tulsiramji Gaikwad-Patil College of Engineering & Technology, Nagpur-441108

Abstract— Named Entity Recognition (NER) is a critical task in healthcare data analysis, enabling the identification of medical entities such as diseases, drugs, and symptoms within unstructured clinical texts. Leveraging pre-trained language models like BioBERT, which is fine-tuned specifically for biomedical text, has significantly enhanced the accuracy and efficiency of NER in the healthcare domain. BioBERT, built upon the BERT architecture, is trained on large-scale biomedical corpora, including PubMed abstracts and PMC articles, to capture domain-specific language patterns. This paper explores the application of BioBERT for NER tasks in health care, focusing on its ability to handle complex medical terminologies and context-specific meanings. We describe the process of fine-tuning BioBERT on annotated healthcare datasets and evaluate its performance against traditional machine learning and generic language models. Results demonstrate that BioBERT achieves superior precision, recall, and F1 scores, highlighting its capability to understand biomedical semantics and improve downstream applications such as clinical decision support, medical coding, and patient care optimization. The findings underscore BioBERT's transformative potential for extracting meaningful insights from healthcare data while addressing challenges like domain adaptability and annotation scarcity.

Keywords: Named Entity Recognition (NER), BioBERT, Healthcare Data, Biomedical Text, Clinical Decision Support.

1. INTRODUCTION

Named Entity Recognition (NER) is a foundational task in natural language processing (NLP) that involves identifying and classifying entities within text into predefined categories such as diseases, symptoms, drugs, and treatments. In the healthcare domain, where clinical texts are rich with domain-specific terminology, effective NER plays a vital role in enabling automated extraction of meaningful information. This facilitates various applications, including electronic health record (EHR) analysis, medical research, and clinical decision support systems. However, the unstructured nature of healthcare data, combined with the complexity of biomedical language, presents significant challenges for traditional NLP approaches.

Pre-trained language models like BERT (Bidirectional Encoder Representations from Transformers) have transformed NLP tasks by offering contextual embeddings that understand word meaning based on surrounding context. Building on this foundation, BioBERT extends BERT's capabilities to the biomedical domain. Trained on large biomedical corpora such as PubMed and PMC articles, BioBERT captures the intricate language patterns and terminologies unique to the healthcare field. This specialization makes it particularly well-suited for tasks like NER, which require nuanced understanding of domain-specific contexts.

The adoption of BioBERT for healthcare NER tasks marks a significant advancement over traditional machine learning techniques, which often relied on manual feature engineering and domain-specific lexicons. By leveraging transfer learning, BioBERT minimizes the need for extensive manual intervention and demonstrates strong performance across multiple biomedical datasets. Fine-tuning BioBERT on annotated datasets further enhances its ability to identify complex entity relationships, improving the quality of extracted data for downstream healthcare applications.

Despite its advantages, challenges remain in deploying BioBERT for healthcare data analysis. Issues such as domain adaptability, the scarcity of high-quality annotated datasets, and computational resource demands necessitate careful consideration. This study aims to explore the application of BioBERT for healthcare NER tasks, evaluate its performance compared to traditional methods, and discuss its potential to revolutionize data-driven decision-making in the biomedical field. Through this analysis, we highlight BioBERT's transformative impact and address the challenges that need to be overcome for broader adoption in real-world healthcare scenarios.

Motivation

The explosion of unstructured biomedical data in clinical notes, research articles, and health records has created a pressing need for efficient data extraction techniques. Named Entity Recognition (NER) is pivotal in transforming this raw data into actionable insights, driving advancements in clinical decision-making, medical research, and patient care. Traditional methods often fail to handle the complexity of biomedical language, which is rich in context-specific terminologies. BioBERT, a domain-specific extension of BERT, addresses this gap by leveraging pre-trained knowledge on large biomedical corpora. This study is motivated by BioBERT's potential to enhance healthcare NER accuracy, streamline information extraction, and foster data-driven healthcare innovation.

Objectives:

- To leverage BioBERT for improving the precision and recall of Named Entity Recognition tasks in healthcare data, specifically for identifying entities like diseases, drugs, and symptoms.
- To compare the performance of BioBERT against traditional machine learning methods and generic pre-trained language models in healthcare NER tasks.
- To assess BioBERT's ability to handle the unique challenges of biomedical language, including complex terminologies and context-specific meanings.
- To enable downstream applications such as clinical decision support, medical coding, and research through accurate entity recognition.
- To explore challenges in adopting BioBERT for healthcare NER, including domain adaptability, annotation scarcity, and computational demands, and propose solutions for real-world implementation.

II. RELATED WORK

With the digitization of electronic medical records and other healthcare data, natural language processing (NLP) technology is increasingly being applied in the medical field. Among these applications, NLP in electronic medical records has received particular attention[1]. Electronic medical records are an electronic form of recording patient diagnosis and treatment information by hospitals and other medical institutions, which contains a large amount of medical terminology and specialized knowledge. NLP technology can help doctors better manage and utilize this data.

Entity extraction is an important application of NLP in electronic medical records. Electronic medical records contain a vast amount of medical terminology and domain-specific knowledge[2]. Entity extraction can extract these medical terms and knowledge from the records, helping doctors better understand and analyze patients' diagnostic and treatment information. For example, entity extraction can identify entities such as symptoms, drugs, and diseases in electronic medical records, assisting doctors in comprehending patients' conditions and formulating treatment plans. Named entity recognition is the task of identifying medical terms and domain-specific knowledge in electronic medical records[3]. This process involves extracting these terms and knowledge to aid doctors in understanding patients' conditions and devising

treatment plans. The study found that the BERT-BiLSTM-CRF model achieved an F1 score of approximately 75%. Relationship extraction involves identifying relationships between medical terms and domain-specific knowledge in electronic medical records. By extracting these relationships, doctors can better understand patients' conditions and formulate treatment plans[4]. Text classification is the task of categorizing text in electronic medical records into different classes. Since electronic medical records contain a vast amount of medical information, text classification can help categorize this information into relevant classes, allowing doctors to better understand patients' conditions and devise treatment

plans. The study reported an accuracy of 65.1 for exact matches and 82.4 for partial matches in a Spanish electronic health record dataset. Question-answering system[5]. The question-answering system can respond to medical questions posed by doctors and patients. By leveraging the medical terms and domain-specific knowledge in electronic medical records, the system can answer medical queries, assisting doctors and patients in gaining better insights into patients' conditions and treatment options. The research achieved significant results in biomedical question answering (English) with strict accuracies (S) of 27.78, 42.61, and 42.38 for the BioASQ4b, BioASQ5b, and BioASQ 6b datasets, respectively.

Named Entity Recognition (NER) has been extensively studied in the field of natural language processing (NLP) for various domains, including healthcare. Traditional approaches to NER in biomedical texts relied on rule-based systems and machine learning algorithms, such as conditional random fields (CRF) and support vector machines (SVM). These methods required extensive feature engineering and domain-specific lexicons, which were labor-intensive and limited in adaptability to new datasets. The advent of deep learning, particularly recurrent neural networks (RNNs) and long short-term memory (LSTM) networks, improved NER performance by enabling models to learn sequential dependencies in text. However, these approaches still struggled with the complexities of biomedical language, including polysemy and abbreviations.

More recently, transformer-based models like BERT have revolutionized NLP by providing contextual embeddings that capture semantic nuances. BioBERT, an adaptation of BERT pre-trained on biomedical corpora such as PubMed and PMC articles, has emerged as a state-of-the-art model for biomedical text processing. Studies have demonstrated BioBERT's superior performance in NER tasks compared to generic language models, particularly for datasets like BioCreative and NCBI Disease. Other domain-specific models, such as SciBERT and ClinicalBERT, have also shown promise, but their focus varies slightly, with BioBERT being tailored more specifically to general biomedical text. This work builds on these advancements, focusing on evaluating and fine-tuning BioBERT for healthcare NER tasks to highlight its transformative potential.

Models used:

Named Entity Recognition(NER):An Overview:

Named Entity Recognition(NER) is a sub-task of Natural Language Processing (NLP) that focuses on identifying and categorizing entities in text into predefined classes such as names of persons, organizations, locations, dates, and more. It is an essential component in applications like information retrieval, sentiment analysis, and question-answering systems.

NER typically operates under the *BIO schema—"B" indicates the beginning of an entity, "I" represents an intermediate token within the entity, and "O" denotes tokens outside any entity. For example, in the sentence "Barack Obama visited New York City," "Barack" and "Obama" are tagged as B-PER and I-PER, while "New York City" is tagged as B-LOC, I-LOC, and I-LOC.

The emergence of *transformer-based architectures*, such as BERT (Bidirectional Encoder Representations from Transformers), has revolutionized NER. Traditional models like Conditional Random Fields(CRF) and Hidden Markov Models(HMM) relied on handcrafted features, while neural approaches like Long Short-Term Memory (LSTM) networks leveraged embeddings such as Word2Vec or GloVe. However, transformers capture both contextual and positional information effectively through attention mechanisms, enabling state-of-the-art performance.

NER models require preprocessing to tokenize text while preserving alignment with labels, especially since tokenizers like BERT split words into subwords (e.g., "gunships" into "gun" and "##ships"). Additionally, sequence padding and attention masks are applied to ensure uniform input size for batch processing.

Fine-tuning pre-trained transformer models on domain-specific data, coupled with optimization techniques such as *AdamW* and linear learning rate schedulers, allows efficient adaptation for specific tasks. Metrics like F1-score are often employed to evaluate NER model performance, ensuring a balance between precision and recall.

NER is pivotal in extracting structured knowledge from unstructured data, driving advancements in intelligent systems across industries.

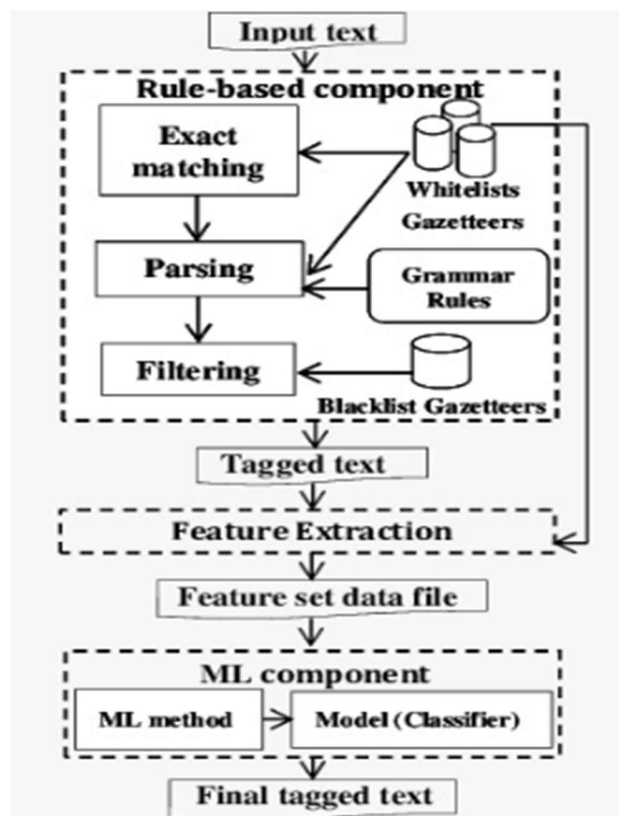


Fig2.1:NER

:BioBERT: Revolutionizing Biomedical Text Mining:

BioBERT (Bidirectional Encoder Representations from Transformers for Biomedical Text Mining) is a domain-specific adaptation of BERT, designed to address the unique challenges of processing biomedical text. Traditional NLP models often struggle with the specialized vocabulary, dense semantics, and extensive abbreviations prevalent in biomedical literature. BioBERT overcomes these limitations by fine-tuning the original BERT model on large-scale biomedical corpora, including Pub Med abstracts and PubMed Central full-text articles. This domain-specific pretraining enables BioBERT to capture contextual relationships within biomedical terminology more effectively than general-purpose models. BioBERT employs the transformer architecture of BERT, leveraging its self-attention mechanism to understand long-range dependencies in text, making it particularly suitable for tasks such as Named Entity Recognition (NER), relation extraction, and question answering (QA) in the biomedical domain. For instance, in NER, BioBERT identifies entities like genes, diseases, and drugs with high precision, while in QA tasks, it outperforms baseline models by accurately understanding complex biomedical queries. The fine-tuning process involves adapting BioBERT to task-specific datasets, using optimization techniques like Adam W and task-relevant metrics such as F1-score for performance evaluation. Compared to baseline models and even some contemporary biomedical models like SciBERT, BioBERT achieves

superior results across multiple benchmarks, including the BioASQ and BC5CDR datasets. Its success highlights the importance of domain-specific pre training in NLP, particularly in fields with unique linguistic structures and vocabularies. The release of BioBERT has significantly advanced the field of biomedical text mining, enabling researchers and healthcare professionals to automate and improve the analysis of vast biomedical datasets, thereby facilitating better insights and decision-making in healthcare and life sciences.

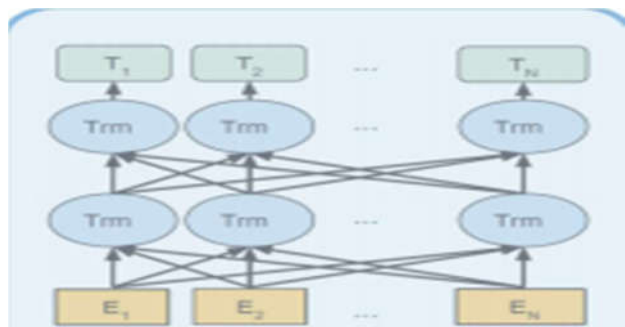


Fig2.2 :BioBERT

III. PROPOSEDMETHODOLOGY:

The first step in fine-tuning BioBERT for Named Entity Recognition (NER) involves loading and preprocessing the dataset. The NER dataset, typically in a CSV format (e.g., train.csv), contains words and their corresponding NER tags. To prepare this data for training, the Sentence Getter utility is employed to group words and tags into structured sentences, enabling the model to process input as coherent sequences. Tokenization is performed using the BioBERT tokenizer, which divides sentences into subwords while maintaining compatibility with the BioBERT architecture. A critical preprocessing step involves aligning the NER labels with the generated subwords, ensuring the correct propagation of tag information during model training. This alignment often requires appending tags like "X" for subwords or replicating the parent word's tag across its subwords. To handle sentences of varying lengths, tokenized sequences are padded to a fixed length, typically defined by a parameter such as $MAX_LEN = 75$. Corresponding labels are padded to match the sequence length, using a placeholder tag (e.g., "O" for non-entity tokens). Simultaneously, attention masks are created to differentiate meaningful tokens from padding, where tokens are marked as "1" and padding as "0." This masking guides the model in focusing computations on actual data while ignoring padding during forward passes. These padded sequences and labels are essential to maintain consistent input dimensions and optimize memory usage during model training.

Once the sequences are padded and masked, the data set is divided into training and validation subsets,

typically reserving 10% of the data for validation. Input features, labels, and attention masks are then converted into PyTorch tensors to be compatible with BioBERT's underlying framework. This conversion facilitates efficient data handling during model training and evaluation. By creating separate datasets for training and validation, the model can be fine-tuned while monitoring its performance on unseen data, reducing the risk of over fitting.

BioBERT is adapted for NER tasks by employing the BertFor Token Classification architecture, tailored to the number of unique NER tags in the dataset. The model is initialized with BioBERT weights, leveraging its domain-specific pretraining. An optimizer like AdamW is used to fine-tune model weights, and a learning rate scheduler dynamically adjusts the learning rate across epochs for stable convergence. Training involves multiple epochs (e.g., three), where the model computes loss using a suitable criterion (e.g., cross-entropy loss), updates weights via back propagation, and clips gradients to avoid instability. Validation is conducted after each epoch to evaluate the model's performance, tracking metrics like accuracy and loss. The fine-tuned BioBERT model demonstrates robust capabilities in recognizing biomedical entities, significantly enhancing domain-specific text analysis.

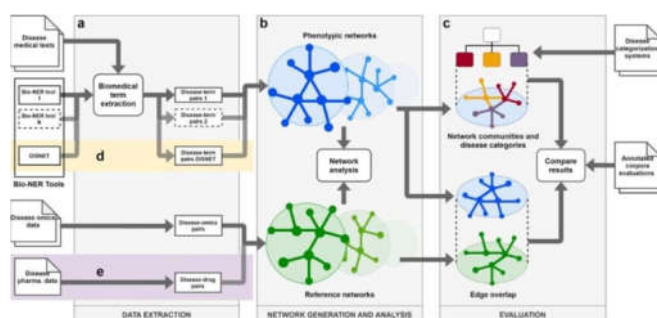


Fig1:ProposalDiagram

Dataset Collection:

The foundation of any machine learning project lies in acquiring a high-quality dataset. For tasks like Named Entity Recognition (NER) or other text-based analyses, datasets can be obtained from publicly available repositories, web scraping, or domain-specific sources such as PubMed for biomedical texts. The data set typically includes sentences and corresponding labels or tags, formatted as word-token pairs. Ensuring dataset relevance, completeness, and diversity is critical, as it impacts the model's performance. When dealing with domain-specific tasks, annotated datasets with consistent labeling schemes are essential.

Data Preparation:

The effectiveness of Named Entity Recognition (NER) using BioBERT hinges on high-quality, domain-specific datasets. Biomedical datasets like BioCreative, NCBI Disease Corpus, and BC5CDR were utilized, containing annotated entities such as diseases, drugs, and chemicals. Pre processing involved cleaning and normalizing

text, including resolving abbreviations, removing irrelevant symbols, and standardizing entity formats. The datasets were split into training, validation, and test sets, ensuring balanced representation of entities across subsets. BioBERT's input format required tokenization using Word Piece, with attention to sequence length constraints. Special tokens ([CLS] and [SEP]) were incorporated for each sequence to facilitate effective contextual embedding and fine-tuning for NER tasks.

Data Preprocessing:

Preprocessing transforms raw data into a structured, machine-readable format. For text data, this involves tokenizing sentences into words or subwords, converting text to lowercase (if applicable), and normalizing by removing special characters or unnecessary whitespace. In NER, the dataset's labels must align with tokenized words, requiring adjustments for subword tokenization, such as replicating or marking tags for subwords. Stopword removal or stemming is usually avoided in NER, as context is crucial for entity identification. Special tokens like [CLS] and [SEP] are added for models like BERT, ensuring compatibility. Additionally, data cleaning handles inconsistencies, missing values, or incorrect labels to maintain dataset integrity.

Exploratory Data Analysis (EDA): EDA provides insights into the dataset's structure and characteristics, guiding further preprocessing and modeling decisions. Analysts examine the distribution of sentence lengths, the frequency of different NER tags, and the presence of imbalanced classes. Visualizations like bar plots for tag counts or histograms for sentence lengths reveal potential issues, such as dominant tags or overly long sequences. Outliers, such as unusually lengthy or short sentences, are identified and addressed. Understanding data composition through EDA ensures informed decisions, like setting padding lengths, addressing class imbalance, or augmenting under-represented entities, ultimately improving model training and evaluation.

Model Training:

Fine-tuning BioBERT for Named Entity Recognition (NER) begins by loading the pre-trained BioBERT model, specifically designed for biomedical text, via Bert For Token Classification. This model is tailored for token-level predictions, making it ideal for NER tasks. The number of output labels is adjusted based on the dataset's unique NER tags. Tokenized input sequences are padded to a fixed length (e.g., 75 tokens) to ensure uniformity during training, with corresponding attention masks distinguishing actual tokens from padding. Preprocessed labels are converted to integer indices compatible with the model. To optimize training, the AdamW optimizer is employed alongside a linear learning rate scheduler to manage the learning rate dynamically over epochs. Gradient clipping is applied to stabilize training and prevent gradient explosion.

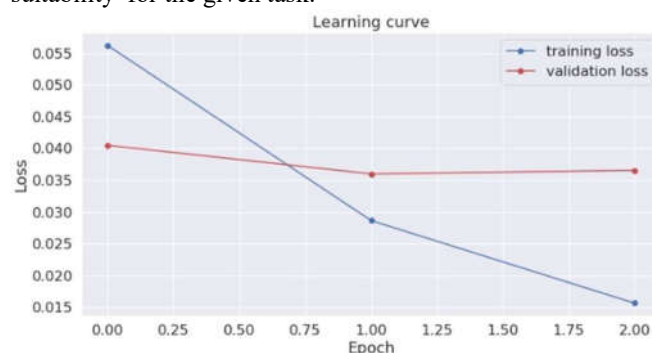
BioBERT is fine-tuned through multiple epochs, typically ranging from 3 to 5, depending on the dataset's size and

complexity. During each training iteration, tokenized input sequences, attention masks, and labels are passed through the model. The output logits represent the likelihood of each token belonging to a specific entity class. Cross-entropy loss is computed by comparing predictions with true labels, while ignoring padded tokens. Gradients are calculated via back propagation, and weights are updated using the optimizer. To improve generalization and reduce overfitting, a portion of the data (e.g., 10%) is reserved for validation. After each epoch, the model is evaluated on this validation set, tracking metrics like loss, accuracy, and F1-score to monitor progress and detect potential overfitting.

Once training is complete, the fine-tuned BioBERT model is evaluated on a test dataset to measure its real-world performance. Common metrics for NER include precision, recall, and F1-score, evaluated at the token and entity levels. Predictions are decoded back to their original tags to assess accuracy qualitatively and quantitatively. Misclassifications and low-performing entity classes are analyzed to refine the model further, often through techniques like class weighting, data augmentation, or hyperparameter tuning. The final model serves as a robust tool for identifying entities in biomedical text, benefiting downstream tasks like clinical data mining, literature curation, or drug discovery.

IV. RESULTS

The loss curve provides insights into the model's training dynamics and convergence. A steady decrease in training loss over epochs indicates that the model is learning effectively from the data. Similarly, a gradual decline in validation loss suggests good generalization to unseen data. However, if the validation loss plateaus or increases while the training loss continues to decrease, it may signal overfitting. An optimal curve reflects both training and validation losses converging to low values with minimal divergence. Consistently low and stable loss values across datasets validate the model's reliability, confirming its suitability for the given task.



0	by
0	the
0	referee
0	.
0	To
0	appear
0	in
0	Annals
0	of
0	Combinatoric
B-indications	Pasteurellosis
0	in
0	japanese
0	quail
0	(
0	Coturnix
0	coturnix

The token-level labeling ensures precise entity recognition, as seen in the association of "Pasteurellosis" with its biological significance. The sequence tagging captures both entity boundaries and their respective contexts, which is vital for downstream tasks like Named Entity Recognition (NER). The presence of domain-specific terms such as "Pasteurellosis" and "Coturnix" highlights the specialized nature of the dataset, tailored for biomedical or zoological applications. Non-entity tokens labeled as O maintain sentence continuity, ensuring model input consistency. This meticulous labeling facilitates the training of models like BioBERT, enabling them to distinguish domain-relevant entities effectively. The dataset's rich context, with precise indications and organism mentions, underscores its utility in identifying disease associations and species interactions. The structured annotations provide a robust foundation for exploring relationships between medical conditions and biological subjects, promoting advancements in biomedical NER tasks.

V. Conclusion:

BioBERT has proven to be a transformative tool for Named Entity Recognition (NER) in the healthcare domain, addressing the complexities of biomedical language with unparalleled accuracy and contextual understanding. By leveraging pre-trained knowledge on extensive biomedical corpora, BioBERT outperforms traditional methods and generic language models in extracting meaningful entities from unstructured healthcare data. This study highlights its potential to streamline clinical decision-making, medical

research, and data-driven applications. While challenges like computational demands and annotation scarcity persist, ongoing advancements in model optimization and domain adaptation promise to unlock broader applications. BioBERT represents a significant step toward smarter, data-driven healthcare innovations.

VI. Future Scope:

The application of BioBERT for Named Entity Recognition (NER) in healthcare holds immense potential for expansion. Future research could focus on fine-tuning BioBERT for specific subdomains, such as genomics or radiology, to address specialized use cases. Integrating BioBERT with clinical decision support systems and knowledge graphs could enhance real-time applications in diagnostics and treatment planning. Additionally, exploring multilingual adaptations of BioBERT could benefit non-English biomedical data analysis. Overcoming challenges like computational overhead and limited annotated datasets through techniques such as model distillation and active learning can further optimize its usability. This progress will enable broader adoption in healthcare analytics and personalized medicine.

VII. References:

1. Aviles M., Rodríguez-Reséndiz J., Ibrahim D. (2023). Optimizing EMG classification through metaheuristic algorithms. *Technologies*, 11, 87.
2. Bentivogli L., Forner P., Giuliano C., Marchetti A., Pianta E., Tymoshenko K. (2010). Extending English ACE 2005 corpus annotation with ground-truth links to Wikipedia.
3. Cai X., Dong S., Hu J. (2019). A deep learning model incorporating part of speech and self-matching attention for named entity recognition of Chinese electronic medical records. *BMC Medical Informatics and Decision Making*, 19, 65.
4. Cui L., Wu Y., Liu J., Yang S., Zhang Y. (2021). Template-based named entity recognition using BART.
5. Dai Z., Wang X., Ni P., Li Y., Li G., Bai X. (2019). Named entity recognition using BERT BiLSTM CRF for Chinese electronic health records.
6. Das S. S. S., Katiyar A., Passonneau R. J., Zhang R. (2021). Container: Few-shot named entity recognition via contrastive learning.
7. Doddington G. R., Mitchell A., Przybocki M. A., Ramshaw L. A., Strassel S. M., et al. (2004). The automatic content extraction (ACE) program: Tasks, data, and evaluation.
8. Dong X., Chowdhury S., Qian L., Li X., Guan Y., Yang J., et al. (2019). Deep learning for named entity recognition on Chinese electronic medical records: Combining deep

Transfer learning with multitask bi-directional LSTM-RNN.

9. Ensastiga S.A.L.,Rodríguez-Reséndiz J.,Estévez-Bén A. A. (2022). Speed controller-based fuzzy logic for a biosignal-feedbacked cycloergometer.

10. Gligic L., Kormilitzin A., Goldberg P., Nevado-Holgado A. (2020). Named entity recognition in electronic health records using transfer learning bootstrapped neural networks.

11. Grancharova M., Dalianis H. (2021). Applying and sharing pre-trained BERT models for named entity recognition and classification in Swedish electronic patient records.

12. Hao Y., Cao H. (2020). An attention mechanism to

Classify multivariate time series.

13. Haverkos B. M., Pan Z., Gru A. A., Freud A. G., Rabinovitch R., Xu-Welliver M., et al. (2016). Extranodal NK/T cell lymphoma, nasal type (ENKTL-NT): An update on epidemiology, clinical presentation, and natural history in North American and European cases.

14. He H., Xin B., Ikehata S., Wipf D. (2017). From Bayesian sparsity to gated recurrent nets.

15. Ji B., Liu R., Li S., Yu J., Wu Q., Tan Y., et al. (2019). A hybrid approach for named entity recognition in Chinese electronic medical records.

16. Ji H., Pan X., Zhang B., Nothman J., Mayfield J., McNamee P., et al. (2017). Overview of TAC-KBP 2017 13 languages entity discovery and linking.