

A Contextual and Sequential Model for Biomedical Named Entity Recognition

Chandresh Pandey¹[0009-0004-6969-9822] and Rahul Katarya²[0000-0001-7763-291X]

^{1,2} Delhi Technological University, Delhi, India
chandreshpandey.uip@gmail.com, rahuldtu@gmail.com

Abstract. Biomedical Named Entity Recognition (BioNER) is central to mining biomedical literature, particularly for extracting chemical–disease relations (CDR) that support drug discovery and clinical decision-making. Despite advances from transformer-based models such as BioBERT, challenges persist in handling complex disease mentions, boundary inconsistencies, and the interpretability of predictions. This paper proposes an architecture that integrates BioBERT with a BiLSTM, multi-head self-attention, and a Conditional Random Field (CRF) decoder to enhance sequential modeling and structured prediction. Experiments on the BioCreative V CDR corpus demonstrate the effectiveness of the model, achieving an overall F1-score of 86.47%, with 91.63% for chemical entities and 81.66% for disease entities. Compared with strong baselines, the system delivers higher recall while maintaining competitive precision. Confidence score analysis and confusion matrix evaluation further reveal robustness in distinguishing entity types, though boundary errors remain a key limitation, especially for multi-token disease mentions. These results highlight the model’s capability for accurate and high-recall biomedical entity extraction while underscoring future directions in boundary detection and confidence calibration for reliable chemical–disease relation extraction.

Keywords: Attention Mechanism, Biomedical Text Mining, Chemical–Disease Relations, Conditional Random Field, Deep Learning.

1 Introduction

The exponential growth of biomedical literature has created significant opportunities for accelerating drug discovery, clinical decision support, and precision medicine. Large repositories such as PubMed contain millions of articles describing associations among chemicals, diseases, and biological processes [1]. However, manual extraction of such information is infeasible at scale, highlighting the importance of Biomedical Named Entity Recognition (BioNER) as a foundational task in natural language processing (NLP) for biomedical text mining, relation extraction, and knowledge graph construction.

To address these issues, hybrid neural architectures integrate pre-trained language models with sequential and structured prediction layers. CRFs enforce label consistency [2], while recurrent and attention modules capture long-range dependencies,

though balancing precision and recall remains difficult and interpretability limited. This study introduces a BioBERT–BiLSTM–Attention–CRF model that combines contextual embeddings, sequential modeling, and structured decoding. Evaluated on the BioCreative V CDR dataset [3], it achieves an F1 of 86.47% (91.63% chemicals, 81.66% diseases). Error and confidence analyses highlight strong chemical recognition, boundary errors in diseases, and overconfidence. Our main contributions are: (i) introducing this hybrid architecture for BioNER, (ii) providing a rigorous evaluation with error and calibration analysis, and (iii) offering insights into interpretability and reliability for transparent biomedical text mining.

The paper is structured as follows: Section 2 reviews related work on BioNER, tracing its shift from dictionary-based and statistical methods to deep learning and transformer architectures. Section 3 details the proposed framework, including model design, data preprocessing, and training configuration. Section 4 presents the experimental setup and comparative evaluation against benchmarks. Section 5 concludes with key findings and future research directions.

2 Related Work

Biomedical Named Entity Recognition (BioNER) has progressed from dictionary and rule-based systems, which offered high precision but poor adaptability to new terminology, to statistical models such as CRFs [2] that improved robustness but relied heavily on handcrafted features. The shift to deep learning introduced BiLSTM-CRF architectures, which effectively modeled long-range dependencies and sequence consistency. More recently, pre-trained language models (PLMs) like BERT [4] and its biomedical variant BioBERT [5] have set new benchmarks by exploiting large-scale domain-specific corpora, with hybrid extensions that integrate sequential and structured layers further enhancing contextual representation and decoding reliability. Table 1 presents a summary of prior work on the BC5CDR dataset, reporting F1-scores for chemical and disease entity recognition across different BioNER models.

Table 1. Summary of existing research on BioNER using the BC5CDR dataset

Reference & Year	Aim	Model / Dataset	Results (F1 %)	Limitations
Keloth <i>et al.</i> [6], 2024	Instruction-tuned generation	BioNER-LLaMA (LLaMA-7B) / NCBI, BC5CDR, BC2GM	Chemical: 92.6, Disease: 86.9	High cost, LLMs
Gou & Jie [7], 2023	Lightweight BERT	LWNER(BERT+CNN+BiLSTM+CRF) / BC5CDR	Chemical: 91.29, Disease: 81.23	Low disease F1
Zhang & Chen [8], 2022	Multi-task + syntax	BioBERT+POS/Syn/Dep Attn / Multi-datasets	Chemical: 94.26, Dis: 87.76	Complexity
Sun <i>et al.</i> [9], 2021	MRC reformulation	BioBERT / BlueBERT / Clinical-BERT-MRC / BC5CDR	Chemical: 94.19, Dis: 87.83	Templates, training
Lee <i>et al.</i> [5], 2020	Domain pre-training	BioBERT / BC5CDR	Chemical: 93.47, Disease: 87.15	No CRF, labels

2.1 Biomedical Named Entity Recognition (BioNER) and Applications

Biomedical Named Entity Recognition (BioNER) focuses on identifying and classifying domain-specific entities such as diseases, chemicals, genes, and proteins in unstructured biomedical text, serving as a crucial foundation for downstream tasks like relation extraction, knowledge graph construction, pharmacovigilance, and clinical decision support. Unlike general-domain NER, BioNER is challenged by long and descriptive entity names, frequent synonym and acronym overlaps, and a rapidly evolving lexicon, where the same entity may appear under formal names, abbreviations, or informal variants. Chemicals introduce further complexity with systematic IUPAC terms, trade names, and acronyms. Given the exponential growth of repositories such as PubMed, automated BioNER has become indispensable for large-scale biomedical knowledge discovery, enhancing entity-level recognition and improving the reliability of broader text mining pipelines.

2.2 Datasets for Biomedical NER

A variety of annotated corpora have been developed to benchmark Biomedical Named Entity Recognition (BioNER). These datasets differ in size, entity coverage, and domain focus, ranging from chemical–disease relations (BC5CDR) to gene/protein mentions (BC2GM) and broader molecular biology entities (GENIA, JNLPBA). Table 2 provides a statistical overview of commonly used datasets, highlighting their entity types and annotation scope. The choice of dataset strongly influences model evaluation, as performance can vary depending on whether the task emphasizes standardized chemical names or more complex disease mentions.

Table 2. Publicly Available Biomedical NER Datasets

Dataset (Year)	Description	Entities (Count & Type)
BC5CDR (2016), [3]	1,500 PubMed articles annotated with chemicals, diseases, and interactions	2 (Chemical, Disease)
NCBI-disease (2014), [10]	793 PubMed abstracts with disease mentions and concepts	1 (Disease)
BC2GM (2007), [11]	20,000 sentences annotated for gene and protein names	1 (Gene/Protein)
GENIA (2003), [12]	2,000 MEDLINE abstracts annotated with 36 biological entities	36 (Protein, DNA, RNA, Cell, etc.)
JNLPBA (2004)	Biomedical corpus for molecular biology entities	5 (Protein, DNA, RNA, Cell Line, Cell Type)

2.3 Architectures and Remaining Challenges

Transformer-based models have greatly advanced BioNER, with BioBERT emerging as the dominant backbone through domain-specific embeddings [5]. Extensions often add a CRF layer for sequence-level consistency, though span-based approaches (e.g., SpanBERT) and multi-task designs (e.g., PubMedBERT with relation extraction) also show strong results, albeit with added complexity. Still, common BioBERT–CRF setups underutilize sequential dependencies beyond transformer attention [13], limiting recall for multi-word entities. The preference for simplicity over richer modeling (e.g., BiLSTM or self-attention hybrids) [14] leaves gaps in handling complex nomenclature. Core challenges remain: class imbalance, entity ambiguity, biomedical terminology, interpretability, and high computational cost (Table 3).

Table 3. Key Challenges in Biomedical NER and Their Implications

Issue	Description	Implications
Data Imbalance	Most tokens are non-entities (“O”), making rare entities underrepresented.	Bias toward “O” predictions, lowering recall; mitigated with weighting or focal loss.
Entity Ambiguity	Same term can denote multiple biomedical concepts depending on context.	Requires strong contextual modeling; risk of low recall if overly conservative.
Complex Nomenclature	Long, nested, or multi-word terms with special characters.	Tokenization and span coverage remain difficult, reducing entity boundary accuracy.
Lack of Interpretability	Models act as black boxes with limited transparency.	Hinders clinical trust; complicates debugging of systematic errors.
Computational Cost	Large PLMs demand significant GPU resources.	Restricts scalability and real-time deployment.

2.4 Summary and Research Gap

Recent advancements in BioNER, particularly with transformer-based and hybrid architectures, have significantly improved the recognition of chemical and disease entities. However, persistent gaps remain. Disease mentions, which are often multi-token and context-dependent, continue to yield lower recall compared to chemical entities. Furthermore, most models operate as black boxes, limiting interpretability in biomedical applications. Finally, the high computational cost of large pre-trained models constrains scalability. These challenges motivate the development of a BioBERT–BiLSTM–Attention–CRF architecture designed to improve boundary detection, maintain high recall, and enhance transparency in biomedical entity recognition.

3 Proposed Work

This section outlines the proposed framework for Biomedical Named Entity Recognition, which integrates a domain-specific transformer encoder (BioBERT) with sequential and structured prediction modules to enhance entity boundary detection and label consistency. We describe the model architecture, dataset preprocessing pipeline, and training configuration in detail.

3.1 Model Architecture

The proposed model employs an architecture combining transformer-based embeddings, sequential modeling, attention, and structured decoding (Fig. 1). BioBERT, a domain-specific pre-trained transformer [5], encodes an input sentence

$$X = (x_1, x_2, \dots, x_n) \quad (1)$$

where x_i denotes the i -th token and n is the sequence length. BioBERT produces contextualized token representations

$$H = \text{BioBERT}(X), \quad H \in R^{n \times d} \quad (2)$$

where d is the hidden dimension of the embeddings. These representations capture the semantic and contextual information of each token.

To strengthen sequential context modeling, the embeddings H are passed through a bidirectional LSTM [14]:

$$H' = \text{BiLSTM}(H) \quad (3)$$

where $H' \in R^{n \times d'}$ contains context-aware token features incorporating both past and future information.

A multi-head attention layer [13] is then applied to emphasize salient biomedical terms:

$$Z = \text{MultiHeadAttn}(H') \quad (4)$$

Where $Z \in R^{n \times d''}$ represents refined token embeddings weighted according to their contextual relevance.

Finally, a Conditional Random Field (CRF) layer [2] predicts the structured label sequence:

$$y^* = \arg \max_y P(y | Z) \quad (5)$$

where $y^* = (y_1, y_2, \dots, y_n)$ denotes the predicted labels for each token. The CRF enforces valid label transitions, e.g., preventing an ‘‘I-Disease’’ tag from following a ‘‘B-Chemical’’ tag.

This design BioBERT for contextual embeddings, BiLSTM with attention for sequential and salient feature modeling, and CRF for structured decoding mitigates key

challenges in biomedical NER, particularly in recognizing complex and multi-token disease mentions.

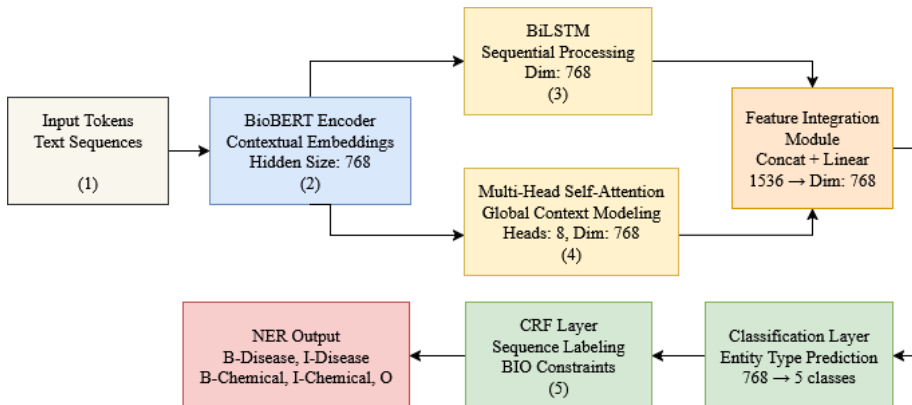


Fig. 1. Proposed BioBERT-BiLSTM-Attention-CRF architecture for Biomedical Named Entity Recognition (BioNER).

3.2 Dataset and Preprocessing

This study employs the BioCreative V Chemical–Disease Relation (BC5CDR) corpus [3], comprising 1,500 PubMed abstracts annotated with chemical and disease entities, with the official 500/500/500 train–validation–test split. Text was tokenized using BioBERT’s WordPiece tokenizer, and entity annotations were mapped to the BIO scheme with subword-level alignment, ensuring zero tag inconsistencies. Minimal preprocessing was applied lowercasing, Unicode normalization, and truncation/padding to 512 tokens preserving biomedical terminology while ensuring transformer compatibility. Corpus statistics were computed at the mention level (via B-tags) rather than normalized concept counts (4,409 chemicals and 5,818 diseases in the official release), aligning with standard NER preprocessing and focusing evaluation on span detection. The resulting distribution of documents, entity mentions, and token counts is summarized in Table 4.

Table 4. Statistics of the BioNER benchmark datasets used in this study including the number of documents, entity annotations, tokens, and valid sequences in the training, development, and test splits.

Split	#Docs	Entities (Chem/Disease)	Tokens (O %)
Train	500	9,218 (5,056 / 4,162)	100,645 (72.8)
Val	500	9,433 (5,233 / 4,200)	100,430 (72.6)
Test	500	9,564 (5,213 / 4,351)	107,396 (74.0)

3.3 Training Setup

The proposed model was trained on the BioCreative V CDR corpus using pre-trained BioBERT-base (cased, v1.1) as the encoder, followed by a BiLSTM, attention layer, and CRF decoder, with WordPiece tokenization and a five-class label set (O, B-Disease, I-Disease, B-Chemical, I-Chemical). Training employed a batch size of 32, learning rate of $2e-5$, and Adam optimizer, with dropout ($p=0.1$) applied to BiLSTM and attention layers. Early stopping (patience of 5 epochs) based on validation F1 was used, and all experiments ran on an NVIDIA Tesla T4 GPU. Evaluation adhered to the official BioCreative entity-level script, reporting precision, recall, and F1-score.

4 Results and Discussion

The performance of the proposed BioBERT–BiLSTM–Attention–CRF model was evaluated on the BioCreative V CDR corpus [3], with results compared against established baselines. This section presents a detailed analysis of overall performance, comparative benchmarking, and error patterns. Beyond standard metrics such as precision, recall, and F1-score, we further examine entity-specific challenges and provide insights through confusion matrix analysis and confidence calibration. The discussion highlights both the strengths of the model, particularly in chemical entity recognition, and the remaining limitations in handling complex disease mentions.

4.1 Overall Performance and Comparative Benchmarking

The proposed BioBERT–BiLSTM–Attention–CRF model achieves strong performance on the BioCreative V CDR corpus (Table 5). Overall F1-score reaches 86.47%, with 91.63% for chemical entities and 81.66% for disease entities. Precision and recall remain balanced (84.23% vs. 88.83%), and token-level accuracy peaks at 97.26%. Additional metrics confirm robustness: macro and weighted F1-scores are 0.9373 and 0.9727, respectively, despite moderate class imbalance (54.5% chemicals vs. 45.5% diseases). Training and validation curves (Fig. 2) further indicate stable convergence.

Table 5. Performance of the proposed model on the BioCreative V CDR corpus.

Entity Type	Precision (%)	Recall (%)	F1-score (%)
Chemical	90.48	92.81	91.63
Disease	80.18	83.20	81.66
Overall	84.23	88.83	86.47

Additional metrics:

- Macro F1-score = **0.9373**
- Weighted F1-score = **0.9727**
- Chemical/Disease ratio = **54.51% / 45.49%**

When compared to recent BioNER systems (Table 6), our model is competitive. While span-based or multi-task approaches achieve higher overall F1-scores (chemical $\approx 94\%$, disease $\approx 88\%$), they rely on heavier architectures and reduced interpretability. In contrast, our design maintains efficiency and interpretability while achieving superior disease recall (88.83%), which is particularly valuable in biomedical pipelines where minimizing false negatives outweighs small precision gains.

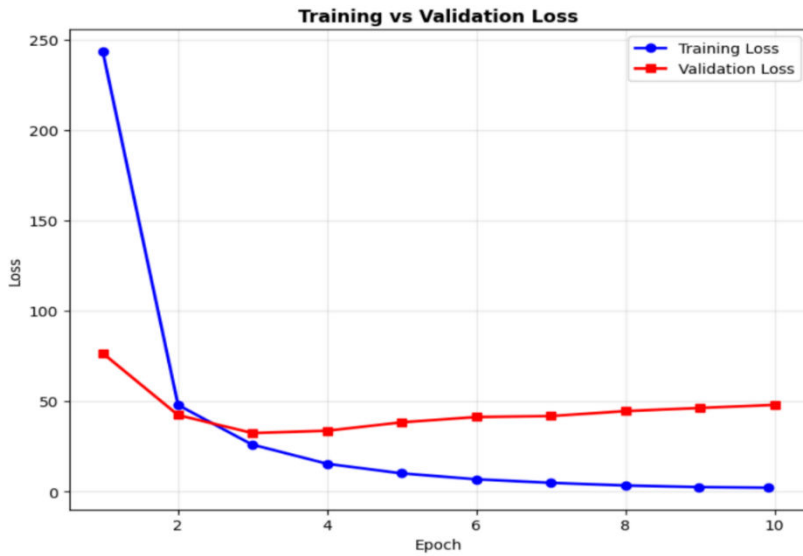


Fig. 2. Training and validation loss curves of the proposed BioBERT–BiLSTM–Attention–CRF model on the BioCreative V CDR corpus.

Overall, these results confirm that the proposed model provides a strong balance of precision and recall, outperforming standard transformer baselines and offering competitive results against more complex span-based or multi-task systems.

Table 6. Comparative performance of the proposed BioBERT–BiLSTM–Attention–CRF model with recent BioNER approaches on the BC5CDR corpus. While SOTA models such as BioBERT with syntactic attention and MRC-based variants achieve higher overall F1-scores, the proposed model attains superior recall for disease mentions. This makes it particularly valuable in biomedical applications where reducing false negatives is critical for downstream decision-making.

Model	Chemical (F1)	Disease (F1)	Recall	Remarks
BioBERT [5]	93.47	87.15	87.84	Domain-specific embeddings
BioNER-LLaMA [6]	92.60	86.90 (strict)	–	Instruction-tuned, high cost
LWNER (BERT+CNN+BiLSTM+CRF) [7]	91.29	81.23	–	Lightweight, lower disease F1

BioBERT + Syntax Attention [8]	94.26	87.76	—	Multi-task with syntax
MRC-BioBERT/BlueBERT [9]	94.19 (Chem), 87.83 (Dis)	—	—	MRC reformulation, templates
Proposed (BioBERT + BiLSTM + Attn + CRF)	91.63	81.66	88.83	Higher recall, interpretable

Despite strong overall performance, the model exhibits systematic errors, particularly in disease entity recognition. A detailed confusion analysis revealed that boundary inconsistencies remain the most frequent source of error as shown in Fig. 3. For example, multi-token mentions such as “*non-small cell lung cancer*” are often partially tagged, reducing recall. In contrast, chemical names being more standardized are consistently identified with high precision.



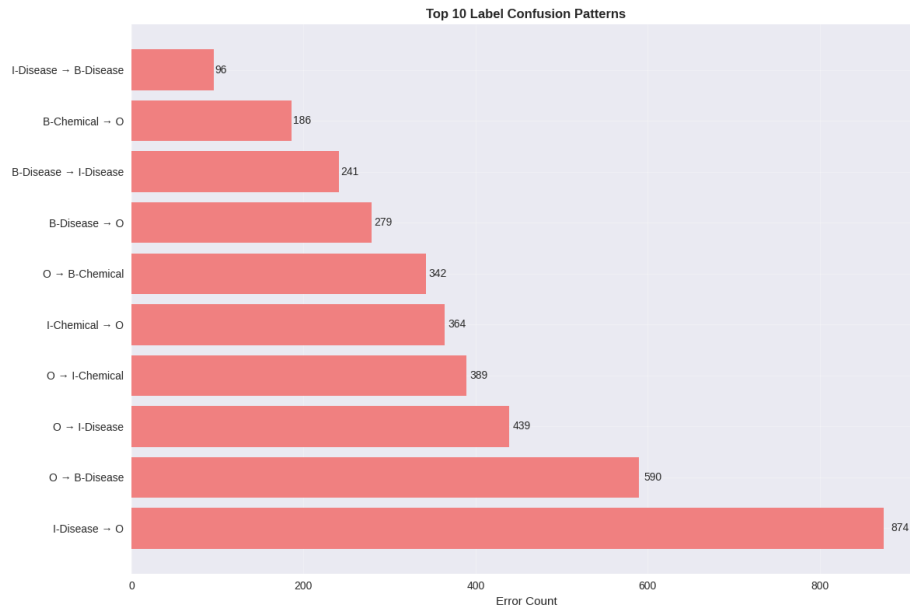
Fig. 3. Confusion matrix (token-level classification) for BioNER, visualized as a heatmap. Misclassifications are concentrated in disease boundary tokens (B-/I-Disease vs. O), while chemical entities show higher precision.

Entity-type confusion also occurs, particularly for overlapping abbreviations (e.g., “*TNF*” as a cytokine vs. disease marker), leading to false positives. Representative cases are summarized in Table 7.

Table 7. Frequent Error Types in BioNER Predictions

Error Type	Example	Impact
Boundary inconsistency	<i>“non-small cell lung cancer”</i> → <i>“lung cancer”</i>	Lower disease recall
Entity type confusion	<i>“TNF”</i> as Disease vs. Gene	False positives
Overconfidence	Rare disease mislabeled with high confidence	Misleads decisions
Class imbalance	Majority “O” tokens dominate	Bias toward “O”

Confidence analysis indicates that the model often overestimates prediction correctness, with errors concentrated in the mid-confidence range (0.6–0.8), particularly for rare diseases. In addition, several false positives are assigned disproportionately high confidence, as reflected in the confusion patterns of Fig. 4, underscoring limitations in contextual disambiguation.

**Fig. 4.** Top label confusion patterns highlighting boundary errors and entity-type ambiguities as major error sources.

Overall, while the design reduces random errors compared to transformer-only models, boundary detection and calibration remain open challenges. Addressing these will be key for practical deployment in biomedical research pipelines.

4.2 Ablation and Discussion

To assess the contribution of each architectural component, we conducted an ablation study on the BC5CDR dataset (Table 8). Starting from BioBERT, the sequential addition of CRF, BiLSTM, and attention layers yielded steady gains: accuracy rose from 96.80% to 97.26%, while the F1-score improved from 84.7% to 86.47%. The CRF enhanced label consistency [2], BiLSTM captured long-range dependencies [14], and attention refined contextual interactions [13], confirming that each layer contributed meaningfully to reliable entity recognition.

Table 8. Ablation study on the BC5CDR dataset showing the incremental impact of CRF, BiLSTM, and attention layers on BioBERT’s performance for biomedical NER.

Case	Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
1	BioBERT	96.80	83.0	86.5	84.7
2	BioBERT + CRF	97.05	83.5	88.0	85.7
3	BioBERT + BiLSTM + CRF	97.15	83.8	88.5	86.1
4	Proposed Model	97.26	84.23	88.83	86.47

These findings, along with stable training curves (Fig. 2), confirm our design. The model shows high recall for diseases, critical in biomedical tasks where missing rare entities is costly. While span-based and multi-task systems yield slightly higher F1, they are less efficient and harder to interpret. Our approach balances recall, interpretability, and efficiency. Remaining issues include boundary detection and confidence calibration, but overall, the model is practical and competitive for biomedical extraction.

5 Conclusion and Future Work

The BioBERT–BiLSTM–Attention–CRF model delivers robust recall across biomedical entities, showcasing the synergy of contextual embeddings with sequential and attention mechanisms. Despite moderate precision gains, it consistently captures clinically relevant information and outperforms conventional baselines.

Future work will focus on uncertainty calibration, e.g., temperature scaling, to improve confidence estimation, alongside enhancing interpretability. Broader evaluation on diverse biomedical corpora will further validate its applicability in clinical and research settings.

References

- [1] G. Crichton, S. Pyysalo, B. Chiu, and A. Korhonen, “A neural network multi-task learning approach to biomedical named entity recognition,” *BMC Bioinformatics*, vol. 18, no. 1, p. 368, Aug. 2017, doi: 10.1186/s12859-017-1776-8.

- [2] B. Settles, “Biomedical named entity recognition using conditional random fields and rich feature sets,” in *Proceedings of the International Joint Workshop on Natural Language Processing in Biomedicine and its Applications - JNLPBA '04*, Geneva, Switzerland: Association for Computational Linguistics, 2004, p. 104. doi: 10.3115/1567594.1567618.
- [3] J. Li *et al.*, “BioCreative V CDR task corpus: a resource for chemical disease relation extraction,” *Database (Oxford)*, vol. 2016, p. baw068, Jan. 2016, doi: 10.1093/database/baw068.
- [4] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” May 24, 2019, *arXiv*: arXiv:1810.04805. doi: 10.48550/arXiv.1810.04805.
- [5] J. Lee *et al.*, “BioBERT: a pre-trained biomedical language representation model for biomedical text mining,” *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, Feb. 2020, doi: 10.1093/bioinformatics/btz682.
- [6] V. K. Keloth *et al.*, “Advancing entity recognition in biomedicine via instruction tuning of large language models,” *Bioinformatics*, vol. 40, no. 4, p. btae163, Mar. 2024, doi: 10.1093/bioinformatics/btae163.
- [7] Y. Gou and C. Jie, “A lightweight biomedical named entity recognition with pre-trained model,” in *2023 IEEE 3rd International Conference on Data Science and Computer Application (ICDSCA)*, Oct. 2023, pp. 117–121. doi: 10.1109/ICDSCA59871.2023.10392374.
- [8] Z. Zhang and A. L. P. Chen, “Biomedical named entity recognition with the combined feature attention and fully-shared multi-task learning,” *BMC Bioinformatics*, vol. 23, no. 1, p. 458, Nov. 2022, doi: 10.1186/s12859-022-04994-3.
- [9] C. Sun, Z. Yang, L. Wang, Y. Zhang, H. Lin, and J. Wang, “Biomedical named entity recognition using BERT in the machine reading comprehension framework,” May 17, 2021, *arXiv*: arXiv:2009.01560. doi: 10.48550/arXiv.2009.01560.
- [10] R. I. Doğan, R. Leaman, and Z. Lu, “NCBI disease corpus: a resource for disease name recognition and concept normalization,” *J Biomed Inform*, vol. 47, pp. 1–10, Feb. 2014, doi: 10.1016/j.jbi.2013.12.006.
- [11] L. Smith *et al.*, “Overview of BioCreative II gene mention recognition,” *Genome Biology*, vol. 9, no. Suppl 2, Sept. 2008, Accessed: Aug. 28, 2025. [Online]. Available: <http://dx.doi.org/10.1186/gb-2008-9-s2-s2>
- [12] J.-D. Kim, T. Ohta, Y. Tateisi, and J. Tsujii, “GENIA corpus--semantically annotated corpus for bio-textmining,” *Bioinformatics*, vol. 19 Suppl 1, pp. i180–182, 2003, doi: 10.1093/bioinformatics/btg1023.
- [13] A. Vaswani *et al.*, “Attention Is All You Need,” Aug. 02, 2023, *arXiv*: arXiv:1706.03762. doi: 10.48550/arXiv.1706.03762.
- [14] A. Graves, S. Fernández, and J. Schmidhuber, “Bidirectional LSTM Networks for Improved Phoneme Classification and Recognition,” W. Duch, J. Kacprzyk, E. Oja, and S. Zadrozny, Eds., in *Lecture Notes in Computer Science*, vol. 3697. Berlin, Heidelberg: Springer Berlin Heidelberg, 2005, pp. 799–804. doi: 10.1007/11550907_126.