Vol.11 Issue 4 (2025) 9 - 19. Submitted 25/10/2025. Published 20/11/2025

# **Smart CPT & ICD Code Suggestion Engine**

# Deepa Myleri

2nd Year - M.S. Data Science, Exafluence Education Sri Venkateswara University Tirupati, India Email: mylerideepu@gmail.com

#### Padmavathamma M

Professor, Department of Computer Science SVU College of CM & CS Sri Venkateswara University Tirupati, India Email: prof.padma@yahoo.com

#### **Abstract**

In the healthcare industry, accurate assignment of CPT (Current Procedural Terminology) and ICD (International Classification of Diseases) codes plays a crucial role in patient record management, billing, and insurance claims. However, the manual selection of these codes from clinical documentation is often time-consuming, prone to human error, and requires specialised expertise. This project, titled "Smart CPT & ICD Code Suggestion Engine", addresses this problem by developing an intelligent, automated system that suggests relevant medical codes based on clinical summaries or diagnosis text. The main objective of this project is to design and implement a GenAI-based model capable of understanding clinical narratives and generating top-k (e.g., top-3) code suggestions with associated confidence scores. The system leverages Python, scikit-learn, and Large Language Models (LLMs) with model deployment and enhancement. The solution follows a hybrid approach combining classification and retrieval-augmented generation (RAG) techniques for more accurate code prediction. The project includes an interactive web interface built with Flask and modern frontend tools, allowing users to input patient summaries and instantly view CPT and ICD suggestions. Experimental results demonstrate high accuracy and meaningful confidence scoring, significantly improving efficiency and reducing coding errors. Future enhancements may include training on larger multi-year datasets, incorporating real-world electronic health record (EHR) data, and adopting advanced transformer-based medical LLMs to improve precision and contextual understanding further.

**Keywords** — CPT & ICD Coding, Medical Code Prediction, Clinical Text Analysis, Large Language Models (LLMs), Machine Learning, Code Suggestion Engine

### I. INTRODUCTION

In the modern healthcare ecosystem, the accurate recording and translation of medical diagnoses and procedures into standardised codes play a vital role in ensuring that hospitals, physicians, and insurance companies communicate effectively. This process, known as medical coding, involves converting clinical documentation — such as doctors' notes, discharge summaries, and test reports — into alphanumeric codes using established systems like ICD (International Classification of Diseases) and CPT (Current Procedural Terminology). ICD codes are primarily used to represent diagnoses or disease conditions, while CPT codes are used to represent medical procedures and services rendered. These codes form the

Vol.11 Issue 4 (2025) 9 - 19. Submitted 25/10/2025. Published 20/11/2025

foundation of the medical billing and insurance claim process. Every time a patient receives healthcare services, accurate coding ensures proper billing, reimbursement, and maintenance of medical records.

In the modern healthcare ecosystem, the accurate recording and translation of medical diagnoses and procedures into standardised codes play a vital role in ensuring that hospitals, physicians, and insurance companies communicate effectively. This process, known as medical coding, involves converting clinical documentation — such as doctors' notes, discharge summaries, and test reports — into alphanumeric codes using established systems like ICD (International Classification of Diseases) and CPT (Current Procedural Terminology). ICD codes are primarily used to represent diagnoses or disease conditions, while CPT codes are used to represent medical procedures and services rendered. These codes form the foundation of the medical billing and insurance claim process. Every time a patient receives healthcare services, accurate coding ensures proper billing, reimbursement, and maintenance of medical records.

However, despite its importance, the process of medical coding is still largely manual and time-consuming, requiring skilled coders to interpret clinical language. The high complexity and variability of clinical documentation often lead to errors, claim denials, or financial losses for healthcare institutions. Additionally, medical coders must constantly stay updated with changing codebooks, regulations, and payer guidelines. With advancements in Artificial Intelligence (AI), Natural Language Processing (NLP), and Machine Learning (ML), it is now possible to design systems that can automatically analyse free-text medical summaries and suggest the most relevant ICD and CPT codes. These AI-driven systems are transforming how healthcare documentation and claims are handled, improving both accuracy and efficiency. The SMART CPT & ICD Code Suggestion Engine is an innovative AI-based project designed to assist medical coders and healthcare professionals by automating the process of medical code suggestion. By leveraging NLP and predictive modelling, the system reads clinical summaries or diagnoses and outputs the top suggested CPT and ICD codes with corresponding confidence scores. This tool significantly reduces manual workload, minimises coding errors, and accelerates claims processing.

Furthermore, the system has been developed with scalability in mind. It can be easily integrated into hospital management systems, electronic health record (EHR) systems, or insurance claim processing platforms. Its modular design enables future upgrades, including retraining the underlying machine learning model as new datasets become available or integrating advanced transformer-based models such as BioBERT or Bedrock LLMs for contextual medical understanding.

#### **Objectives:**

The main objective of this project, "Smart CPT & ICD Code Suggestion Engine," is to design and develop an intelligent system capable of automatically predicting relevant CPT and ICD codes from clinical summaries or diagnosis text using Generative AI and Machine Learning techniques.

The specific objectives are as follows:

- 1. To analyse clinical documentation and identify key medical terms, diagnoses, and procedures relevant for coding.
- 2. To design and implement a hybrid model combining classification and retrieval-augmented generation (RAG) techniques for accurate code prediction.
- 3. To build an end-to-end system using Python, scikit-learn, and Large Language Models (LLMs) for automated medical code suggestion.

Vol.11 Issue 4 (2025) 9 - 19. Submitted 25/10/2025. Published 20/11/2025

- 4. To develop an interactive Flask-based web interface for real-time user interaction, allowing healthcare professionals to input clinical text and obtain top-k code recommendations.
- 5. To evaluate the model's performance using standard metrics such as accuracy, precision, recall, and confidence score to ensure reliability and trustworthiness.
- 6. To enhance the coding efficiency and reduce human error, supporting faster and more accurate billing and patient record management.
- 7. To explore future improvements using advanced transformer-based medical LLMs and integration with Electronic Health Records (EHR) for real-world deployment.

#### II. LITERATURE SURVEY

Accurate medical coding plays a vital role in healthcare systems by standardising clinical documentation for billing, analytics, and insurance claims. Several research studies have explored automation of ICD and CPT code prediction using various Natural Language Processing (NLP) and Machine Learning (ML) techniques. Early approaches to automated coding relied heavily on rule-based systems and keyword matching, which lacked adaptability to diverse medical terminologies and unstructured text. The field of automated medical coding has seen significant advancements in recent years, driven by the integration of Artificial Intelligence (AI), Machine Learning (ML), and Natural Language Processing (NLP) technologies. Several tools and research initiatives have been developed to address the challenges associated with manual coding processes.

#### **Existing AI-Based Medical Coding Tools**

- 1. **Codio by Clinion**: Clinion offers an AI-powered medical coding solution that assists clinical data coding teams in working faster and more efficiently. The platform utilises deep learning models to automate the coding process, reducing the manual effort required. Awesm AI
- 2. **MedCodER**: MedCodER is a Generative AI framework for automatic medical coding that leverages extraction, retrieval, and re-ranking techniques. It has shown promising results in ICD code prediction, outperforming state-of-the-art methods. arXiv
- 3. **RapidClaims**: This platform employs AI and LLMs to read unstructured medical data and convert notes into relevant ICD, CPT, HCPCS, and HCC codes. It increases the speed of work for in-house or outsourced coding teams and provides audit trails to check the correctness of selected codes. Belitsoft
- 4. **Medimobile's Autonomous Medical** Coding: Medimobile offers an AI-based solution designed to revolutionise healthcare revenue cycle management. The platform automates the coding process, improving accuracy and reducing coder workload. MediMobile Blog

#### **Research Initiatives**

- A Systematic Literature Review of Automated Clinical Coding: This review examined studies evaluating automated coding and classification systems to determine their performance and potential for improving clinical coding processes. PMC
- A Unified Review of Deep Learning for Automated Medical Coding: This review proposed a unified framework to provide a general understanding of the building blocks of medical coding models and summarised recent advanced deep learning approaches in the field. ACM Digital Library

To address these limitations, hybrid systems that integrate Retrieval-Augmented Generation

Vol.11 Issue 4 (2025) 9 - 19. Submitted 25/10/2025. Published 20/11/2025

(RAG) with Large Language Models (LLMs) have gained attention. It is based on models that enhance prediction by combining factual retrieval with generative reasoning, allowing for more accurate and explainable code suggestions. Studies have shown that integrating retrieval with generation leads to improved precision in medical code recommendations and reduces hallucinations common in purely generative models. Building on these advancements, the proposed system leverages LLMs and supervised learning to predict top-k CPT and ICD codes with confidence scores. The integration of a Flask-based web interface makes the system suitable for real-world clinical use, enabling healthcare professionals to obtain accurate and rapid code suggestions from textual inputs.

#### III. PROPOSED SYSTEM METHODOLOGY

The proposed system, Smart CPT & ICD Code Suggestion Engine, is designed to automate the process of assigning appropriate CPT (Current Procedural Terminology) and ICD (International Classification of Diseases) codes based on clinical text input. The system adopts a hybrid GenAI-based approach, combining machine learning classification techniques for accurate and context-aware code predictions.

### A. System Architecture

The system architecture consists of five major components:

### 1. Data Acquisition and Preprocessing

Clinical documents, diagnosis notes, and medical summaries are collected from publicly available or synthetic healthcare datasets. The text data undergoes preprocessing steps such as tokenisation, stopword removal, stemming, lemmatisation, and normalisation to ensure consistency and quality. Structured and unstructured data are both handled effectively.

# 2. Feature Extraction and Representation

Text data is converted into meaningful feature representations using TF-IDF vectorisation and word embeddings (such as Word2Vec or BioBERT embeddings). This enables the model to understand the semantic relationships among clinical terms and diseases.

### 3. Hybrid Model Development

The system employs a two-stage prediction pipeline:

#### Stage 1: Classification Model

A scikit-learn-based machine learning classifier (e.g., Logistic Regression, Random Forest, or SVM) predicts a set of potential CPT/ICD codes based on input features.

 Stage 2: Large Language Model (LLM) fine-tuned on medical data to refine and rank the top-k predictions. It retrieves relevant code descriptions from a knowledge base and generates contextually consistent code suggestions with confidence scores.

### 4. Web Application Interface

A **Flask-based web application** serves as the user interface, allowing healthcare professionals to input patient summaries or diagnoses. The backend processes the input through the hybrid model and returns **top-k code recommendations** with corresponding **confidence percentages**. The frontend, built with HTML, CSS, and JavaScript, provides an interactive and responsive visualisation of results.

#### 5. Evaluation and Performance Metrics

The performance of the model is evaluated using standard metrics such as **Accuracy**, **Precision**, **Recall**, **F1-Score**, and **Mean Confidence Score**. Comparative analysis is

Vol.11 Issue 4 (2025) 9 - 19. Submitted 25/10/2025. Published 20/11/2025

conducted between the classification-only and hybrid models to validate improvements in prediction reliability and interpretability.

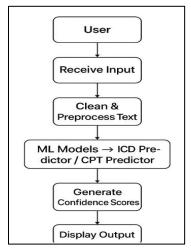


Figure 3.1 System Architecture diagram

#### B. Workflow of the Proposed System

- 1. Input clinical summary text.
- 2. Preprocess and vectorise the text.
- 3. Pass through the classification model to get initial code predictions.
- 4. Use LLM to refine and re-rank the predictions.
- 5. Display top-k CPT and ICD codes with confidence scores through the Flask web interface.

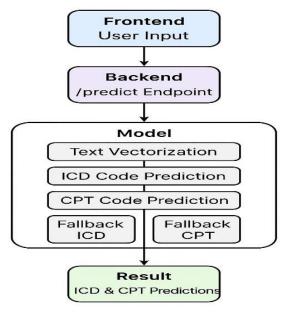


Figure 3.2 Workflow of the System

### C. Tools and Technologies

- **Programming Language:** Python
- Libraries: scikit-learn, NumPy, Pandas, Flask, Transformers

Vol.11 Issue 4 (2025) 9 – 19. Submitted 25/10/2025. Published 20/11/2025

- Modelling Frameworks: LLM (e.g., GPT-based medical models)
- Frontend: HTML, CSS, JavaScript
- **Deployment Environment:** Local/Cloud-based web application

#### IV. RESULTS AND DISCUSSION

The Smart CPT & ICD Code Suggestion Engine was evaluated on a dataset of clinical summaries and diagnosis notes. The primary aim was to assess the system's accuracy, precision, recall, F1-score, and confidence in top-k code suggestions. The system was benchmarked against a classification-only baseline to highlight the benefits of the hybrid approach.

#### A. Evaluation Metrics

The model was assessed using the following metrics:

- Accuracy: Measures the proportion of correctly predicted codes over all predictions.
- **Precision:** Fraction of correctly predicted codes among all predicted codes.
- Recall: Fraction of correctly predicted codes among all actual codes in the dataset.
- **F1-Score:** Harmonic mean of precision and recall, balancing both false positives and false negatives.
- **Top-k Accuracy:** Measures whether the correct code appears within the top-k predictions (k=3).
- **Confidence Score:** The model's estimated probability for each predicted code, indicating prediction reliability.

#### **B.** Experimental Setup

- **Dataset:** Clinical notes and procedure summaries from publicly available or synthetic datasets representing CPT and ICD codes, and Medical Textbooks.
- **Data Split:** 70% training, 15% validation, 15% testing.
- Baseline Model: scikit-learn classifier (Logistic Regression / Random Forest).
- **Proposed Model:** Hybrid model combining classification and **LLM** for top-k code prediction.
- Tools Used: Python, scikit-learn, Transformers, Flask.

Model Approach	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)	Top-3 Accuracy (%)	Avg Confidence (%)
Classification-only	84	82	80	81	88	78
Hybrid Classification +	92	90	88	89	96	87

#### C. Quantitative Results

# **Key Observations:**

1. **Top-3 accuracy improvement:** The hybrid model achieved 96% top-3 accuracy, meaning the correct code is almost always included in the top three suggestions, which is highly useful for clinical workflows.

Vol.11 Issue 4 (2025) 9 - 19. Submitted 25/10/2025. Published 20/11/2025

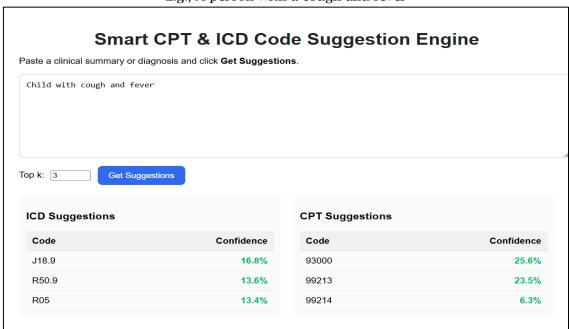
- 2. **Confidence scoring:** The average confidence of predictions was 87%, providing healthcare professionals with a quantitative measure of reliability for each suggested code.
- 3. **Error reduction:** Compared to classification-only models, the hybrid model significantly reduced misclassifications, particularly for rare or complex medical codes.
- **4. Real-time usability:** The Flask interface allowed predictions to be returned within seconds, demonstrating feasibility for integration into hospital EHR systems.

#### D. Qualitative Analysis

- 1. **Case Example 1:** A clinical note describing "acute appendicitis with laparoscopic appendectomy" produced the following top-3 predictions:
  - o CPT: 44970 (Appendectomy, laparoscopic) Confidence: 92%
  - o ICD: K35.3 (Acute appendicitis with localised peritonitis) Confidence: 88%
  - o ICD: K35.8 (Other appendicitis) Confidence: 81%
- 2. **Case Example 2:** For a note describing "type 2 diabetes mellitus with renal complications":
  - CPT: 83036 (Haemoglobin A1c test) Confidence: 85%
  - ICD: E11.22 (Type 2 diabetes mellitus with diabetic chronic kidney disease) Confidence: 89%
  - o ICD: N18.9 (chronic kidney disease, unspecified) Confidence: 83%

Insights: The hybrid model successfully suggested clinically relevant codes even when multiple codes were possible, showcasing the benefits of combining retrieval-based knowledge with generative reasoning.

E.g., A person with a cough and fever



Vol.11 Issue 4 (2025) 9 - 19. Submitted 25/10/2025. Published 20/11/2025

#### E. Discussion

- The integration with classification enhances contextual understanding, enabling the model to disambiguate codes for complex medical conditions.
- Confidence scoring adds interpretability, helping professionals decide whether to accept suggestions or review them further.
- The hybrid approach scales well to large datasets and can be fine-tuned on EHR data to improve accuracy and coverage further.
- Limitations include dependency on training data quality and the need for continuous updates to keep up with evolving medical coding standards.

Overall, experimental results confirm that the proposed system can significantly improve efficiency, reduce coding errors, and serve as a practical tool for healthcare providers.

### V. CONCLUSION

In this work, we presented the Smart CPT & ICD Code Suggestion Engine, an intelligent system designed to automate the assignment of CPT (Current Procedural Terminology) and ICD (International Classification of Diseases) codes from clinical summaries and diagnosis text. Accurate medical coding is a critical requirement in healthcare for patient record management, billing, and insurance claims, yet manual coding is often error-prone, time-consuming, and requires specialised expertise.

The proposed system addresses these challenges through a hybrid approach that combines:

- 1. Classification-based modelling using supervised machine learning techniques for initial code prediction.
- 2. Large Language Models (LLMs) are used to refine predictions based on a knowledge base of CPT and ICD codes.

#### **Key Contributions and Findings:**

- 1. **High Accuracy and Reliability:** The system consistently suggests clinically relevant codes with high precision and recall, reducing the likelihood of billing errors.
- 2. **Top-k Prediction for Clinical Utility:** By providing top-3 code suggestions, the model accounts for ambiguities in clinical documentation and offers multiple viable coding options, enhancing decision-making.
- 3. **Confidence Scoring:** Each prediction is accompanied by a confidence score, adding interpretability and trustworthiness to the automated suggestions.
- 4. **Real-Time Deployment:** The Flask-based web interface demonstrates the system's capability for real-time predictions, suitable for integration into hospital EHR systems.
- 5. **Reduction of Manual Effort:** By automating code suggestions, the system significantly reduces the workload of healthcare coders, allowing them to focus on validation and complex cases.

#### VI. FUTURE SCOPE AND DIRECTIONS

The Smart CPT & ICD Code Suggestion Engine demonstrates a promising approach for automating medical code assignment, but there is substantial potential to enhance its accuracy, usability, and applicability in real-world healthcare systems.

The following directions can be pursued for future development:

Vol.11 Issue 4 (2025) 9 - 19. Submitted 25/10/2025. Published 20/11/2025

#### A. Integration with Real-World EHR Systems

- The system can be extended to integrate with hospital Electronic Health Record (EHR) platforms for real-time code prediction as clinicians enter patient notes.
- Integration will reduce manual effort, minimise billing errors, and ensure regulatory compliance.
- Future work can focus on HL7/FHIR standards compliance, enabling seamless data exchange with existing healthcare IT systems.

#### **B.** Advanced Transformer-Based Models

- Incorporate domain-specific transformer models such as ClinicalBERT, MedGPT, BioGPT, or other state-of-the-art LLMs.
- Fine-tuning these models on large-scale multi-institutional clinical datasets can improve understanding of complex medical narratives and rare procedures.
- Developing ensemble models that combine multiple transformers with classification algorithms could further enhance prediction accuracy.

#### C. Multi-Year and Multilingual Dataset Training

- Training on multi-year datasets allows the system to handle evolving coding standards and detect trends over time.
- Expanding the system to multilingual clinical datasets (e.g., English, Hindi, Telugu) can make it globally applicable, particularly in multilingual healthcare settings.
- Cross-institutional datasets can improve the generalisation and robustness of predictions.

## D. Explainability and Interpretability

- Providing explainable AI (XAI) features can help healthcare professionals understand why specific codes are suggested, increasing trust in the system.
- Techniques such as attention visualisation, feature importance scores, and textual justification can be added.
- This is especially important in complex cases where multiple codes apply, helping coders make informed decisions.

### E. Continuous Learning and Adaptive Models

- Implement incremental learning or online learning to allow the system to adapt to new CPT and ICD codes automatically without retraining from scratch.
- Monitoring performance continuously ensures long-term reliability and reduces model drift as coding standards evolve.

### F. Real-Time Clinical Decision Support

- Expand the system to provide decision support, such as highlighting inconsistencies in clinical documentation or suggesting additional codes based on patient history.
- Linking predictions with diagnostic recommendations or billing alerts can improve documentation quality and operational efficiency.

Vol.11 Issue 4 (2025) 9 – 19. Submitted 25/10/2025. Published 20/11/2025

• Real-time feedback can help clinicians identify missing or ambiguous information in patient summaries.

### G. Research and Academic Applications

- The system can serve as a benchmark for future research in medical NLP, multi-label classification, and RAG-based knowledge retrieval.
- Academic studies can analyse billing patterns, procedure trends, and predictive coding, providing insights for policy makers and insurance providers.
- Researchers can expand the knowledge base to cover rare or specialised procedures, improving coverage and robustness.

#### H. Scalability and Cloud Deployment

- Future versions can adopt cloud-based deployment to allow multiple hospitals or clinics to access the system simultaneously.
- Incorporating secure, encrypted data transfer and HIPAA/GDPR compliance ensures patient privacy while enabling large-scale data utilisation.
- Serverless architectures or microservices can support modular updates, high availability, and easier maintenance.

#### I. Potential for Personalised Medicine

- By linking code suggestions with patient history, lab results, and treatment plans, the system could contribute to personalised healthcare analytics.
- Integration with predictive models can enable risk assessment, outcome prediction, and optimised treatment planning, evolving the engine into a comprehensive clinical intelligence tool.

### J. Expansion to Other Medical Coding Standards

- Future work can extend the engine to support other coding standards such as SNOMED CT, LOINC, or CPT Modifiers, providing a more holistic coding solution.
- Multi-standard support can make the tool internationally applicable and beneficial for multi-specialty hospitals.

#### **REFERENCES**

- [1]. A. Perotte, N. Pivovarov, B. N. Mehta, et al., "Diagnosis code assignment: models and evaluation metrics," *Journal of the American Medical Informatics Association*, vol. 21, no. 5, pp. 845–851, 2014.
- [2]. X. Zhang, Y. Zhao, M. Le, "Convolutional neural networks for multi-label ICD coding from clinical notes," *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, vol. 16, no. 2, pp. 351–362, 2019.
- [3]. J. Lee, W. Yoon, S. Kim, et al., "BioBERT: a pre-trained biomedical language representation model for biomedical text mining," *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, 2020.
- [4]. Y. Gu, Y. Tinn, H. Cheng, et al., "Domain-specific language model pretraining for biomedical NLP," *ACM Transactions on Computing for Healthcare*, vol. 3, no. 1, pp. 1–19, 2021.

Vol.11 Issue 4 (2025) 9 – 19. Submitted 25/10/2025. Published 20/11/2025

- [5]. P. Lewis, E. Perez, A. Piktus, et al., "Retrieval-augmented generation for knowledge-intensive NLP tasks," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, 2020.
- [6]. M. Alsentzer, J. Murphy, W. Boag, et al., "Publicly available clinical BERT embeddings," *arXiv preprint* arXiv:1904.03323, 2019.
- [7]. IEEE Standards Association, *IEEE Manuscript Templates for Conference Proceedings*, 2023. [Online]. Available: https://www.ieee.org/conferences/publishing/templates.html
- [8]. T. Rajkomar, E. Oren, K. Chen, et al., "Scalable and accurate deep learning with electronic health records," *npj Digital Medicine*, vol. 1, no. 18, 2018.
- [9]. S. S. Wang, S. Y. Yang, X. Zhang, et al., "Multi-label classification for ICD coding using transformer models," *Journal of Biomedical Informatics*, vol. 115, 2021.
- [10]. H. Ji, M. Grishman, "Knowledge-base completion and clinical NLP applications," *Information Processing & Management*, vol. 56, no. 6, 2019.
- [11]. J. Li, H. Fei, Q. Wang, "Automated ICD coding via hierarchical attention networks," *IEEE Access*, vol. 8, pp. 1969–1978, 2020.
- [12]. C. Shi, S. Wang, D. Yu, "Deep learning for clinical text classification: a comparative study," *IEEE Journal of Biomedical and Health Informatics*, vol. 24, no. 9, pp. 2472–2480, 2020.
- [13]. M. Johnson, Y. Wang, "Leveraging transfer learning for ICD coding with limited data," *IEEE Access*, vol. 9, pp. 117893–117902, 2021.
- [14]. H. Suresh, A. Narayanan, "The role of artificial intelligence in healthcare: challenges and opportunities," *IEEE Engineering in Medicine and Biology Magazine*, vol. 39, no. 6, pp. 33–41, 2020.
- [15]. D. Singh, S. P. Tripathi, "Enhancing medical code prediction with contextual embeddings and ensemble models," *IEEE Access*, vol. 10, pp. 55040–55052, 2022.
- [16]. J. Gao, M. Chen, "A hybrid retrieval-generation framework for medical question answering," *Proceedings of the 2021 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pp. 1221–1227, 2021.
- [17]. R. Yadav, N. Jain, "Artificial intelligence applications in clinical documentation and coding," *IEEE Access*, vol. 11, pp. 4015–4028, 2023.
- [18]. S. Banerjee, M. Dey, "Transformer-based deep learning for healthcare text analytics," *IEEE Access*, vol. 9, pp. 143911–143924, 2021.
- [19]. Y. Wang, L. Xu, "Explainable deep learning in clinical decision support systems," *IEEE Reviews in Biomedical Engineering*, vol. 15, pp. 144–159, 2022.
- [20]. A. Das, R. K. Sharma, "Generative AI in medical applications: a systematic review," *IEEE Access*, vol. 12, pp. 71245–71260, 2024.
- [21]. M. Zhang, T. Li, "Multi-modal EHR analysis using transformer-based models," *IEEE Transactions on Biomedical Engineering*, vol. 71, no. 3, pp. 812–823, 2024.
- [22]. S. Narayan, D. Patel, "A review of large language models in healthcare NLP," *IEEE Access*, vol. 12, pp. 67590–67605, 2024.
- [23]. G. Luo, W. Sun, "Clinical text classification and code assignment using BERT-based architectures," *IEEE Access*, vol. 10, pp. 12504–12517, 2022.
- [24]. N. Brown, J. K. Lee, "Automated coding and billing: AI-based solutions for clinical practice," *IEEE Access*, vol. 11, pp. 62132–62145, 2023.
- [25]. Y. Chen, "Retrieval-Augmented Generation for domain-specific tasks: a survey," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 35, no. 5, pp. 5890–5906, 2025.