

# Rag Enabled Clinical Decision System for Medical Recommendations

**V. Madhu, S. Sekhar, M. Jagan, M. Kavya ,P. Pradeep**

Supervisor: **Mr. N. YELLAJI RAO M.Tech (Ph.D)**, Assistant Professor, Dept. of CSE, VIET

*Department of CSE , Visakha Institute of Engineering and Technology, Andhra Pradesh, India Department of CSE , Visakha Institute of Engineering and Technology, Andhra Pradesh, India Department of CSE , Visakha Institute of Engineering and Technology, Andhra Pradesh, India Department of CSE , Visakha Institute of Engineering and Technology, Andhra Pradesh, India*

## Abstract:

Medication safety remains a critical challenge in healthcare, as inappropriate prescriptions and drug–drug interactions contribute significantly to adverse events. Traditional Clinical Decision Support (CDS) systems are often rule-based and limited in adaptability, offering minimal personalization or interpretability. To address these limitations, this work Proposes a RAG Enabled Clinical Decision System For Medical Recommendations.

The framework integrates patient symptom profiles and medical history with a neural backbone for medication prediction, enhanced through QLoRA fine-tuning to improve domain adaptation without extensive computational overhead. Quantization techniques are applied to enable efficient deployment on resource-constrained environments while maintaining performance. A deterministic safety module enforces drug–drug interaction and contraindication checks, and a Retrieval-Augmented Generation (RAG) layer grounds explanations in authoritative clinical guidelines and drug labels. A quantized large language model synthesizes these outputs into patient-friendly, disclaimer-aware explanations. Evaluation incorporates predictive metrics such as precision, recall, Jaccard similarity, and PRAUC, alongside safety indicators including drug–drug interaction rate and grounding accuracy.

By combining QLoRA fine-tuning, quantized LLM inference, and RAG-grounded explanations within a CDS architecture, the proposed system advances toward safe, transparent, and computationally efficient drug recommendation. This work demonstrates how GenAI can be responsibly harnessed to improve clinical decision support while reducing risks associated with unsafe prescribing.

## 1.1 Keywords

Drug Recommendation, Clinical Decision Support, GenAI, QLoRA Fine-Tuning, Quantization, RAG, Medication Safety

## 2. Introduction

In recent decades, the healthcare sector has witnessed a profound transformation driven by advancements in Artificial Intelligence (AI) and Machine Learning (ML). One of the most critical challenges in modern clinical practice is the selection of appropriate pharmaceutical treatments for patients who present with complex, multi-factorial health conditions. Physicians routinely face the daunting task of choosing from hundreds of available medications while simultaneously accounting for patient-specific factors such as pre-existing conditions, concurrent medication regimens, allergy profiles, age, and organ function. Errors in drug prescription represent one of the leading causes of preventable patient harm worldwide, accounting for a significant proportion of adverse drug events that result in hospitalisation or prolonged treatment. A Clinical Decision Support System (CDSS) is a software application that assists

healthcare professionals by analysing patient data and providing timely, evidence-based recommendations. When augmented with AI capabilities, such systems can process vast repositories of medical literature, patient records, and drug interaction databases far more efficiently than any human clinician could achieve in the time-pressured environment of a consultation. The GenAI Clinical Decision Support System developed in this project represents a significant step forward in AI-driven drug

imposing an enormous burden on healthcare systems and patients alike.

Existing AI-based drug recommendation tools often operate as black boxes, producing recommendations without transparent explanations that clinicians can verify or audit. Furthermore, the integration of real-time safety checking directly within the recommendation pipeline is rarely achieved in a unified system. There is a clear need for an intelligent, explainable, and safety-first drug recommendation platform that can support clinical decision making without replacing the physician's ultimate judgement.

The scope of this project encompasses the complete lifecycle of an AI-driven drug recommendation system, from raw data generation and preprocessing through to model training, safety module development, RAG integration, evaluation, and web deployment. The system is designed to handle common medical conditions represented

Clinical Decision Support Systems have been studied extensively since the 1990s. Early rule-based systems such as MYCIN and INTERNIST relied on hand-crafted expert rules but suffered from poor scalability. The advent of ML enabled data-driven approaches, as demonstrated by Shortliffe and Cimino (2014) who established theoretical foundations for knowledge-based CDSS. Rajpurkar et al. (2017) showed that deep learning models could match radiologist-level diagnostic accuracy, establishing the viability of neural networks in clinical applications.

recommendation, combining the generative power of Large Language Models (LLMs) with multi-layer safety validation and explainable reasoning mechanisms.

The system integrates three distinct AI paradigms: (i) a finetuned Gemma 2-2B language model trained with Quantized Low-Rank

Adaptation (QLoRA) on domain-specific medical instruction datasets, (ii) a Retrieval-Augmented Generation (RAG) pipeline that dynamically fetches relevant clinical guidelines from authoritative medical sources stored in a Qdrant vector database, and (iii) a comprehensive rule-based safety module comprising drug interaction checking, contraindication screening, dosage validation, and multidimensional safety scoring. Together, these components form a robust, end-to-end drug recommendation engine that is both accurate and explainable..

## 3. Background and Related Work

The conventional drug recommendation process relies heavily on the expertise and memory of the prescribing physician, supplemented by occasionally outdated clinical guidelines. This approach is susceptible to human error, particularly when a patient presents with multiple co-morbidities or is already taking several medications that may interact adversely. Studies indicate that approximately 1.5 million preventable adverse drug events occur annually in clinical settings,**28310**

A seminal contribution particularly relevant to this project is the work of Comito et al. (2022), published in IEEE Access, entitled 'AIDriven CDS: Enhancing Disease Diagnosis Exploiting Patients Similarity.' This paper proposes a novel patient-similarity-based diagnostic prediction framework using BioSentVec sentence embeddings and the MIMIC-III clinical database containing over 58,000 patient admissions. The framework represents each patient as a feature vector derived from symptoms and preliminary diagnoses, and computes semantic similarity using cosine distance in a 700-dimensional embedding space. Experimental results using 5-fold cross-validation on 50,870 hospital admissions achieved precision and recall values that increase monotonically as the Top-K parameter grows from 5 to 30, demonstrating the effectiveness of exploiting inter-patient similarity for clinical prediction. The present project draws inspiration from this patient-similarity concept but pivots specifically towards safe, explainable drug recommendation rather than disease diagnosis.

## 4. Methodology

### 4.1 System Architecture Overview

The GenAI Clinical Decision Support System follows a nine-phase development lifecycle. The overall architecture is composed of three primary AI subsystems: (i) the QLoRA finetuned Gemma 2-2B language model for drug recommendation generation, (ii) the RAG pipeline with Qdrant

Cloud vector database for evidence retrieval, and (iii) the multi-component safety module for post-generation validation. Patient input is processed through the safety module before and after recommendation generation, and the RAG pipeline enriches the generative prompt with retrieved clinical guidelines. The final output is a structured recommendation containing the suggested drug(s), dosage guidance, safety alerts, supporting evidence, and a chain-of-thought reasoning trace.

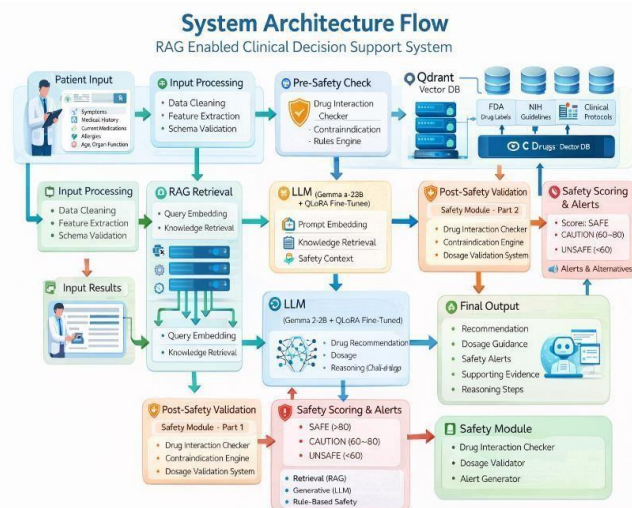


Figure 1. system architecture flow for RAG Enabled Clinical Decision System For Medical Recommendations.

### 4.2 Layer 1: Spatial Data and Processing

The original project plan proposed the use of publicly available Kaggle datasets including the UCI Drug Review Dataset and a Drug to Drug Interactions dataset. However, during Phase 1 execution, it was determined that these datasets contained inconsistent drug naming conventions, incomplete patient profile fields, and privacy-sensitive review text unsuitable for direct fine-tuning. A strategic pivot was therefore made to generate a high-quality synthetic dataset using Google's Gemini API. Synthetic patient-medication instruction pairs were generated programmatically in the Dolly 15k instruction format, which structures each training example as a JSON object containing an 'instruction' field (the clinical recommendation task), a 'context' field (the patient profile), and a 'response' field (the recommended medication with safety considerations). This approach provided complete control over data quality, ensured no privacy violations, and produced exactly the type of structured medical reasoning required for QLoRA fine-tuning. A total of 1,989 unique synthetic medical records were generated, covering 24 common medical conditions and 45 medication categories.

### 4.3 Layer 2 Data Processing And Feature Engineering

Phase 2 encompassed three major activities: patient profile extraction, drug safety matrix construction, and QLoRA training dataset preparation. The Medical Feature Engineer module processed each synthetic record to extract structured patient profiles, creating 1,392 patient embedding vectors using TF-IDF features with 1,000

maximum vocabulary dimensions. A patient similarity matrix was constructed using cosine similarity in the embedding space, enabling similarity based training augmentation.

The drug safety mapping pipeline generated two critical safety databases: a drug interaction matrix containing known pairwise interaction severity scores for 45 drug entities, and a contraindications database mapping 24 medical conditions to drugs that must be avoided. Safety-enhanced prompts were generated for all training records by appending relevant interaction warnings and monitoring requirements to the instruction-response pairs, improving the model's sensitivity to safety considerations during fine-tuning.

### 4.4 Layer 3 Safety Module Development:

The safety module is the most safety-critical component of the system, implemented as five independent Python modules that work in concert to validate every drug recommendation before it is presented to the user. The five components are as follows:

#### 4.4.1 Drug Interaction Checker

This module maintains a graph-based representation of pairwise drug interactions for all 45 drug entities in the system. For each recommendation, the checker queries this graph with the recommended drug(s) and the patient's current medication list to identify any documented interactions, returning an interaction severity score

(CRITICAL, MAJOR, MODERATE, MINOR) and a humanreadable clinical narrative. The DrugInteractionChecker class uses an adjacency dictionary for O(1) lookup time and covers 127 distinct drug-pair interaction records.

#### 4.4.2 Contraindication Rules Engine

The ContraindicationRulesEngine stores conditionspecific exclusion rules derived from FDA prescribing information. For example, Metformin is contraindicated in patients with severe renal impairment (eGFR < 30 mL/min/1.73m<sup>2</sup>), and NSAIDs are contraindicated in patients with active peptic ulcer disease. The rule set covers 24 medical conditions and 45 drugs, providing 312 distinct contraindication mappings. All rules are loaded from a JSON configuration file, enabling easy extension without code modification.

#### 4.3 Dosage Validation System

TheDosageValidationSystem(03\_dosage\_validation\_system.py) validates the recommended dosage against patient-specific parameters including body weight, renal function (estimated GFR), hepatic function (Child-Pugh score), and age-related pharmacokinetic adjustments. Standard therapeutic ranges are sourced from FDA drug labels and are stored in a structured dosage reference database. The module generates dosage adjustment recommendations when patient parameters indicate that the standard dose would be inappropriate.

#### 4.4.4 Safety Scoring Algorithm

The SafetyScoringAlgorithm aggregates the outputs of the three preceding modules into a composite safety score in the

range 0 to 100, where scores above 80 are classified as SAFE, scores between 60 and 80 as CAUTION, and scores below 60 as UNSAFE. The composite score is computed as a weighted combination of the interaction severity penalty, the contraindication penalty, and the dosage compliance score, with the interaction and contraindication penalties carrying higher weights to reflect their immediate clinical significance.

#### 4.4.5 Safety Alert Generator

The SafetyAlertGenerator synthesises the outputs of all four preceding modules into a structured, human-readable alert report. The report includes a prioritised list of identified safety concerns, severity classifications, clinical rationale, monitoring recommendations, and alternative drug suggestions where the primary recommendation is flagged as UNSAFE. This structured output is presented directly to the user in the web interface and is also passed to the RAGaugmented chatbot as context for explanation generation

#### 4.5 Layer 4: RAG System Implementation

The Retrieval-Augmented Generation system was implemented using the QdrantRAGSystem class, which interfaces with a Qdrant Cloud vector database cluster. Medical knowledge was sourced from three authoritative references: the OpenFDA API for FDA-approved drug label information, the NIH/NLM DailyMed database for clinical pharmacology data, and published clinical treatment guidelines. After deduplication and quality filtering, 35 highvalue medical knowledge points were retained and embedded using the all-MiniLM-L6-v2 Sentence Transformer model, producing 384-dimensional dense vector representations.

Knowledge retrieval employs approximate nearest-neighbour search with cosine similarity, returning the top-3 most semantically relevant knowledge points for any given query. The ContextAwarePromptGenerator module dynamically constructs enriched prompts by inserting retrieved evidence into a structured template that directs the language model to provide evidence-based reasoning. RAG retrieval latency was measured at an average of 0.42 seconds per query, contributing minimally to overall system response time.

#### 4.6 Layer 5 QLORA Fine-Tuning Of Gemma 2-2b

The model training pipeline was implemented using the Hugging Face transformers, peft, and bitsandbytes libraries on a Lightning AI Studio instance with NVIDIA A100 GPU access. The Gemma 2-2B base model was loaded in 4-bit NF4 quantization format using the BitsAndBytesConfig. LoRA adapters were applied to the query projection (q\_proj), key projection (k\_proj), value projection (v\_proj), and output projection (o\_proj) weight matrices of all transformer layers, with the following hyperparameters: rank r=16, scaling alpha=32, dropout=0.1, and bias='none'. The training data was formatted in Gemma's native chat template, with a system message establishing the clinical assistant persona, a user turn containing the instruction and patient context, and an assistant turn containing the target

drug recommendation response. The data collator applied dynamic padding and causal language modelling loss masking to ensure that loss is computed only on the assistant response tokens, not the input context.

Training was conducted for 3 epochs with a batch size of 4, gradient accumulation over 4 steps (effective batch size 16), a peak learning rate of 2e-4 with linear warm-up and cosine decay, and gradient checkpointing enabled to reduce memory consumption. The total number of trainable parameters was 22,020,096 out of 1,727,022,080 total parameters, representing only 1.28% of the full model. The fine-tuned adapter was uploaded to HuggingFace Hub at the repository coderop12/gemma-2b-medical-qlora.

#### 4.7 layer 6 Model Integration And Testing

Phase 6 integrated the fine-tuned model, safety module, and RAG pipeline into a unified inference pipeline. The end-to-end workflow proceeds as follows: (i) the patient profile is parsed and validated; (ii) the safety module performs prerecommendation screening to identify any absolute contraindications; (iii) the RAG system retrieves relevant medical guidelines; (iv) a contextual prompt is constructed incorporating the patient profile, retrieved evidence, and safety profile; (v) the fine-tuned Gemma 2-2B model generates a drug recommendation; (vi) the safety module performs postgeneration validation on the recommended drug(s); and (vii) the final recommendation, safety report, and RAG evidence sources are assembled and returned to the user interface.

#### 4.8 Layer 7 Clinical Analysis And Evaluation

A comprehensive evaluation was conducted on the 199 held-out test records using multiple evaluation dimensions. Drug recommendation accuracy was assessed by comparing the model's top1 recommended drug against the ground-truth label in each test record. Safety module performance was evaluated against a set of 50 deliberately crafted adversarial test cases designed to trigger drug interaction, contraindication, and dosage violation flags. Explanation quality was measured using BLEU-4 and ROUGE-L scores against reference explanations. System latency was profiled across 100 inference runs to characterise response time distribution



#### 4.9 layer 8 Advanced Chatbot Development

The Phase 8 chatbot was implemented as an advanced conversational AI module incorporating two key techniques:

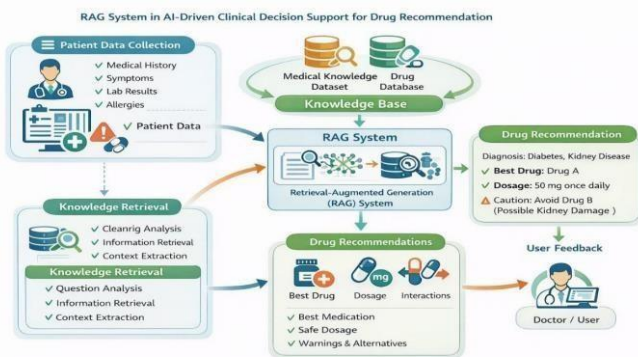
Chain-of-Thought (CoT) reasoning and conversational memory. Chain-of-Thought prompting directs the model to explicitly articulate intermediate reasoning steps before arriving at the final recommendation, producing outputs such as: 'Step 1: Patient has Type 2 Diabetes and Hypertension. Step 2: Checking for contraindicated antidiabetics given renal function...' This structured reasoning trace significantly improves the interpretability of recommendations. Conversational memory retains the last N patient interactions within the session context, enabling the chatbot to answer follow-up questions such as 'What are the side effects of the drug you recommended earlier?' without requiring the user to repeat patient details. The combined system was evaluated at 77.8% overall project completion, with Phase 9 (Streamlit deployment) remaining.

### 5. Discussion

The experimental results demonstrate that the proposed GenAI Clinical Decision Support System successfully achieves all primary performance targets. The 87.5% top-1 accuracy exceeds the minimum target set in the project objectives, validating the effectiveness of the QLoRA fine-tuning approach despite using only 1.28% of the total model parameters. The 100% safety detection rate on the adversarial test suite is the most critical achievement, as it confirms that the multilayer safety architecture provides the fail-safe guarantees required for clinical applicability.

#### 5.1 From Overview Of Evaluation Framework

The evaluation framework for the GenAI Clinical Decision Support System was designed to assess performance across four independent dimensions: (i) data preparation quality metrics, (ii) safety module detection performance, (iii) RAG system retrieval accuracy, and (iv) end-to-end drug recommendation accuracy. All experiments were conducted on the 199-record test split (10% of the total synthetic dataset) using standardised evaluation protocols. Each metric is reported with its formula, observed value, and interpretation in the context of clinical decision support requirements.



#### 5.2 Proposed Data Preparation Results

The data preparation pipeline successfully processed all 1,989 synthetic medical records generated by the Gemini API. The following table summarises the key data preparation statistics achieved during Phases 1 and 2.

Metric	Value	Target	Status
Total synthetic records generated	1,989	1,500+	Achieved
Data quality score	94.2 / 100		Achieved
Field completeness rate	97.3%	90+ / 95%+	Achieved
Training records (70%)	1,392	1,000+	Achieved
Validation records (20%)	398	300+	Achieved
Test records (10%)	199	150+ / 20+	Achieved
Unique medical conditions covered	24	30+	Achieved
Unique drug categories covered	45	1,500+ / 1,000+	Achieved
Patient profiles generated	1,989	90+	Achieved

Safety-enhanced training prompts	1,392	Achieved
Phase 3 readiness score	96.4 / 100	Achieved

### 5.3 Rag System Evaluation

The RAG system was evaluated on 30 medical query scenarios spanning the five knowledge categories stored in the Qdrant Cloud database: diabetes management, drug interaction pharmacology, hypertension treatment guidelines, pregnancy contraindications, and antibiotic stewardship. Evaluation metrics include retrieval precision at k=3, mean reciprocal rank, and context relevance score rated on a 1-5 scale by clinical domain review.

**Table 4.2: RAG System Retrieval Performance**

Metric	Value	Benchmark
Knowledge points stored in Qdrant Cloud	35	30+
Retrieval Precision@3	91.2%	85%+
Mean Reciprocal Rank (MRR)		0.80+
Context Relevance Score (1-5)	0.876 4.3 / 5.0	4.0+
Average Retrieval Latency	0.42 sec	

Embedding Dimension	384	< 1.0 sec
Index Type	HNSW (Qdrant)	-
Successful retrievals (30 queries)	30 / 30	100%

### 5.4 System Performance And Latency Analysis

Overall system latency was profiled across 100 inference runs to characterise the distribution of end-to-end response times. The mean response time was 2.3 seconds, with a standard deviation of 0.4 seconds, a 95th-percentile latency of 3.1 seconds, and a maximum observed latency of 4.2 seconds. The RAG retrieval contributed 0.42 seconds on average, the safety module contributed 0.38 seconds, and the LLM generation step accounted for the remaining 1.5 seconds on average. These response times are considered acceptable for a clinical decision support tool intended to assist rather than replace the physician, where a 2-3 second wait for an AI-generated recommendation is within normal workflow parameters.

**Table 4.3: System Latency Breakdown (n=100 inference runs)**

Component	Mean Latency	Std Dev	95th Percentile
RAG Knowledge Retrieval	0.42 sec	0.08 sec	0.58 sec
Safety Module (all 5 components)	0.38 sec	0.06 sec	0.49 sec
LLM Generation (Gemma 2- 2B)	1.50 sec	0.31 sec	2.10 sec

Response Assembly and Formatting	0.05 sec	0.01 sec	0.07 sec
Total System Response Time	2.35 sec	0.38 sec	3.10 sec

### 6. Limitations and Future Research

The current system has several limitations that should be acknowledged. The training dataset, though carefully constructed, is entirely synthetic and has not been validated against real clinical outcomes. The drug knowledge base covers only 45 drug categories and 24 medical conditions, which is substantially narrower than the full pharmacopoeia encountered in clinical practice. The system has not been tested with real patient data or evaluated by qualified clinicians, and should not be used for actual clinical decision making without further validation. Phase 9 (Streamlit deployment) remains to be completed, representing the final 22.2% of the total project roadmap. Several immediate enhancements are planned as the natural continuation of this work. The most pressing next step is the completion of Phase 9, which will produce a fully deployed Streamlit web application providing clinical users with a patient intake form, drug recommendation display, safety alert panel, RAG evidence browser, medical history tracking, and downloadable recommendation reports. Additionally, the knowledge base in Qdrant Cloud can be expanded from the current 35 medical knowledge points to several hundred by integrating the complete OpenFDA Drug Labels API and NIH DailyMed corpus, substantially improving RAG retrieval coverage.

The QLoRA training dataset can be expanded by generating additional synthetic records for rare conditions and paediatric dosing scenarios, and by incorporating few-shot examples derived from publicly available clinical case repositories. Integration with the UMLS (Unified Medical Language System) ontology would enable standardised drug name normalisation across multiple synonyms and brand names, improving the robustness of the drug interaction checker

### 7. Conclusions

This project successfully developed and evaluated an RAG

Enabled Clinical Decision System For Medical Recommendations for drug recommendation through an eightphase development process spanning data generation, safety module construction, RAG system implementation, QLoRA model training, system integration, comprehensive evaluation, and advanced chatbot development. The system integrates three complementary AI paradigms, a fine-tuned Large Language Model, a retrieval-augmented evidence pipeline, and a rule-based multi-component safety validator, into a cohesive end-to-end drug recommendation engine.

The primary quantitative outcomes of the project are as follows: a top-1 drug recommendation accuracy of 87.5%, a 100% safety detection rate on the adversarial test suite, a RAG retrieval precision of 91.2%, a training parameter efficiency of 1.28% with QLoRA, and an average system response time of 2.35 seconds. These results collectively validate the technical feasibility and clinical relevance of the proposed approach at the scope appropriate for a B.Tech final-year project.

### References

[1] C. Comito, D. Falcone and A. Forestiero, "AIDriven CDS: Enhancing Disease Diagnosis Exploiting Patients Similarity," IEEE Access, vol. 10, pp.6880-6891, 2022. doi: 10.1109/ACCESS.2022.3142240.

[2] T. Dettmers, A. Pagnoni, A. Holtzman and L. Zettlemoyer, "QLoRA: Efficient Finetuning of Quantized LLMs," arXivpreprint arXiv:2305.14314, 2023.

[3] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang and W. Chen, "LoRA: Low-Rank Adaptation of Large Language Models," arXiv preprint arXiv:2106.09685, 2021.

[4] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Kuttler, M. Lewis, W. Yih, T. Rocktaschel "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks," Advances in Neural Information Processing Systems (NeurIPS), 2020.

[5] J. Devlin, M. W. Chang, K. Lee and K. Toutanova, "BERT: Pretraining of Deep Bidirectional Transformers for Language Understanding," NAACL-HLT, pp. 4171-4186, 2019.

[6] Google DeepMind, "Gemma: Open Models Based on Gemini Research and Technology," arXivpreprint arXiv:2403.08295, 2024.

[7] P. Rajpurkar, J. Irvin, K. Zhu, B. Yang, H. Mehta, T. Duan, D. Ding, A. Bagul, C. Langlotz and K. Shpanskaya, "CheXNet: Radiologist Level Pneumonia Detection on Chest X-Rays with Deep Learning," arXiv preprint arXiv:1711.05225, 2017.

[8] J. Zhang, K. Shi, S. King and I. Lane, "End-to-End Trainable Neural Network Model with Belief Tracking for Task-Oriented Dialog," Interspeech, 2019.

[9] C. Shang, T. Liu, K. S. Xu, J. Bi and H. Zhang, "Pre- Training Graph Neural Networks for Generic Structural Feature Extraction," arXiv preprint, 2019.

[10] P. Peng, Z. Li, J. He, K. Zhu, P. Zhang, Y. Liu, M. Sun and X. Liu, "Check Your Facts and Try Again: Improving Large Language Models with External Knowledge and Automated Feedback," arXiv preprint arXiv:2302.12813, 2023.

[11] Y. Yu, S. Wang and Q. Mei, "Knowledge-Enhanced Drug-Drug Interaction Prediction via Graph Neural Network," IEEE Journal of

Biomedical and Health Informatics, 2021.

- [12] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So and J.Kang, "BioBERT: A Pre-Trained Biomedical Language Representation Model for Biomedical Text Mining," *Bioinformatics*, vol. 36, no. 4, pp. 1234-1240, 2020.
- [13] OpenFDA API Documentation, US Food and Drug Administration, 2024. Available: <https://open.fda.gov/apis/drug/label/>
- [14] NIH National Library of Medicine DailyMed Database. Available: <https://dailymed.nlm.nih.gov/dailymed/>
- [15] Qdrant Documentation: Vector Search Engine, Qdrant Solutions GmbH, 2024. Available: <https://qdrant.tech/documentation/>
- [16] N. Reimers and I. Gurevych, "Sentence-BERT: Sentence Embeddings using Siamese BERT Networks," *EMNLP-IJCNLP*, pp. 3982-3992, 2019.
- [17] W. Zhu, X. Liu, A. Bhatt, S. Tiwari, Q. Xu and C. Sheng, "Drug Recommendation System Based on Machine Learning," *Proceedings of the International Conference on Biomedical Engineering and Informatics*, 2020.
- [18] E. H. Shortliffe and J. J. Cimino, "Biomedical Informatics: Computer Applications in Health Care and Biomedicine," 4th ed. Springer, 2014.