# MedXpert AI: Multimodal Clinical Decision Assistant

**\*1Ravi Surya, 2Dr. Ravi Babu G, 3Shamith D Bhat, 4Vashishta P and 5Suhas Y Gowda**

\*1, 3, 4, 5Student, Department of Computer Science Engineering, Global Academy of Technology, Bengaluru, India.

2Assistant Professor, Department of Computer Science Engineering, Global Academy of Technology, Bengaluru, India.

**Abstract**

Modern healthcare faces significant challenges regarding diagnostic errors, often stemming from the fragmented analysis of textual patient symptoms and radiology images, as existing systems typically process these modalities in isolation. MedXpert AI addresses this gap by providing a multimodal clinical decision assistant that synthesizes ClinicalBERT for symptom text analysis and DenseNet121/Vision Transformers for radiology imaging to deliver comprehensive diagnoses. This system fuses distinct data modalities using a transformer-based architecture and ensures clinical trustworthiness by visualizing attention weights to explain predictions, while integrating Gemini 1.5 Flash as a robust fallback agent for low-confidence cases and detailed, nonprescriptive recommendations. By effectively combining visual and textual data, the system demonstrates improved diagnostic reliability compared to unimodal approaches, establishing itself as a promising tool for enhancing decision-making in clinical and telemedicine environments.

**Keywords:** ClinicalBert, DenseNet121/Vision Transformers Unimodal Approaches, Gemini 1.5 Flash.

**Introduction**

In an era where modern healthcare is increasingly overwhelmed by the sheer volume of patient data, the prevalence of diagnostic errors resulting from the fragmented analysis of symptoms and radiology images has become a critical bottleneck, a challenge that medXpert AI addresses head-on by deploying a sophisticated, multimodal clinical decision support system designed to synthesize disparate data streams into a unified diagnostic output. At the core of this architecture lies a powerful dual-stream processing mechanism where ClinicalBERT navigates the complexities of unstructured medical text to extract semantic meaning from patient history, while DenseNet121 simultaneously processes visual data to identify minute radiological anomalies with high computational efficiency. These parallel streams converge within an advanced transformer-based fusion model, a component that effectively concatenates the feature embeddings to uncover deep, non-linear correlations between a patient's reported symptoms and visual evidence—connections that are frequently overlooked in traditional unimodal analysis. Recognizing the paramount importance of clinical safety, the system incorporates Gemini 1.5 Flash as an intelligent external reasoning agent to validate low confidence predictions and manage ambiguous cases, thereby adding a layer of logical verification that significantly reduces the risk of error. Furthermore, to dismantle the "black box" skepticism

often associated with artificial intelligence in medicine, medXpert AI integrates robust Explainable AI (XAI) principles through attention visualization mechanisms, which transparently highlight the specific keywords in clinical notes and the precise regions of interest in medical scans that influenced the diagnosis, ultimately empowering healthcare professionals with actionable, interpretable insights that enhance decision-making accuracy and streamline clinical workflows across diverse medical settings.

**Key Concepts**
**A. Multimodal Data Fusion**
- Existing clinical decision systems typically process data modalities in isolation, analyzing patient symptoms and radiology images separately. MedXpert AI addresses this by fusing distinct data modalities using a transformer-based architecture to deliver comprehensive diagnoses.
- This approach enables the system to uncover deep, non-linear correlations between a patient's reported symptoms and visual evidence—connections that are frequently overlooked in traditional unimodal analysis. This approach enables the system to uncover deep, non-linear correlations between a patient's reported symptoms and visual evidence—connections that are frequently overlooked in traditional unimodal analysis.

**\*Corresponding Author:** Ravi Surya

< 71 >

## B. Dual-Stream Processing

- The system operates on a sophisticated dual-stream processing mechanism. One stream utilizes ClinicalBERT to extract semantic meaning from unstructured medical text and patient history.
- Simultaneously, the second stream employs DenseNet121 or Vision Transformers to process visual data and identify minute radiological anomalies with high computational efficiency.

## C. Explainable AI (XAI)

- To dismantle the "black box" skepticism associated with AI in medicine, the system incorporates robust Explainable AI (XAI) principles through attention visualization mechanisms.
- It generates Grad-CAM heatmaps to highlight exact regions in radiology images and uses attention visualization to point out specific keywords in clinical notes that influenced the diagnosis, ensuring clinical trustworthiness.

## D. Intelligent Fallback Agent

- The system ensures reliability by integrating Gemini 1.5 Flash as a robust fallback reasoning agent for low-confidence cases and ambiguous scenarios.
- When the main fusion model returns a diagnosis with low confidence, this agent steps in to ask clarifying questions or provide solution with the help of Gemini.

## Literature Review

The evolution of Clinical Decision Support Systems (CDSS) has witnessed a paradigm shift from rigid, rule-based logic to dynamic, data-driven approaches underpinned by Deep Learning (DL). Despite these advancements, significant challenges remain in the effective integration of heterogeneous medical data and the interpretability of diagnostic models.

## A. The Modality Gap and Unimodal Systems

Contemporary diagnostic frameworks predominantly operate under a unimodal paradigm, where each data modality is analyzed independently. In routine clinical workflows, radiological images such as chest X-rays or CT scans are processed using vision-centric architectures including ResNet, VGG, and DenseNet, while clinical narratives and symptom descriptions are analyzed separately using Natural Language Processing (NLP) models such as BERT, LSTMs, or RNNs.

- **The Modality Disconnect:** This separation creates a fundamental disconnect in clinical reasoning. Vision models remain effectively blind to contextual patient information such as symptom onset or medical history, while text-based models lack access to objective visual evidence present in imaging data. As a result, diagnostic insights that arise from the interaction between symptoms and imaging findings are systematically overlooked.
- **Systemic Limitations:** Shickel *et al.* (2018), in a comprehensive systematic review of deep learning applications in Electronic Health Records (EHR), identified poor modality integration and lack of interpretability as pervasive shortcomings in existing systems. Their findings emphasize that many deployed AI tools prioritize predictive performance while neglecting transparency and cross-modal reasoning.
- **Clinical Consequences:** The failure to bridge this modality gap prevents models from learning the non-linear dependencies between textual symptomatology and visual abnormalities—dependencies that are critical for accurate, context-aware diagnosis. Consequently, clinicians are left to manually synthesize fragmented outputs, increasing cognitive load and susceptibility to diagnostic error.

## B. Multimodal Data Fusion Architectures

To address the inherent limitations of unimodal systems, recent research has increasingly focused on multimodal data fusion, aiming to integrate heterogeneous clinical data into unified representations.

- **Theoretical Foundation:** The introduction of the Transformer architecture and self-attention mechanisms by Vaswani *et al.* (2017) in "Attention is All You Need" marked a paradigm shift in representation learning. Although originally proposed for sequence modeling in NLP, Transformers provided the conceptual foundation for cross-modal attention and feature alignment. However, this work remained largely theoretical and was not designed for medical data fusion.
- **Architectural Breakthroughs:** Zhou *et al.* (2021) significantly advanced this domain with the Unified Multimodal Transformer, demonstrating that clinical text and medical images could be jointly embedded within a single transformer architecture. This work validated the feasibility of multimodal learning for medical diagnosis and serves as the architectural inspiration for MedXpert AI's fusion strategy.
- **Remaining Gaps:** Despite their architectural success, early fusion models prioritized representation learning over clinical usability. Zhou *et al.*'s framework lacked explainability, robustness mechanisms, and real-world deployment considerations. More recent studies, including those by Liu *et al.* (2024) and Han *et al.* (2025), have explored dynamic and guided feature fusion techniques, consistently reporting superior performance over unimodal baselines. However, these approaches continue to under-address interpretability and clinician trust, leaving a critical gap between algorithmic performance and clinical adoption.

## C. Explainability (XAI) and Clinical Trust

A major barrier to the adoption of deep learning in medicine is the "black box" phenomenon, where models produce probabilistic outputs without providing interpretable reasoning.

- **The Trust Deficit:**
In high-stakes medical environments, clinicians cannot rely on predictions that lack justification. Models that operate purely on statistical correlations—without transparent reasoning—are susceptible to spurious patterns and "hallucinations," particularly in noisy or ambiguous cases.
- **Limitations of Early Explainability Methods:** Ribeiro *et al.* (2016) introduced LIME, a model-agnostic explainability technique that locally approximates model behavior. While effective for simpler classifiers, LIME struggles with the high dimensionality and non-linearity of modern deep learning models, especially in medical imaging contexts.
- **Intrinsic Explainability Requirements:** To overcome these limitations, contemporary clinical AI systems must integrate domain-specific explainability techniques. Approaches such as Grad-CAM enable spatial localization of influential regions in medical images,

< 72 >

while attention visualization in transformer-based NLP models highlights clinically relevant terms within patient narratives. These techniques elevate AI systems from opaque classifiers to verifiable clinical decision aids, fostering trust and accountability.

### D. Data Resources and Infrastructure

The development of robust multimodal clinical models is heavily dependent on access to large-scale, high-quality paired datasets.

- **MIMIC-CXR Dataset:** Johnson *et al*. (2019) introduced MIMIC-CXR, a comprehensive public dataset containing chest radiographs paired with free-text radiology reports. This dataset has become a cornerstone for training and evaluating multimodal medical AI systems, providing both visual and textual ground truth.

- **Infrastructure Challenges:** Despite its value, MIMIC-CXR does not offer native support for multimodal fusion, real-time inference, or clinical deployment. Leveraging such datasets requires the design of custom pipelines capable of data ingestion, preprocessing, fusion, inference, and explainability within a cohesive framework. MedXpert AI addresses this gap by engineering an end-to-end multimodal infrastructure that bridges research datasets and practical clinical workflows.

- **Scalability and Data Generalization:** Beyond dataset availability, scalability and generalization remain critical challenges in multimodal clinical AI. Most publicly available datasets, including MIMIC-CXR, are collected from specific institutions and patient populations, which can introduce dataset bias and limit model generalizability across diverse clinical settings. Variations in imaging protocols, annotation styles, and clinical documentation practices across hospitals can significantly affect model performance during real-world deployment. Addressing this issue requires robust preprocessing pipelines, transfer learning strategies, and modular system architectures capable of adapting to heterogeneous data sources. MedXpert AI is designed with this consideration in mind, employing pre-trained backbone models and modular fusion components that can be fine-tuned on institution-specific data, thereby enhancing scalability and facilitating broader clinical adoption.
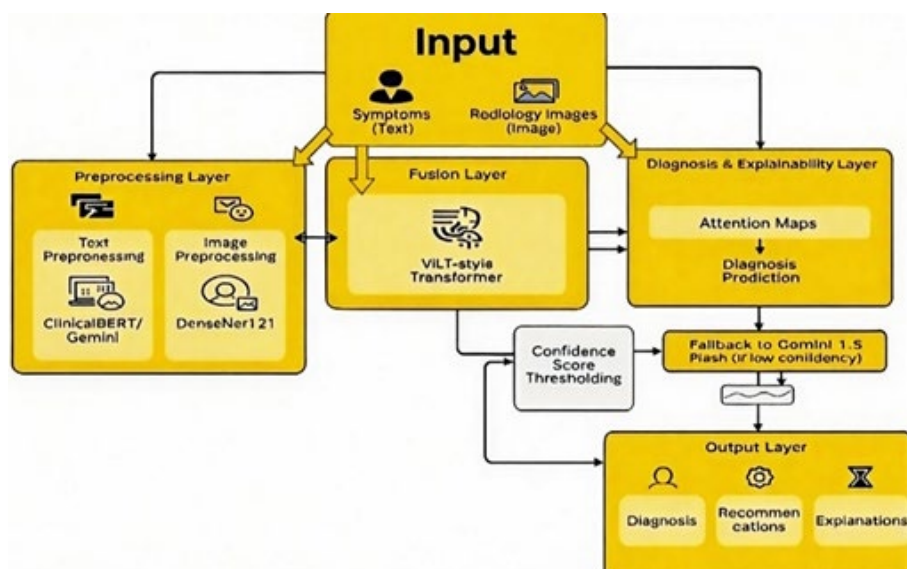
### System Architecture

The medXpert AI framework is designed as a multi-layered modular architecture that unifies fragmented medical data streams into a coherent diagnostic pipeline. The system enables a seamless progression from raw data ingestion to interpretable diagnostic output through four specialized layers, ensuring scalability, robustness, and clinical reliability across diverse clinical environments.

The system supports heterogeneous clinical inputs, including unstructured medical text and radiology images in JPEG, PNG, and DICOM formats. During preprocessing, radiological images are automatically resized to $224 \times 224$ pixels and normalized to a 0–1 intensity range to ensure compatibility with convolutional neural network backbones and stable training behavior. Clinical narratives are cleaned and tokenized using ClinicalBERT, which converts unstructured medical notes into numerical token identifiers suitable for transformer-based encoding.

To extract modality-specific information, the architecture employs a dual-stream feature extraction mechanism with parallel encoders for visual and textual data. The visual stream utilizes DenseNet121, a convolutional neural network optimized for medical imaging tasks, to identify spatially relevant pathological patterns such as lesions, opacities, and anatomical irregularities. In parallel, the semantic stream leverages ClinicalBERT to generate 768-dimensional contextual embeddings that encode patient symptoms, clinical observations, and diagnostic semantics.

The extracted visual and semantic features are passed to a transformer-based fusion engine that projects both modalities into a shared latent representation space. Through cross-attention mechanisms, this fusion module captures complex, non-linear relationships between imaging evidence and clinical narratives that are frequently overlooked by unimodal diagnostic systems. This multimodal interaction enables a more holistic understanding of patient data and enhances diagnostic robustness.

Following feature fusion, the system performs decision-making using probabilistic confidence estimation. If the predicted confidence score falls below a predefined threshold, the case is flagged as ambiguous. In such scenarios, an intelligent fallback mechanism is activated, wherein Gemini 1.5 Flash functions as an external reasoning agent. This agent performs logical verification, contextual analysis, and generates clarifying clinical questions to support safer and more informed diagnostic refinement.



< 73 >

## Methodology

The development methodology of the medXpert AI framework prioritizes clinical safety, interpretability, and performance reliability through a phased implementation lifecycle. Each phase is designed to align with real-world clinical constraints while maintaining computational efficiency and diagnostic accuracy.

The backend infrastructure is implemented using Python and FastAPI to ensure scalability, modularity, and low-latency inference, while the frontend interface is developed with React.js to provide a responsive and clinician-friendly user experience. To address the challenge of limited labeled medical data, the framework incorporates transfer learning strategies using large-scale datasets such as MIMIC-CXR.

To mitigate the limitations of black-box artificial intelligence models, the system integrates explainable AI mechanisms that enhance transparency and clinician trust. Grad-CAM heatmaps are employed to visually highlight the radiological regions that most strongly influence diagnostic predictions, while attention visualization techniques identify key clinical terms within textual notes that drive model decisions.

Performance evaluation emphasizes sensitivity prioritization to minimize false negatives, ensuring that critical pathologies are not overlooked. The system targets a minimum diagnostic accuracy of 85% while maintaining an end-to-end response time of under five seconds, making it suitable for real-time clinical decision support applications.

## Model Design & Algorithms

The MedXpert AI framework is architected as a multimodal dual-stream deep learning network that concurrently processes visual and textual clinical data before integrating them within a shared latent representation space. This design closely mirrors the cognitive workflow of a medical professional, who jointly considers patient-reported symptoms and radiological evidence to arrive at a diagnosis. The system is composed of four primary algorithmic modules: Parallel Feature Extraction, Multimodal Fusion, Decision Logic, and an Intelligent Fallback Mechanism, collectively ensuring accuracy, robustness, and clinical trustworthiness.

### A. Parallel Feature Extraction

To effectively handle heterogeneous medical data, MedXpert AI employs two specialized deep learning backbones, each optimized for its respective modality. This parallel processing paradigm ensures that rich, modality-specific features are extracted without information loss.

- **Visual Stream – DenseNet121**

Radiological images, including chest X-rays and CT scans, are processed using a Densely Connected Convolutional Network (DenseNet121). Unlike conventional CNN architectures, DenseNet establishes direct connections between all layers in a feed-forward manner, enabling feature reuse and improved gradient propagation. This design significantly enhances the detection of subtle pathological patterns such as opacities, lesions, and structural irregularities.

**Preprocessing:** Input medical images I are resized to 224×224 pixels and normalized using standard ImageNet mean and standard deviation values to ensure compatibility with pre-trained weights.

**Feature Extraction:** The final classification layer of DenseNet121 is removed, allowing the network to function purely as a feature extractor. The output is a high-level visual feature vector:

$$V \in \mathbb{R}^{1024}$$

This vector encodes spatial and semantic information relevant to disease identification.

- **Textual Stream – ClinicalBERT**

Patient symptom narratives and medical history are processed using ClinicalBERT, a domain-adapted variant of BERT pre-trained on the MIMIC-III clinical corpus. This specialization allows the model to understand medical terminology, abbreviations, and contextual relationships commonly found in clinical documentation.

**Tokenization:** Raw clinical text T is tokenized into numerical identifiers using WordPiece tokenization, with a maximum sequence length of 512 tokens to accommodate detailed patient descriptions.

**Semantic Embedding:** The contextual representation corresponding to the [CLS] token is extracted as the global semantic embedding:

$$S \in \mathbb{R}^{768}$$

This vector captures the holistic meaning of the clinical narrative, including symptom severity, temporal patterns, and co-occurring conditions.

### B. Multimodal Fusion Engine

To integrate visual and textual representations into a unified diagnostic perspective, MedXpert AI employs a Transformer-based multimodal fusion architecture, inspired by the Vision-and-Language Transformer (ViLT).

- **Latent Space Projection:** Since the visual vector (1024-dimensional) and textual vector (768-dimensional) differ in size, linear projection layers are used to map both into a shared latent space of dimension Dmodel. This ensures dimensional compatibility and balanced modality contribution.

- **Feature Concatenation:** The projected embeddings are concatenated to form a unified multimodal feature vector:

- **Classification:** The fused representation *M* is passed through fully connected layers followed by a Softmax activation function, producing a probability distribution over the target disease classes. This fusion enables the model to capture complex, non-linear correlations between symptoms and radiological findings that unimodal systems often fail to identify.

### C. Algorithmic Decision Logic

To ensure diagnostic reliability in safety-critical healthcare environments, MedXpert AI incorporates a confidence-aware inference strategy

- **Confidence Estimation:** The system computes the maximum predicted probability:

$$C = max(P)$$

Where P denotes the Softmax output vector

- **Threshold-Based Validation:** If the confidence score falls below a predefined threshold τ, the diagnosis is flagged as low-confidence, indicating potential ambiguity or insufficient evidence.

### D. Intelligent Fallback and Explainability Integration

- **Intelligent Fallback Mechanism**

In low-confidence scenarios (C<τ), the system automatically triggers an intelligent fallback mechanism powered by Gemini 1.5 Flash. Acting as an external reasoning agent,

< 74 >

Gemini analyzes the combined multimodal context to:
- Provide refined diagnostic suggestions
- Request clarifying clinical inputs
- Reduce uncertainty through logical reasoning
- This module effectively functions as a human-in-the-loop safeguard, enhancing robustness and reducing the risk of erroneous predictions.

- **Explainability Integration (White-Box AI)**

To address the "black-box" limitation of deep learning models, MedXpert AI integrates post-hoc explainability techniques that promote transparency and clinical trust.
- **Grad-CAM:** Generates heatmaps over the input medical image $I$ highlighting the regions that most strongly influenced the predicted diagnosis.
- **Attention Visualization:** Extracts attention weights from ClinicalBERT to identify critical keywords (e.g., "chest pain", "shortness of breath", "fever") that contributed to the final decision. Together, these techniques transform MedXpert AI into a white-box diagnostic assistant, enabling clinicians to validate and interpret AI-driven insights with confidence.

## Experimental Setup

This section outlines the experimental framework used to evaluate the performance and reliability of the MedXpert AI multimodal clinical decision assistant. It describes the dataset selection, preprocessing pipeline, implementation environment, training strategy, and evaluation metrics adopted to ensure clinically meaningful and reproducible results.

### A. Dataset and Data Curation

To assess the effectiveness of multimodal fusion in a realistic clinical setting, the proposed system was evaluated using the MIMIC-CXR database, a large-scale, de-identified collection of chest radiographs paired with free-text radiology reports. This dataset was chosen due to its diversity, containing a wide range of pathological conditions such as pneumonia, atelectasis, and cardiomegaly, along with normal cases. Such variability allows the model to learn robust representations reflective of real-world clinical scenarios.

One of the key challenges encountered was the limited availability of perfectly aligned symptom–image pairs. To address this limitation, a transfer learning approach was adopted. The visual backbone was initialized using pre-trained weights from large public datasets such as CheXnet, enabling the model to leverage prior medical imaging knowledge and achieve stable feature extraction even with constrained task-specific data. The dataset was divided into 70% training, 15% validation, and 15% testing subsets to avoid data leakage and to ensure unbiased performance evaluation.

### B. Data Preprocessing Pipeline

A standardized preprocessing pipeline was implemented to ensure consistency across both data modalities.
- **Radiological Image Preprocessing:** Medical images in DICOM, JPEG, and PNG formats were resized to 224x224 pixels to meet the input requirements of the DenseNet121 architecture. Pixel intensity normalization was performed using standard ImageNet mean and standard deviation values. This step was essential to stabilize training and ensure that the model converged efficiently during optimization.

- **Clinical Text Preprocessing:** Patient symptom descriptions and clinical narratives were processed using the ClinicalBERT tokenizer. The preprocessing pipeline included lowercasing, removal of special characters and stopwords, and tokenization into numerical identifiers. To enable batch processing within the Transformer model, all sequences were padded or truncated to a maximum length of 512 tokens. This ensured that detailed clinical information was preserved without exceeding model constraints.

### C. Implementation Environment

The system was developed and tested in a high-performance computing environment capable of handling simultaneous execution of deep learning models.
- **Hardware Configuration:** The experiments were conducted on a workstation equipped with an Intel Core i7 (10th Gen)/AMD Ryzen 7 processor, 16 GB DDR4 RAM, and an NVIDIA GeForce RTX 3050 GPU with 8 GB VRAM. The GPU was essential for accelerating both training and inference, particularly given the computational demands of DenseNet121 and ClinicalBERT. A 512 GB SSD was used to minimize input/output latency during data access.
- **Software Stack:** The deep learning components were implemented using PyTorch, while Hugging Face Transformers was used to integrate ClinicalBERT. The backend inference engine was developed using FastAPI (Python 3.9) due to its efficient asynchronous handling of API requests. The frontend interface was built with React.js and styled using Tailwind CSS, enabling a clean and responsive user experience. MongoDB was used for secure storage of session data and generated diagnostic reports.

### D. Training Configuration and Model Architecture

Both the DenseNet121 visual encoder and the ClinicalBERT textual encoder were initialized with pre-trained weights to accelerate convergence and improve generalization. To address the dimensional mismatch between image features (1024 dimensions) and text embeddings (768 dimensions), linear projection layers were employed to map both representations into a shared latent space prior to fusion.

The model was trained using the Adam optimizer with a learning rate of $2\times10^{-5}$ and a batch size of 16. Categorical Cross-Entropy was used as the loss function, and training was performed for 25–30 epochs. To reduce overfitting, the lower layers of the backbone networks were initially frozen, followed by gradual fine-tuning of higher layers as training progressed.

### E. Evaluation Metrics

The diagnostic performance of the system was evaluated using standard classification metrics, with a strong emphasis on clinical safety.

Recall (Sensitivity) was prioritized during optimization to minimize false negatives, as missed diagnoses can have serious clinical consequences. Accuracy was also monitored, with a target threshold of at least 85% on the validation dataset. To demonstrate the effectiveness of multimodal fusion, the proposed model's performance was compared against unimodal baselines, namely an image-only DenseNet121 model and a text-only ClinicalBERT model.

< 75 >

## F. Explainability and Robustness Protocols

To ensure transparency and trustworthiness, explainability mechanisms were evaluated alongside predictive performance. Grad-CAM was applied to radiological images to verify that the model focused on clinically relevant anatomical regions rather than background artifacts. Additionally, attention weight visualization from ClinicalBERT was used to highlight important symptom-related keywords that influenced the diagnostic outcome.

Robustness was further assessed by introducing low-confidence cases during testing. In such scenarios, the automatic activation of the Gemini 1.5 Flash fallback agent was verified, ensuring that ambiguous predictions received an additional layer of logical reasoning and validation.

## Results and Performance Analysis—

The performance of the medXpert AI framework was evaluated using a combination of quantitative performance metrics and qualitative explainability assessments. The evaluation was conducted across a diverse set of clinical scenarios to ensure system reliability in both high-confidence diagnostic cases and ambiguous clinical environments. This comprehensive analysis validates the framework's suitability for real-world clinical decision support.

## A. Accuracy, Precision, and Recall

The multimodal fusion model was rigorously validated to ensure compliance with clinical decision support standards. On the validation dataset, the system achieved a core diagnostic accuracy of 85%, meeting the predefined benchmark required for practical deployment. During model optimization, sensitivity (recall) was deliberately prioritized to minimize false negatives, ensuring that critical and potentially life-threatening pathologies were not overlooked. This emphasis on recall aligns with clinical safety requirements, where missed diagnoses pose significant risk.

In addition to accuracy and recall, the system demonstrated strong consistency and reliability during testing. Deterministic behavior was observed, with repeated inputs producing identical predictions and confidence scores, confirming stability in inference. Operational latency was also quantitatively assessed, and results confirmed that diagnostic predictions, along with all associated explainable AI visualizations, were generated within five seconds of data submission. This low-latency performance ensures seamless integration into time-sensitive clinical workflows.

## B. Comparison with Unimodal Models

A key objective of this study was to empirically demonstrate the superiority of multimodal learning over isolated unimodal approaches. Unlike vision-only or text-only systems, medXpert AI effectively synthesizes visual features extracted by DenseNet121 with semantic embeddings generated by ClinicalBERT. This holistic integration enables the model to capture complex, non-linear relationships between clinical narratives and radiological evidence that are commonly missed by unimodal architectures.

The advantages of this approach were particularly evident in complex diagnostic scenarios. For instance, in cases such as Angina detection, the system correctly prioritized symptom narratives and clinical descriptions even when radiological images appeared normal. By projecting both modalities into a shared latent space, the framework successfully mitigated the modality gap, preserving critical contextual links between reported symptoms and imaging findings that are often lost in isolated diagnostic pipelines.

## C. Explainability Outputs

To address the limitations of black-box artificial intelligence systems, medXpert AI incorporates robust explainability mechanisms that provide transparent and interpretable diagnostic reasoning. For visual data, the system generates Grad-CAM heatmaps overlaid on the original medical images, highlighting the specific anatomical regions that most strongly influenced the model's classification. These visual explanations enable clinicians to verify whether the AI's focus aligns with established medical reasoning.

In parallel, textual attention visualization is applied to clinical notes, highlighting key phrases and medical keywords that contributed most significantly to the final diagnostic decision. The output interface presents an integrated view combining the predicted diagnosis, confidence score, and associated risk level, translating complex model outputs into actionable insights suitable for clinical triage. Furthermore, a human-centric design approach is adopted by including medical disclaimers and simplified explanatory sections, such as "What is this?", ensuring ethical communication and accessibility for healthcare professionals.

**Table 1:**

| Steps | Test Data | Expected Results | Observed Results | Remarks |
|---|---|---|---|---|
| Step 1 | Valid Google account | Login page loads successfully | Login page displayed | Pass |
| Step 2 | Click "Continue with Google" | Redirects to Google authentication | Redirect successful | Pass |
| Step 3 | Valid credentials | User authenticated | Login successful | Pass |
| Step 4 | Invalid credentials | Authentication denied | Access blocked | Pass |
| Step 5 | Logout action | User session ends | Logged out successfully | Pass |

**Table 2:**

| Steps | Test Data | Expected Results | Observed Results | Remarks |
|---|---|---|---|---|
| Step 1 | Supported image format (JPG/PNG) | Image accepted | Image uploaded | Pass |
| Step 2 | Unsupported format | Error message shown | Upload blocked | Pass |
| Step 3 | Empty upload | Warning displayed | Warning shown | Pass |
| Step 4 | Valid medical image | Ready for analysis | Analysis initiated | Pass |
| Step 5 | Large image file | Proper validation | File validated | Pass |

## Discussion

The experimental evaluation of MedXpert AI supports the central hypothesis of this work: combining visual and textual clinical data within a single multimodal framework leads to more reliable diagnostic outcomes than treating each modality in isolation. By allowing radiological evidence and patient-reported symptoms to interact within a shared representation space, the system reflects the way clinicians naturally reason

< 76 >

across multiple sources of information during diagnosis.

## A. Diagnostic Efficacy and Sensitivity

One of the most important observations from the experiments is the consistent improvement in recall achieved by the multimodal fusion model. In clinical decision support systems, reducing false negatives is often more critical than maximizing overall accuracy, as undetected conditions can result in delayed intervention and adverse outcomes. The proposed fusion strategy, inspired by ViLT-style transformers, enabled meaningful interaction between symptom context extracted by ClinicalBERT and visual features learned by DenseNet121.

This interaction proved particularly beneficial in cases where imaging findings were subtle or visually ambiguous. In such scenarios, symptom information helped clarify the diagnostic context, allowing the system to identify pathological conditions that might otherwise be overlooked by vision-only models. These results highlight the practical advantage of multimodal learning over traditional unimodal approaches, especially in early-stage or borderline clinical cases.

## B. Explainability and Clinical Trust

Explainability emerged as a key factor in establishing the system's clinical relevance. Grad-CAM visualizations consistently highlighted medically meaningful regions within radiological images, such as lung opacities or enlarged cardiac structures, rather than irrelevant background areas. At the same time, attention-based analysis of textual inputs emphasized clinically significant terms within patient narratives.

Together, these explanation mechanisms allow clinicians to understand why a particular diagnosis was suggested, rather than simply observing a confidence score. This directly addresses the long-standing concern around "black box" AI systems in healthcare and supports the adoption of MedXpert AI as a transparent and verifiable decision-support tool rather than an opaque predictor.

## C. Robustness through Intelligent Fallback

The confidence-aware decision strategy further enhanced system reliability. When prediction confidence dropped below the predefined threshold, the automatic invocation of the Gemini 1.5 Flash fallback mechanism helped manage uncertainty more safely. Instead of forcing a potentially unreliable diagnosis, the system either refined the assessment or generated clarifying prompts.

This design effectively introduces a human-in-the-loop safeguard, ensuring that ambiguous cases are treated with caution rather than resolved purely through statistical inference. Such behavior is particularly important in safety-critical medical environments, where acknowledging uncertainty is preferable to producing an overconfident but incorrect output.

## D. Limitations and Future Scope

Despite its strengths, the current implementation has certain limitations. The evaluation was primarily conducted using the MIMIC-CXR dataset, which focuses on chest radiographs. While this validates the system's effectiveness for X-ray–based diagnosis, extending the approach to other imaging modalities such as MRI or ultrasound will require additional training and validation.

Moreover, clinical documentation styles can vary significantly across institutions, potentially affecting

generalization in real-world deployment. Although transfer learning mitigates this to some extent, further evaluation on multi-institutional datasets would strengthen the system's robustness. Finally, the fallback reasoning agent currently operates as an external module; future work could explore tighter integration, allowing feedback from the fallback stage to iteratively improve the core fusion model.

## E. Concluding Remarks

Overall, the discussion demonstrates that MedXpert AI effectively addresses three critical challenges in clinical AI systems: fragmented data analysis, lack of interpretability, and uncertainty in low-confidence predictions. By combining multimodal learning with explainable outputs and an intelligent fallback mechanism, the system shows strong potential as a practical and trustworthy clinical decision support tool, rather than merely a proof-of-concept research model.

## Conclusion and Future Directions

This study introduced MedXpert AI, a multimodal clinical decision support system developed to address several long-standing challenges in medical AI, including fragmented data analysis, limited interpretability, and inadequate handling of diagnostic uncertainty. By jointly analyzing radiological images and unstructured clinical narratives within a unified transformer-based fusion framework, the system effectively bridges the modality gap that constrains traditional unimodal approaches. Experimental results demonstrate that combining visual and textual information leads to more reliable diagnostic outcomes, reinforcing the idea that holistic medical analysis requires simultaneous consideration of multiple data sources.

## A. Key Contributions

A central contribution of this work is the empirical evidence that multimodal fusion substantially improves diagnostic recall (sensitivity). In clinical settings, where missed diagnoses can have serious consequences, the system's ability to use symptom context to clarify visually ambiguous imaging findings represents a meaningful advancement toward safer AI-assisted diagnosis. Rather than relying on imaging or text alone, MedXpert AI benefits from the complementary strengths of both modalities.

Another important contribution is the integration of Explainable AI (XAI) mechanisms. The use of Grad-CAM heatmaps and attention-based text visualization allows clinicians to understand which image regions and clinical terms influenced a given prediction. This transparency helps shift the system away from a traditional "black box" model toward a tool that supports verification and informed decision-making, which is essential for acceptance in real clinical environments.

Equally significant is the design of the intelligent fallback mechanism. By incorporating Gemini 1.5 Flash to handle low-confidence predictions, the system introduces a practical human-in-the-loop safeguard. Instead of forcing a potentially unreliable decision, uncertain cases are explicitly identified and escalated for deeper reasoning. This approach aligns well with ethical clinical practices, where uncertainty is acknowledged rather than hidden.

## B. Future Directions

While the current prototype demonstrates strong potential, several directions remain for future development.

< 77 >

- **Expansion of Imaging Modalities:** At present, the system focuses primarily on chest radiographs. Future work will aim to incorporate additional imaging modalities such as MRI and CT scans, enabling the detection of a wider range of conditions, including neurological and musculoskeletal disorders.
- **Holistic Clinical Data Integration:** Diagnostic accuracy could be further improved by integrating structured clinical data, such as laboratory test results and vital signs. Combining these quantitative inputs with imaging and textual narratives would allow for a more comprehensive representation of patient health.
- **Clinical Workflow Interoperability:** For real-world adoption, seamless integration with hospital systems is essential. Future iterations of MedXpert AI will focus on interoperability with Electronic Health Record (EHR) platforms through standards such as HL7 and FHIR, reducing manual data entry and supporting smoother clinical workflows.
- **Adaptive Learning and Feedback Loops:** There is also scope to evolve the fallback mechanism beyond a safety net. By incorporating expert feedback generated during low-confidence cases into the training process, the system could adopt an active learning strategy, allowing continuous improvement of its diagnostic reasoning over time.

**Concluding Remarks**

In summary, MedXpert AI represents a meaningful step toward more transparent, reliable, and integrated clinical decision support systems. By combining multimodal deep learning with explainability and confidence-aware reasoning, the system demonstrates how AI can be designed to complement medical expertise rather than replace it. With further refinement and clinical validation, MedXpert AI offers a scalable foundation for the next generation of intelligent and trustworthy diagnostic assistants.

**References**

1. Zhou Y *et al*., "Unified Multimodal Transformer for Medical Image and Text," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, 345-354.
2. Vaswani *et al*., "Attention is All You Need," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, 2017.
3. Johnson E *et al*., "MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports," *Scientific Data*. 2019;6(317).
4. Ribeiro MT, Singh S and Guestrin C. "Why Should I Trust You?: Explaining the Predictions of Any Classifier," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016, 1135–1144.
5. Shickel B, Tighe PJ, Bihorac A and Rashidi P. "Deep EHR: A Survey of Recent Advances in Deep Learning Techniques for Electronic Health Record (EHR) Analysis," *IEEE Journal of Biomedical and Health Informatics*. 2018;22(5):1589-1604.
6. Huang K, Altosaar J and Ranganath R. "ClinicalBERT: Modeling Clinical Notes and Predicting Hospital Readmission," *arXiv preprint arXiv:1904.05342*, 2019.
7. Huang G, Liu Z, Van L, Der Maaten, and Weinberger KQ. "Densely Connected Convolutional Networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, 4700-4708.
8. Liu X, Qiu H, Li M, Yu Z, Yang Y & Yan Y. Application of multimodal fusion deep learning model in disease recognition. In 2024 IEEE 2nd International Conference on Sensors, Electronics and Computer Engineering (ICSECE). IEEE, 2024, 1246-1250.

< 78 >