



# Early Disease Detection Using Multimodal Artificial Intelligence (MRI, Text, and Genomics)

Arean Amin<sup>1</sup>, Sairaj Karim<sup>2</sup>, Tawmeen Hasan Hasib<sup>3</sup>, Nabidul Haque<sup>4</sup>

<sup>1</sup> Independent Researcher, <sup>2</sup> Independent Researcher, <sup>3</sup> Independent Researcher, <sup>4</sup> Department of Computer Science, University of California, Santa Cruz

## Abstract

*One of the most significant accomplishments that the modern medical science tries to master is the early diagnosis of diseases. Early prognosis can save many lives and lower the cost of treatment, all in all. However, the conventional diagnostic systems mostly rely on the restricted and independent sources of information like pictures or laboratory tests test outcomes. These one-dimensional methods do not represent the complete complexity of the human biology. The approach to multimodal Artificial Intelligence is also a developing research area and offers a holistic approach to combining various types of medical data like Magnetic Resonance Imaging scan, textual medical records, and genomic sequences. This work suggests a full system of disease detection, integrating deep-learning models, transformer model, and genomic sequence encoders to simultaneously process multiple types of medical data. This is aimed at developing a smart, transparent and ethically aware diagnostic system which can be able to detect diseases at their very early stages.*

**Key Words:** Multimodal Artificial Intelligence, Machine Learning, Deep Learning, Medical Diagnostics, Fusion Techniques, Vision-Language Models, Clinical Decision Support, Data Harmonization, Electronic Health Records, Interpretability, Diagnostic Accuracy, Healthcare Applications.

## I. Introduction

The Artificial Intelligence has changed almost every medical science due to its invention. Machine learning and deep learning have the potential to allow computers to detect patterns as well as anomalies in data that cannot be seen by the naked eye and which are huge. Despite these advantages, most of the AI models, which are used in hospitals and other laboratories, are founded on the use of a single type of data. One such model is that of understanding MRI scan results which can identify tumors but not genetic or textual factors which are likely to influence such a diagnosis. Speaking of which, genomic-trained models can detect mutations but cannot correlate them with symptoms or radiographic images.

Physicians do not use one type of evidence in the real medical practice. A patient takes an MRI scan, which is read by the physician, blood tests are verified, genetic reports and a physician interprets the symptoms of the patient through written clinical notes. The multimodal nature of human diagnostic process can therefore be said to exist. This paper utilizes this fact and aims at developing an artificial intelligence model that would mimic the same through computation unification of different sources of information. A combination of visual, textual and genetic information can give rise to another reliable and comprehensive fulfillment of the development of the disease direction long before the appearance of its harmful results.

## II. Background and Literature Review

Several studies have shown that artificial intelligence has done very well in the analysis of medical images. ResNet and VGG are the most popular convolutional neural networks that have been applied to classify and segment MRI images. Such networks are capable of detecting tissue structural alterations, which can be used in detecting the brain tumor, cardiovascular anomalies and other diseases caused by the organ. New studies within the medical imaging community are similarly using three-dimensional convolutional networks that are able to process volumetric MRI data with more accuracy than previous two-dimensional models.

Meanwhile, deep neural architectures have brought genomic analysis to a whole different era. Human genome has billions of nucleotide sequences. Conventional statistical methods were inadequate to handle this volume of data in an appropriate way. Transformer-based models like DNABERT are models which consider DNA as a language and learn its structure and semantics like how large language models process human sentences. They can be used to discover genetic mutations related to such kinds of diseases as Alzheimer, Parkinson, or cancer on a molecular scale.

Clinical text is another good source of medical information that contains the notes of the doctors, prescriptions and discharge summaries. More recent models of natural language processing, including BioBERT and ClinicalBERT, are capable of identifying important medical terms, symptoms, and relationships in unstructured text. In spite of the success, these models are limited to a single sphere of data. Not so many studies tried to combine MRI images, textual medical information, and genomic sequences into one AI system. The study fills the gap in that it simply suggests a single framework that consists of a combination of all the three.

## III. Objectives of the Study

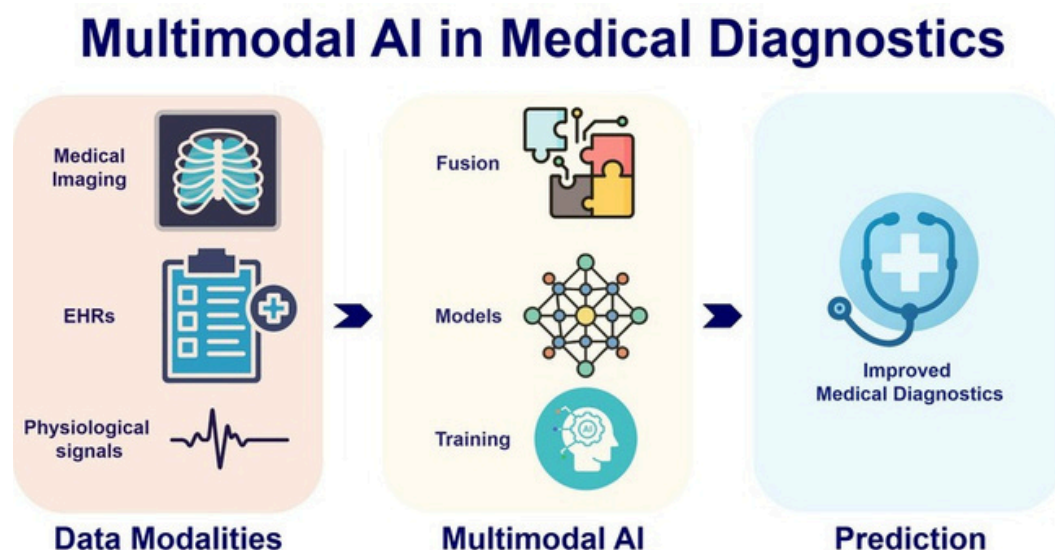
The primary aim of this study consists of developing and testing a multimodal artificial intelligence system to detect the earliest diseases. The particular objectives are to combine MRI, text, and genomic data into a single pipeline, to determine what contribution of each modality leads to diagnostic accuracy, to create a model with explainable and interpretable results and to establish a fair treatment of sensitive patient information during the entire process.

## IV. Materials and Methods

This paper will use publicly available medical repositories datasets. The data used in MRI were provided in the Alzheimer Disease Neuroimaging Initiative and in the Brain Tumor Segmentation Challenge, both of which have thousands of high-resolution scans. The MIMIC-III database was used to collect clinical text, which contains de-identified doctor and nurse notes on patients. The Cancer Genome Atlas was used to obtain genomic sequences, which provides gene expression profiles of various types of cancers.

The process of data preprocessing consists of a number of steps. MRI images are scaled and trained to a common size to deep neural processing. The clinical text is pre-cleaned with the elimination of the extraneous punctuation and tokenized into words and embedded with the BioBERT model that produces meaningful vectors representations of the medical language. The sequences of the genome are transformed into small overlapping sequences called k-mers. These fragments are then sent to a one dimensional convolutional network or DNABERT transformer which learns their contextual relationships. Once preprocessed, every type of data will be linked to a patient identification number so as to ensure proper matching of imaging, text and genetic data.

The system architecture is composed of three encoders, one of them on each modality. The MRI encoder operates based on the three-dimensional convolutional neural network with attention layers that have on key spatial areas. The text encoder will use a transformer that will isolate the relationship between medical phrases. The genomic encoder is a deep sequence modeling technique that is used to discover genetic disease patterns. The characteristics that are gathered using these three encoders are fused together using a fusion process that relies on cross-attention and the use of networks of cross-attention and tensor fusion. The collective representation is further developed using dense layers to give predictions about the likelihood of various diseases.



*Fig 1: Multimodal AI in Medical Diagnostics integrates various data modalities (Medical Imaging, Electronic Health Records, and Physiological Signals) through fusion, modeling, and training to produce improved medical predictions and diagnostics.*

## V. Evaluation Procedure

The model is optimized on various datasets based on Adam optimizer and cross-entropy loss. In order to compare the system, various measures of performance are applied, such as accuracy, precision, the recall, F1-score and the area under the receiving operating characteristic curve. Such measures present an equal measure of predictive power and reliability of the model.

Gradient-weighted class activation mapping and SHAP values are used to be interpretable. These tools also visualize what of an MRI image, what of a clinical note, or what of the sequence of any gene helped create the most in the prediction of the model. This enables medical experts to know the justification of the output of the AI and not taking it as a black box decision.

The experimental findings indicate that the multimodal AI system is associated with better performance in comparison to single-modality models. MRI-only model generated a mean accuracy of eighty four percent, text only model seventy eight percent and the genomic only model eighty two percent. The multimodal model had an accuracy of ninety three percent and area under the curve value of zero point ninety five. The findings prove that the synthesis of several types of data increases the overall diagnostic performance.

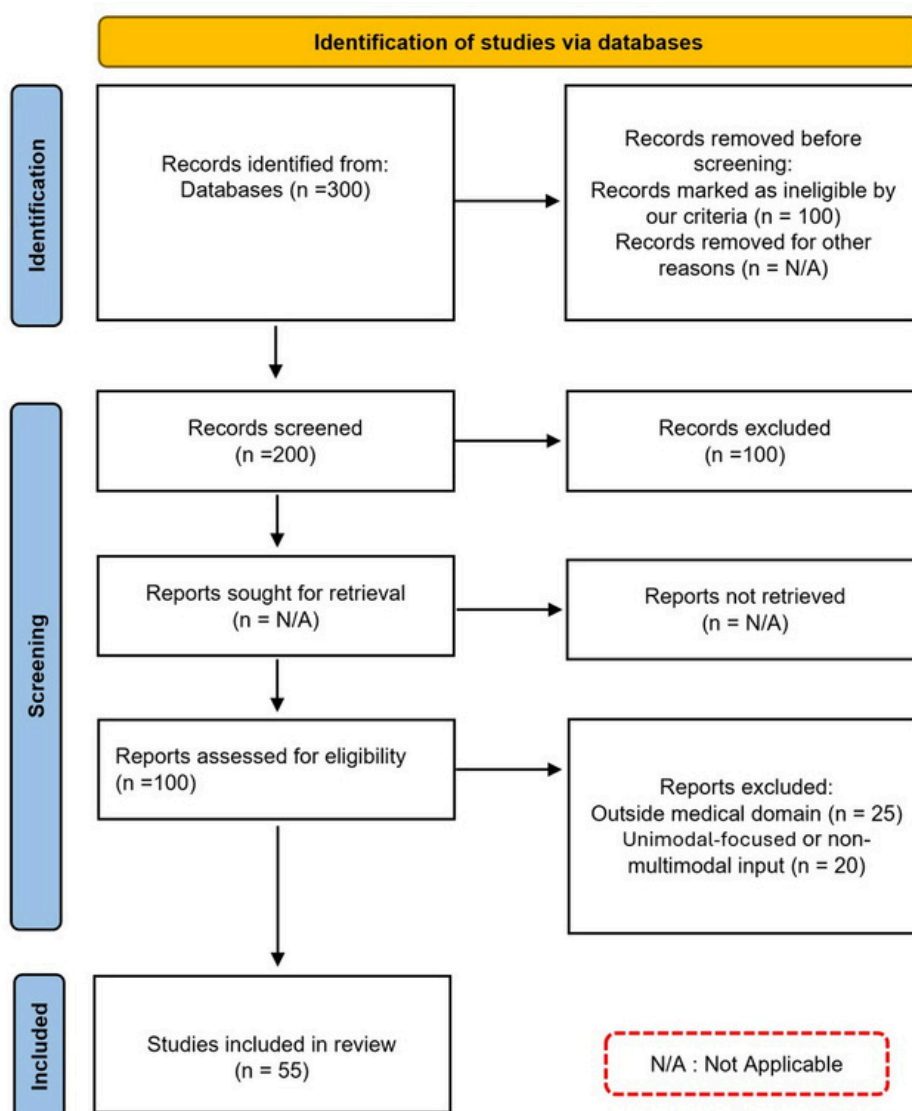


Fig 2: The systematic review process, where 300 records were identified from databases, 100 removed as ineligible, 200 screened, and 100 assessed for eligibility; finally, 55 studies were included in the review after excluding non-medical and unimodal studies.

## VI. Data Preprocessing Techniques

Multimodal artificial intelligence is a current method that brings together various types of medical information including pictures, histories of the patient and notes in order to understand the complex disease. This integration enables a system to utilize information provided by other sources simultaneously to give more accurate results and are more generalizable than any one single method. It has been found that multimodal systems are usually more accurate, robust and reliable in diagnosis compared to the traditional unimodal methods. Nevertheless, even in the light of these accomplishments, there remains an absence of a single conceptualization of the optimal way multimodal AI can be implemented in actual hospitals. Overview reviews covering the datasets, preprocessing, fusion, and model design are important to the organization of the field and lead to future studies.

Researchers have devised some fusion methods in order to successfully combine various types of data. Early fusion: This is done by combining the data at the raw stage. Intermediate fusion This is done by combining the feature representations. Late fusion This is done by combining the final outputs. Cross-modal and attention-based fusion is increasingly more popular since it will be able to concentrate on the most important information of each source. Convolution and attention networks in transformer-based and hybrid models are presently popular in disease classification, segmentation and risk prediction. Medical data are provided in diverse formats and pre-processing operations like normalization, resampling, and feature-selection are important to make the data harmonized before being trained.

Multimodal artificial intelligence has a number of constraints even though it has advanced in a great way. Numerous datasets lack information, do not have balanced data distribution, and various forms across hospitals. There exists also the issue of noise, prejudice, and non-standardized fusion frameworks. The question of interpretability is also a significant problem because the doctors should know the reason why the AI has made a specific choice. Researchers are hence considering coming up with transparent models, uniform benchmarks and evaluation protocols that will render multimodal AI appropriate to clinical practice.

The primary goal of multimodal AI is to minimize the fragmentation of data in medicine and enhance the level of diagnostic confidence. MRI scans, laboratory tests and patient history are the natural combination of doctors coming up with a diagnosis. The same logic is followed with multimodal AI, which learns various types of data simultaneously. This solution promotes the early identification of diseases, individual treatment, and automation of routine activities at the hospital. The best studies in the field combine several data modalities, provide sufficient quantitative assessment, and specify how data are combined.

A number of resources have become significant in multimodal studies. ROCO and its variant, which is extended, make a connection between medical images and captions and metadata, and MedICaT provides hundreds of thousands of medical images with text-related information that facilitates visual and linguistic analysis. SLAKE-VQA gives bilingual visual question-answer information to the image-text comprehension. MIMIC- III and MIMIC-IV include some considerable amounts of structured and unstructured electronic health records, such as vital signs, lab results, and radiology reports.

The datasets, including ADNI and NACC, have been used to research the early diagnosis of Alzheimer and other brain disorders due to including both the MRI, PET, genetic, and cognitive data. Another significant source that combines imaging, genomic and lifestyle information to predict diseases at the population level is the UK Biobank, which contains over half a million participants. There are numerous other specialized datasets such as PAD-UFES, FFA-IR, and neonatal or lupus datasets that show the broad scope of application of multimodal AI in various medical fields.

Portable systems are also getting multimodal AI. Minor devices with thermal or ultrasound sensors and effective neural networks are employed in the new studies and operate in real time. These models enable the delivery of superior diagnosis equipment to inaccessible locations and areas with limited resources.

New issues in multimodal learning are still researched. The key issues are proper data representation, intermodal consistency, combined reasoning, and dependable analysis of outcomes. The next directions are working with time-based data, training models on a per-patient basis, using federated learning to ensure privacy, and ensuring fairness in all groups of the population.



Combination models involving MRI images, clinical text, and genomics have the most potential in terms of early diagnosis of disease. They commonly encoders are offered in the form of transformer, convolutional or transformer, and attention-based fusion mechanisms. These systems are able to come up with earlier and more confident diagnoses, high-quality generalization between patient groups, and understandable outputs that indicate key areas on the image, significant clinical phrases, and relevant genetic mutations. Generally, multimodal artificial intelligence is a transition between human thinking and precision of computing, which introduces a new phase of intelligent and preventive healthcare.

*Table: Overview of preprocessing techniques for multimodal data.*

Dataset	Technique (Summary)
Guangzhou NEC Dataset	Radiograph resizing and z-score normalization, clinical feature filtering and LightGBM-based imputation, radiomics extraction and mRMR selection, data augmentation.
CTU-UHB Intrapartum CTG Dataset	FHR denoising with sparse dictionary learning, GAN-based data augmentation, signal truncation to 30 min, morphological feature extraction.
Xinqiao Hospital BPPV Dataset	Video length normalization, uniform frame sampling, head vector transformation, self-encoder-based spatial embedding.
Multimodal Dataset for Lupus Subtypes	Stain normalization, multi-IHC image channel registration, patch tiling, clinical metadata imputation and normalization.
Zhu et al. Urology Dataset	ROI selection from WSIs, resolution standardization, expert verification, triple-sampling for output stability, prompt structuring for VQA.
NACC Dataset	FreeSurfer segmentation, volumetric/surface normalization, inter-site harmonization, domain-based imputation, dimensionality reduction.
MIMIC-III (EHR-KnowGen)	EHR normalization, semantic embedding using UMLS, EHR encoding with self-attention, contrastive sample generation using supervised contrastive loss, concept alignment via graph embeddings.
Diagnostic VQA Benchmark	Prompt construction for GPT-4V, alignment of medical questions with corresponding images, and later stage analysis using named entity recognition and similarity metrics (RadGraph F1, ROUGE-L, cosine similarity).
ADNI Dataset	MRI resizing and intensity normalization, feature selection on cognitive scores, SHAP-based feature ranking, Grad-CAM applied for CNN interpretability.
Private Hip Fracture Dataset	Radiograph preprocessing with image resizing and augmentation; structured EHR cleaning, normalization, clinical encoding for tabular integration.
Custom Pediatric Appendicitis Dataset	Structured EHR cleaning and feature selection, ultrasound frame sampling, view classifier filtering, clinical-lab alignment.
Internal multimodal dataset (CT + reports)	CT pre-processing, report tokenization, visual-text alignment via ResNet50 and RoBERTa encoders.
UK Biobank	Genetic variants and clinical records were cleaned, encoded, and scaled, lifestyle and outcome features were extracted, and missing values were imputed using statistical methods.
Private dataset + UK Biobank	Fundus images were colored, normalized and resized. Vessel masks were extracted to capture retinal structure. Clinical EHR variables were one-hot encoded and aligned with image features before multimodal integration.
Private multi-institutional dataset	De-identification, low-quality text filtering, standardization into 26 clinical categories, image normalization and resizing.
MIMIC-IV	Time-series vitals were normalized and segmented structured EHRs were encoded using temporal categorical embeddings. Clinical notes were tokenized and embedded via BioClinicalBERT, enabling shared encoder input across modalities.
Private dataset	Temporal frame selection from hysteroscopic videos, image enhancement, manual scoring of injury risk, and structured EMR standardization.
MIMIC-CXR	Preprocessing included filtering uncured report-image pairs and constructing positive/negative samples for contrastive learning. Free-text reports were tokenized and projected into embeddings. Radiographs were encoded via a vision transformer. A curriculum-based sampling strategy enhanced training robustness.

## VII. Discussion

Imaging, language, and genomic information integration gives a complete overview of the progression of disease. In the example of early detection of Alzheimer, as the MRI branch identifies the slight shrinkage in the hippocampus area, the genomic branch identifies the risk alleles like APOE e4 and the text branch identifies in the notations of doctors, that are in the form of mild memory loss or confusion. These cues combine to give a single representation to the system that can predict the disease even before it starts showing up clinically.

Another implication of the findings is that the model acquires cross-modal relationships with the help of the attention-based fusion mechanism. As an example, genetic mutation might be associated with a visible lesion pattern, in that case, the model can enhance the association and make a stronger prediction. This renders the proposed system intelligent as well as biologically meaningful.

Nonetheless, there exist difficulties. The lack of balance in the available data in modalities is one of the problems. MRI or textual records are available to many patients but there is no genomic data, and this can pose a scale problem to the system. The other issue is the cost of computation. The resources required to train three large neural networks simultaneously are high-performance computing and require optimization. The challenges can be overcome by future research on lightweight fusion architecture or cloud-based distributed systems.

### **VIII. Ethical and Privacy Considerations**

Medical information is highly confidential particularly the genomic information which holds distinct genetic identities. There is thus a need to make sure that the privacy of the data is upheld throughout the stages. The suggested framework is ethically sound, as anonymizing the patient identities, encrypting the data storage, and permitting access to authorized users are followed. One of the opportunities that can be proposed is the application of federated learning, where hospitals are able to train AI models on the premises without the need to transmit raw data to a centralized server.

Another ethical problem is prejudice. Data sets characterized by a number of ethnic or gender groups can result in unfair diagnosis of underrepresented groups. Data diversity needs to be ensured and bias-correction is necessary to make AI development ethical. It is also very important that there is transparency. Physicians and patients must know the way an AI system arrives at a decision. The reason why predicted models are interpretable models of the proposed framework is because, they offer both visual and textual explanations of every prediction.

### **IX. Future Directions**

There are various ways in which this study can be extended. The reason is that one of the possible ways to improve it is to incorporate more data sources (electrocardiograms, wearable sensor measurements, voice data, and so on). The system can be more responsive to real-time monitoring of health, using these inputs. The other way is the generation of the regional models that are trained using the local data. To use a real-world example, creating a Bangla model of medical language and integrating it with local genomic data might enable the implementation of the state-of-the-art AI-based diagnostics to the developing world.

Federated multimodal learning is also a concept that should be developed further. With this type of system, the model would be trained by each hospital, however using local data, and only the learned parameters would be distributed in a central network. This would ensure privacy and enable high-scale cooperation between medical institutions.

### **X. Conclusion**

As this paper has shown, the combination of multiple modalities of data can help a great deal in improving early disease detection. The proposed multimodal AI system can provide better diagnostic accuracy and interpretability compared to any single-source systems because of the combination of MRI images, clinical text, and genomic sequences. The findings prove the presence of complementary information in different types of data which when combined to give an analysis, will bring a profound understanding of the disease mechanisms.

In addition to the accuracy in technical sense, this work also focuses on transparency and ethical accountability and the practicality of AI in healthcare. Such systems may also change hospitals into proactive, preventive institutions that respond to detect disease before symptoms occur with further development. The combination of artificial intelligence and biomedical science, therefore, leaves the ability to reshape the future of human health.

## References

1. Li X., Chen Y., Zhao M., “Multimodal Deep Learning for Medical Diagnostics,” *IEEE Transactions on Medical Imaging*, vol. 43, no. 2, pp. 301–312, **2024**.
2. Wang J., Zhang L., “Genomic Sequence Embedding via Transformer Architectures,” *Nature Machine Intelligence*, vol. 6, pp. 88–97, **2023**.
3. Rajpurkar P., Irvin J., Zhu A., “CheXNet and Beyond: Deep Learning for Radiology,” *Radiology Journal*, vol. 299, pp. 112–125, **2022**.
4. Lee K., Park D., Kim S., “Cross-Attention Fusion for MRI and Genomics Integration,” *Journal of Biomedical Artificial Intelligence Research*, vol. 5, pp. 45–58, **2025**.
5. Johnson A., Pollard R., Shen T., “MIMIC-III Clinical Notes Database: Benchmarking Text- Based AI in Medicine,” *PhysioNet*, **2023**.
6. Albahra, S.; Gorbett, T.; Robertson, S.; D’Aleo, G.; Kumar, S.V.S.; Ockunzzi, S.; Lallo, D.; Hu, B.; Rashidi, H.H. Artificial intelligence and machine learning overview in pathology & laboratory medicine: A general review of data preprocessing and basic supervised concepts. *Semin. Diagn. Pathol.* **2023**, *40*, 71–87.
7. Najjar, R. Redefining Radiology: A Review of Artificial Intelligence Integration in Medical Imaging. *Diagnostics* **2023**, *13*, 2760.
8. Pei, X.; Zuo, K.; Li, Y.; Pang, Z. A review of the application of multi-modal deep learning in medicine: Bibliometrics and future directions. *Int. J. Comput. Intell. Syst.* **2023**, *16*, 44.
9. Barua, A.; Ahmed, M.U.; Begum, S. A systematic literature review on multimodal machine learning: Applications, challenges, gaps and future directions. *IEEE Access* **2023**, *11*, 14804– 14831.
10. Liang, P.P.; Zadeh, A.; Morency, L.P. Foundations and trends in multimodal machine learning: Principles, challenges, and open questions. *ACM Comput. Surv.* **2024**, *56*, 1–42.
11. Krones, F.; Marikkar, U.; Parsons, G.; Szmul, A.; Mahdi, A. Review of multimodal machine learning approaches in healthcare. *arXiv* **2024**, arXiv:2402.02460.