



COMPREHENSIVE REVIEW: CANCER TYPE DETECTION USING ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING

Prateek Sikarwar¹, Saurabh Singh², Aman Singh³, Rohit Sharma⁴

UG Student, Department of Computer Science & Engineering, Raja Balwant Singh Engineering Technical Campus,
Agra, India¹

UG Student, Department of Computer Science & Engineering, Raja Balwant Singh Engineering Technical Campus,
Agra, India²

Assistant Professor, Department of Computer Science & Engineering, Raja Balwant Singh Engineering Technical
Campus, Agra, India³

Assistant Professor, Department of Computer Science & Engineering, Raja Balwant Singh Engineering Technical
Campus, Agra, India⁴

Abstract: In recent years, the use of Artificial Intelligence in healthcare is increasing, mainly for detecting cancer. Early and correct diagnosis of cancer plays a crucial role for better treatment results and for lowering the number of deaths. This work describes a deep learning system that is designed to identify and categorize different cancers like brain tumor, skin cancer, lung cancer, and breast cancer, using medical imaging data.

The system uses Convolutional Neural Networks (CNNs), which helps in analyzing images and automatically extract important information from images. Different datasets are collected and organized into groups such as benign, malignant, and normal. Before training, images undergo preparation steps such as resizing, normalization, and data augmentation to improve results and reduce overfitting.

Different CNN models are trained for each cancer type with TensorFlow and Keras frameworks. The performance of these models is measured using metrics such as accuracy and loss. Experimental results show that some models achieved high accuracy (approximately 85–90%), while others demonstrated moderate performance due to challenges such as limited dataset size and class imbalance.

To enhance usability, a simple and interactive graphical user interface (GUI) is developed, allowing users to upload medical images and obtain real-time predictions along with confidence scores. Additionally, an invalid image detection mechanism is incorporated to prevent incorrect predictions for unrelated inputs, thereby improving system reliability.

Overall, this paper demonstrates the effectiveness of deep learning in cancer detection while also highlighting key challenges such as data limitations and model generalization. The proposed system can serve as a foundational framework for future research and can be further improved using advanced architectures, larger datasets, and real-time deployment strategies for practical healthcare applications.

Keywords: Cancer Detection, Medical Imaging, Artificial Intelligence, Machine Learning, Deep Learning, Convolutional Neural Networks (CNN), Healthcare Analytics, Tumor Classification, Clinical Decision Support Systems.

I. INTRODUCTION

Cancer is a serious global health issue and early detection is key to improving the chances of recovery. Over the past few years healthcare has increasingly adopted AI and ML to help doctors in identifying diseases and making decisions, specially through deep learning in medical imaging. Medical imaging techniques such as MRI, CT scans, and



dermoscopic images play an important role in detecting different forms of cancer. Manual analysis of medical images by doctors can be time-consuming and mistakes can occur, which is why deep learning-based automated systems are used for faster and more accurate results.

In this paper, a deep learning-based system is developed for the detection and classification of multiple types of cancers, including brain tumor, skin cancer, lung cancer, and breast cancer. The system uses Convolutional Neural Networks (CNN), which are highly effective for image classification tasks. Separate models are trained for each type of cancer using structured datasets.

The main objective of this paper is to build a user-friendly and efficient system that can take an input image and predict whether the image belongs to a specific cancer category or not. Additionally, an invalid image detection mechanism is also implemented to avoid incorrect predictions when unrelated images are provided. This work not only demonstrates the potential of deep learning in medical image analysis but also highlights the challenges such as limited dataset size, class imbalance, and variation in model performance. The system can be further improved by using advanced models, larger datasets, and real-time deployment techniques.

II. LITERATURE REVIEW

In recent years, the field of healthcare has experienced significant transformation due to the advancement of Artificial Intelligence and Machine Learning, particularly in the area of cancer detection and diagnosis. Researchers are actively focusing on enhancing the accuracy and efficiency of prediction systems by applying various techniques, including deep learning, hybrid models, and ensemble learning approaches. This section provides a comprehensive review of different research studies related to cancer prediction systems.

Chen et al. (2024) [1] proposed a multi-modal deep learning approach that combines histopathology images and genomic data for cancer classification. Their model used a hybrid CNN-MLP architecture, where CNN was responsible for extracting image features and MLP handled the structured genomic data. The study showed that combining multiple types of data improves the model's ability to understand complex cancer patterns. This approach also helps in reducing errors and increasing prediction reliability. The authors concluded that multi-modal systems are highly effective and represent the future of advanced cancer diagnosis systems.

Almarri et al. (2024) [2] developed a Breast Cancer Prediction Model (BCPM) using a combination of Decision Tree and Random Forest algorithms. Their model was designed in such a way that it can be used in real-time hospital environments. The system achieved an impressive accuracy of around 98%, which shows its strong performance. The study also highlighted that machine learning models must not only be accurate but also easy to interpret for doctors. This makes their approach practical and suitable for real-world clinical use.

Hajjar et al. (2024) [3] introduced a multi-cancer detection system using ensemble learning methods such as Random Forest, XGBoost, and Gradient Boosting. Their model was trained on a large dataset containing different cancer types. The results showed that ensemble techniques can significantly improve detection accuracy compared to single models. The system was able to identify breast, lung, and colon cancers with high sensitivity. The authors concluded that combining multiple models helps in reducing errors and increasing stability.

Zhou et al. (2024) [4] focused on making AI models more understandable by using Explainable AI techniques such as SHAP and LIME. Their study addressed one of the biggest problems of deep learning, which is the lack of transparency. By explaining how the model makes decisions, doctors can better trust the system. The research showed that explainability is very important in medical applications. This approach helps in building confidence among users and supports better decision-making.

Unger et al. (2024) [5] worked on deep learning models for analyzing cancer genomics and histopathology images. They used CNN for extracting spatial features and RNN for understanding sequential patterns in data. Their results showed that deep learning models perform better than traditional machine learning methods. The combination of different neural networks helped in improving classification accuracy. The study also highlighted the importance of feature extraction in medical data analysis.

Kumar et al. (2024) [6] developed an advanced deep learning framework for multi-cancer image classification using medical imaging data. Their model utilized modern Convolutional Neural Networks (CNN) along with optimized feature



extraction techniques to automatically detect patterns from different cancer types. The system was designed to classify multiple cancers using a single integrated model, making it more efficient compared to traditional single-cancer systems. **Alharbi (2023) [7]** studied the use of machine learning algorithms such as Decision Tree, Random Forest, and Support Vector Machine for cancer classification. The research used gene expression datasets, which are usually very complex and high-dimensional. The results showed that ensemble models like Random Forest performed better in terms of accuracy and stability. The study also emphasized the importance of preprocessing and feature selection. Proper handling of data can greatly improve model performance.

Mokoatle et al. (2023) [8] conducted a comparative study of different machine learning models using clinical and genomic data. Their results showed that hybrid models combining Random Forest and Gradient Boosting achieved the best performance. The study also discussed challenges such as overfitting and data imbalance. They suggested that proper model tuning is necessary to avoid such issues. This research highlights the importance of combining different techniques for better results.

Lee et al. (2023) [9] developed a deep learning system using CNN and Autoencoders for cancer classification. Their model was designed to automatically extract important features from RNA-seq data. This reduced the need for manual feature engineering. The results showed improved accuracy across multiple cancer types. The study also highlighted that deep learning models are very useful for handling high-dimensional data.

Avila and Deepa (2023) [10] proposed a hybrid approach combining CNN with Principal Component Analysis (PCA). PCA was used to reduce the number of features and remove unnecessary data. This helped in improving computational efficiency and reducing overfitting. The model showed better performance compared to traditional algorithms. The study proved that combining statistical techniques with deep learning can enhance results.

Wang et al. (2023) [11] developed an explainable machine learning model using SHAP values. Their approach focused on identifying the most important features that influence predictions. This makes the model more transparent and understandable. The study showed that explainable models are more reliable in medical applications. It also helps doctors to trust AI-based systems.

Kumar and Singh (2022) [12] designed a hybrid deep learning model combining CNN and LSTM. CNN was used for extracting image features, while LSTM handled sequential patterns. Their model achieved an accuracy of 96.4%, which is quite high. The study showed that combining different neural networks can improve performance. This approach is useful for complex medical data.

Alakwaa et al. (2022) [13] proposed a deep neural network model for cancer prediction using gene expression data. Their system integrated feature extraction and classification in one pipeline. This reduced complexity and improved accuracy. The study showed that deep learning models perform better than traditional methods. It also highlighted the importance of handling high-dimensional data properly.

Lu et al. (2022) [14] introduced a hybrid model combining CNN and MLP for multi-modal cancer classification. Their approach used both image data and genomic data. The results showed improved accuracy and better generalization. The study concluded that data integration is very important for cancer prediction systems. This approach helps in capturing more information.

Way et al. (2022) [15] developed a machine learning model for pan-cancer prediction using TCGA dataset. Their model was able to identify patterns across different cancer types. The study showed that machine learning can be used for cross-cancer analysis. This approach is useful for large-scale cancer detection systems.

Kourou et al. (2021) [16] presented a review of machine learning techniques used in cancer diagnosis. The study discussed various algorithms such as SVM, Random Forest, and Neural Networks. It highlighted the importance of feature selection and data preprocessing. The research concluded that proper data handling is very important for achieving good accuracy.

Hossain et al. (2021) [17] proposed a multi-task learning model for cancer classification. Their approach allowed the model to learn from multiple datasets at the same time. This improved generalization and accuracy. The study showed that multi-task learning is useful for complex problems. It also helps in reducing overfitting.



Khairi et al. (2021) [18] compared different deep learning models such as VGG16, ResNet50, and DenseNet121. Their study showed that deeper networks provide better feature extraction. However, they also require more computational power. The research highlighted the trade-off between accuracy and complexity.

Wang et al. (2021) [19] used transfer learning techniques with pretrained models like ResNet50 and InceptionV3. Their approach helped in improving accuracy even with small datasets. The study showed that transfer learning is very useful in medical applications. It reduces training time and improves performance.

Huang et al. (2021) [20] developed a CNN-based model for lung cancer detection using CT scan images. Their model was able to detect small tumor regions with high accuracy. The study showed that deep learning is very effective for image-based cancer detection. It also highlighted the importance of high-quality data.

III. RESEARCH GAP

After analyzing our paper on multi-disease cancer detection (Brain, Skin, Breast, Lung) using deep learning models, several research gaps have been identified:

First, the current system uses separate models for each disease, which increases complexity and makes deployment difficult. There is no unified model that can handle multiple diseases efficiently in a single pipeline. This creates a gap in scalability and real-world usability.

Second, although the models provide decent accuracy, they sometimes give high-confidence predictions even for irrelevant or unknown images. This indicates a lack of proper unknown/invalid image detection, which is very important in medical applications to avoid misleading results.

Another important gap is related to dataset size and diversity. The datasets used in the paper are relatively small and limited in variation. Because of this, the models may not generalize well to real-world medical images, leading to overfitting and unstable predictions.

Additionally, the paper mainly focuses on accuracy as the evaluation metric. Other important metrics such as precision, recall, F1score, and confusion matrix are not fully analyzed. This creates a gap in understanding model performance, especially for critical classes like malignant tumors.

There is also a limitation in model architecture. The paper uses basic CNN models, but more advanced techniques like transfer learning (ResNet, VGG, EfficientNet) are not utilized. This limits the potential improvement in accuracy and robustness.

From a system perspective, the current implementation works as a local GUI-based application. However, it lacks real-time deployment, cloud integration, and scalability features, which are essential for practical healthcare applications.

Finally, although an attempt has been made to handle invalid inputs, the approach is not fully reliable. The model still sometimes predicts on non-medical images, showing the need for a more robust validation or hybrid detection mechanism.

IV. CHALLENGES AND FEATURES OF CANCER TYPE DETECTION

A. Challenges

During the implementation of the proposed cancer detection system, several challenges were identified that impacted the performance and reliability of the models.

- **Dataset Limitations:** The datasets used for brain tumor, skin cancer, breast, and lung cancer detection were limited in size and lacked diversity. Due to this, the models were unable to generalize effectively on unseen data. Additionally, some datasets had class imbalance, which affected prediction accuracy.
- **Accuracy Variation Across Models:** Different models showed inconsistent performance. For instance, the skin cancer model achieved relatively high accuracy, while the lung cancer model showed significantly lower accuracy. This variation made it difficult to maintain consistency across the system.
- **Overfitting Issue:** In several cases, training accuracy increased rapidly while validation accuracy remained lower. This indicates that the model was overfitting the training data and failing to generalize well.



- **Invalid Image Prediction Problem:** Initially, the system predicted results even for unrelated or unknown images with high confidence. This is a critical issue in medical applications, as incorrect predictions can lead to misleading conclusions.
- **Multi-Model Integration Complexity:** Separate models were developed for different cancer types. Combining these models into a single system increased complexity and required additional logic for proper prediction handling.
- **Computational Constraints:** The training process was performed on CPU-based systems without GPU support, which increased execution time and limited the ability to use more advanced deep learning models.

B. Features

The proposed system incorporates several important features that enhance its usability and effectiveness.

- **Multi-Disease Detection Capability:** The system can detect multiple types of cancer, including brain tumor, breast, skin cancer, and lung cancer, using dedicated models for each category.
- **Deep Learning-Based Approach:** Convolutional Neural Networks (CNNs) are used for image classification, enabling automatic feature extraction and improved performance compared to traditional methods.
- **User-Friendly Graphical Interface:** A GUI-based application is developed where users can upload images and receive predictions easily. The interface is designed to be simple, visually appealing, and user-friendly.
- **Real-Time Prediction System:** The system provides quick predictions for input images and displays the result along with confidence score in real time.
- **Invalid Image Detection Mechanism:** An additional validation step is implemented to identify unknown or irrelevant images and prevent incorrect predictions, thereby improving system reliability.
- **Modular System Design:** Each cancer detection model is developed independently, making the system modular and allowing easy future enhancements or additions.
- **Model Saving and Deployment Support:** Trained models are saved and reused for prediction, eliminating the need for retraining and making the system efficient for deployment.

Table I: Shows various features with their challenges

S.no	Author & Year	Methodology	Key Features	Challenges
1.	Chen et al. (2024)	Multi-modal DL (CNN + MLP) using histopathology + genomic data	High accuracy across multiple cancer types; data fusion improves prediction	Requires large datasets; high computational cost
2.	Almarri et al. (2024)	Decision Tree + Random Forest (BCPM)	Achieved 98% accuracy; suitable for real-time hospital use	Limited generalization on diverse datasets
3.	Hajjar et al. (2024)	Ensemble models (RF, XGBoost, Gradient Boosting)	High sensitivity for multi-cancer detection	Complex training and tuning
4.	Zhou et al. (2024)	XAI using SHAP & LIME	Improves interpretability and clinical trust	Extra computational overhead
5.	Unger et al. (2024)	CNN + RNN for imaging & genomic data	Extracts spatial & temporal features effectively	Needs large labeled datasets



6.	Alharbi (2023)	ML models (DT, RF, SVM)	Ensemble improves robustness & accuracy	Imbalanced dataset handling issue
7.	Mokoatle et al. (2023)	RF + Gradient Boosting hybrid	High accuracy with clinical data	Overfitting & data scarcity
8.	Lee et al. (2023)	CNN + Autoencoder (RNA-seq data)	Automatic feature extraction	High computation cost
9.	Avila and Deepa (2023)	CNN + PCA	Reduces overfitting; efficient computation	Risk of losing important features
10.	Wang et al. (2023)	Explainable ML with SHAP	Improves transparency & trust	Complex interpretation
11.	Kumar and Singh (2022)	CNN + LSTM hybrid	Captures spatial & temporal features; high accuracy	Complex architecture
12.	Alakwaa et al. (2022)	Deep Neural Network	Integrated feature extraction + classification	Overfitting risk
13.	Lu et al. (2022)	CNN + MLP (multi-modal)	Better accuracy using data fusion	Integration complexity
14.	Way et al. (2022)	ML on TCGA dataset	Identifies molecular patterns across cancers	Limited interpretability
15.	Kourou et al. (2021)	Review of ML models (SVM, RF, ANN)	Highlights importance of preprocessing	Lack of standard evaluation
16.	Hossain et al. (2021)	Multi-task learning framework	Better generalization & performance	Complex model design



17.	Khairi et al. (2021)	CNN (VGG16, ResNet50, DenseNet)	Deep models give better features	Requires large datasets
18.	Alakwaa et al. (2021)	Deep Neural Network	Improves prediction accuracy	High computational cost
19.	Wang et al. (2021)	CNN + Transfer Learning	Works well with small datasets	Model dependency on pretrained data
20.	Huang et al. (2021)	CNN (VGG16, ResNet50)	Accurate lung nodule detection; reduced training time	Requires high-quality CT data

Table I. describes characteristics of study participants: Author(s), Methodology, Key Features, and Challenges Used in the above table.

V. DATASET BASED COMPARISON

This section introduces each dataset used in the study, detailing their source, size, key attributes, and any specific characteristics, describes the specific methods and statistical tests applied to the data for analysis and comparison.

Table II: Shows Various Dataset Used and Their Evaluation.

S.no.	Author (Year)	Methodology	Dataset Used	Evaluation (%)
1.	Chen et al. (2024)	CNN + MLP (Multi-modal DL)	Histopathology + Genomic Dataset	95%
2.	Almarri et al. (2024)	Decision Tree + Random Forest	Breast Cancer Dataset (Benchmark)	98%
3.	Hajjar et al. (2024)	Ensemble (RF, XGBoost, GB)	Multi-cancer Dataset	94%
4.	Zhou et al. (2024)	Explainable AI (SHAP, LIME)	Cancer Genomic Dataset	91%
5.	Unger et al. (2024)	CNN + RNN	Histopathology + Genomic Data	93%
6.	Alharbi (2023)	DT, RF, SVM	Gene Expression Dataset	90%



7.	Mokoatle et al. (2023)	RF + Gradient Boosting	Clinical + Genomic Dataset	92%
8.	Lee et al. (2023)	CNN + Autoencoder	RNA-seq Dataset	94%
9.	Avila and Deepa (2023)	CNN + PCA	Gene Expression Dataset	93%
10.	Wang et al. (2023)	Explainable ML (SHAP)	Genomic Dataset	91%
11.	Kumar and Singh (2022)	CNN + LSTM	Cancer Imaging Dataset	96.4%
12.	Alakwaa et al. (2022)	Deep Neural Network	Gene Expression Dataset	95%
13.	Lu et al. (2022)	CNN + MLP	CT Scan + Genomic Data	96%
14.	Way et al. (2022)	ML Model	TCGA Dataset	92%
15.	Kourou et al. (2021)	ML (SVM, RF, ANN)	Multiple Cancer Datasets	90%
16.	Hossain et al. (2021)	Multi-task Learning	Multi-cancer Dataset	94%
17.	Khairi et al. (2021)	CNN (VGG16, ResNet50)	Histopathology Dataset	93%
18.	Alakwaa et al. (2021)	Deep Neural Network	Gene Expression Dataset	92%
19.	Wang et al. (2021)	CNN + Transfer Learning	Medical Image Dataset	95%
20.	Huang et al. (2021)	CNN (VGG16, ResNet50)	CT Scan Lung Dataset	94%

Table II. describes characteristics of study participants: Author(s), Methodology, Dataset, and Evaluation Metrics used in the analysis.

VI. CONCLUSION

Our research focuses on refining cancer diagnosis through the integration of Deep Learning and AI. By training Convolutional Neural Networks (CNNs) on diverse medical datasets, we successfully created a framework to identify malignancies in the brain, skin, lungs, and breasts. The data suggests that these tools don't just speed up process-they provide doctors with a reliable second opinion that could significantly improve early intervention.



We didn't just build a model; we built a complete user-ready system. By streamlining the entire process from data prep to the user interface, we made sure the technology is actually usable for people who aren't tech experts. The system worked great for skin and breast cancer detection, mostly because those datasets were quite strong. On the other hand, our lung cancer predictions weren't quite as sharp, mainly because the available data was a bit sparse and skewed. It was a clear reminder of how much the quality of data impacts the final outcome.

We ran into a few tough realities while conducting this research. The biggest headaches were definitely the small sample sizes and the constant battle against overfitting, not to mention how tricky it is to juggle different cancers. One thing became very clear in medicine, a high accuracy percentage doesn't tell the whole story because you cannot afford to be wrong. We added an invalid image detection feature to make the system more dependable, and while that helped a lot, there's definitely still room to grow and improve.

Ultimately, this work proves that AI has a permanent place in the future of healthcare, especially when it comes to spotting early warning signs that might be missed manually. But at the end of the day, it is a tool for doctors, not a substitute for them. By handling the data-heavy side of diagnostics, the system allows professionals to focus on what matters most—making informed, life-saving decisions with a bit more confidence.

The next phrase of this project is all about making the system truly 'hospital-ready'. We plan to move past current limitations by feeding the model much larger datasets and experimenting with more advanced deep learning setups. By adopting tougher evaluation benchmarks, we want to make sure the tool is reliable and consistent enough for doctors to actually trust it in a real-world healthcare setting.

REFERENCES

- [1] X. Chen et. al., "Multi-modal deep learning for cancer type classification using genomic and image data," *IEEE Transactions on Biomedical Engineering*, vol. 71, no. 5, pp. 1452–1463, 2024.
- [2] B. Almarri et. al., "The BCPM method: Decoding breast cancer with machine learning," *BMC Medical Imaging*, vol. 24, no. 1, pp. 1–12, 2024.
- [3] M. Hajjar et. al., "Machine learning approaches in multi-cancer early detection," *Information*, vol. 15, no. 2, pp. 1–15, 2024.
- [4] L. Zhou et. al., "Explainable AI for cancer prediction: SHAP and LIME," *Artificial Intelligence in Medicine*, vol. 148, pp. 102–115, 2024.
- [5] S. Unger et. al., "Deep learning for cancer genomics and histopathology," *IEEE Access*, vol. 12, pp. 55678–55690, 2024.
- [6] A. Alharbi, "Comparative analysis of ML models for cancer classification," *International Journal of Intelligent Computing*, vol.12, no. 3, pp. 210–220, 2023.
- [7] M. M. Mokoatle et. al., "ML techniques for cancer detection," *BMC Bioinformatics*, vol. 24, no. 1, pp. 1–13, 2023.
- [8] M. Lee et. al., "Deep learning with genomic data for cancer subtype," *Biology*, vol. 12, no. 4, pp. 1–14, 2023.
- [9] P. Avila et. al., "Deep learning model for cancer classification," *Journal of Population Therapeutics*, vol. 30, no. 2, pp. 89–98, 2023.
- [10] Z. Wang et. al., "Explainable AI framework for cancer prediction," *Artificial Intelligence in Medicine*, vol. 135, pp. 102–110, 2023.
- [11] R. Kumar et. al., "Hybrid CNN-LSTM model for cancer prediction," *Expert Systems with Applications*, vol. 195, pp. 116–128, 2022.
- [12] Y. Lin et. al., "ML in diagnosis of cancer," *Translational Cancer Research*, vol. 11, no. 6, pp. 1800–1810, 2022.
- [13] Y. Lu et. al., "Hybrid CNN + MLP framework," *Computers in Biology and Medicine*, vol. 145, pp. 105–115, 2022.
- [14] G. P. Way et. al., "ML detects pan-cancer activation," *Cell Reports*, vol. 2, no. 3, pp. 456–468, 2012.
- [15] K. Kourou et. al., "ML applications in cancer prediction," *Computational and Structural Biotechnology Journal*, vol. 19, pp.554–565, 2021.
- [16] M. A. Hossain et. al., "Multi-task learning for cancer classification," *Bioinformatics*, vol. 37, no. 14, pp. 2005–2012, 2021.
- [17] M. T. Khairi et. al., "Deep learning breast cancer detection," *Diagnostics*, vol. 11, no. 5, pp. 1–12, 2021.
- [18] F. M. Alakwaa et. al., "Deep learning predicts cancer types," *Frontiers in Genetics*, vol. 12, pp. 1–10, 2021.
- [19] Y. Wang et. al., "CNN-based cancer detection using transfer learning," *IEEE Journal of Biomedical and Health Informatics*, vol. 25, no. 7, pp. 2675–2684, 2021.
- [20] J. Huang et. al., "Lung cancer detection using deep CNN models," *IEEE Access*, vol. 9, pp. 34567–34578, 2021.