# Early-Stage Detection of Colorectal Cancer using Image Classification

Ashlin Santhosh
*Division of Computer Science and Engineering*
*Karunya Institute of Technology and Sciences*
Coimbatore, Tamil Nadu, India
ashlinsanthosh@karunya.edu.in

Anusha Bamini A M
*Division of Computer Science and Engineering*
*Karunya Institute of Technology and Sciences*
Coimbatore, Tamil Nadu, India
anushabamini@karunya.edu

Rajeswari M
*Division of Computer Science and Engineering*
*Karunya Institute of Technology and Sciences*
Coimbatore, Tamil Nadu, India
rajeswari@karunya.edu

Brindha D
*Division of Computer Science and Engineering*
*Karunya Institute of Technology and Sciences*
Coimbatore, Tamil Nadu, India
brindha@karunya.edu

Chitra R
*Division of Computer Science and Engineering*
*Karunya Institute of Technology and Sciences*
Coimbatore, Tamil Nadu, India
chitrar@karunya.edu

Stewart Kirubakaran S
*Division of Computer Science and Engineering*
*Karunya Institute of Technology and Sciences*
Coimbatore, Tamil Nadu, India
stewart@karunya.edu

*Abstract –Early detection of colorectal cancer is important as it is one of the most common and deadliest types of cancer worldwide. This project proposes an image classification model based on deep learning architecture that could detect colorectal cancer in its preliminary stages by using an EfficientNetB0. By using the Colorectal Cancer Whole Slide Images (WSI) dataset, a classification model using the WSI dataset to classify tissue samples into six distinct categories of varying grades of cancerous growths and normal. In this model, transfer learning is included to maintain computational efficiency while enhancing accuracy. A user interface with Streamlit, which has real-time image upload, analysis, and prediction display, was also developed. The predictions come with confidence scores, and the Interface also provides educational content for each cancer stage to help raise awareness and prevent the disease. Applications with potential for medical diagnosis and as an educational tool for early detection of colorectal cancer are possible for this system.*

**Keywords: Colorectal Cancer, Early Detection, Image Classification, EfficientNetB0, Deep Learning**

## I.INTRODUCTION

Colorectal cancer (CRC) is one of the major global health crises, causing about 900000 deaths annually and being responsible for about 10 percent of all cancer cases. However, detection early is important, as CRC does not usually present with symptoms in its early stages, early intervention is key in helping reduce treatment outcomes and survival. Today, although effective, traditional diagnostic methods such as colonoscopy and histopathological analysis are invasive, costly, and prone to human errors. That calls for the development of automated, noninvasive diagnostic tools for early CRC detection.

Recent research has been conducted recently on the application of artificial intelligence (AI) in the analysis of medical images in cancer detection and classification. Specifically, Convolutional Neural Networks (CNNs) have strong performance in feature extraction from histopathological images with the ability for label discrimination between cancerous and noncancerous tissues. Here, an extremely efficient CNN model called EfficientNetB0 was used as the basis of an image classification system to identify and classify CRC at early stages. Our model trains using transfer learning of pre-trained ImageNet weights and fine-tuning with special CRC data so our model can rely on general image features as well as its medical pattern characteristics.

To enable the mainstream use of visual analytics for population data, a system that includes data preprocessing, transfer learning, and an interactive user interface was developed with Streamlit. The image is resized, normalized, and classified into one of six categories corresponding to different cancer stages and normal tissue. Users can upload images, view predictions with primary and secondary confidence scores, and access educational content related to each cancer stage and recommended preventative health measures. One especially notable contribution of the interface is that the dual confidence display gives users a new type of interpretability in the form of understanding the model uncertainty.

## II. LITERATURE SURVEY

Focusing on research in early diagnosis of colorectal cancer (CRC) using AI and machine learning reflects substantial progress made in quantifiable early diagnosis. Combined, they point to machine learning models, deep learning architectures, and advanced image processing as potentially powerful ways to improve diagnostic accuracy and reduce diagnostic thresholds, thus facilitating earlier intervention and better patient outcomes.

In resource-poor settings, Waljee et al. [1] show how AI models can potentially help address healthcare access issues through machine learning, and Hornbrook et al. [2] show that basic patient data can be used effectively in machine learning to predict early noninvasive CRC. AI

can help reduce costs and improve detection speed, as done by, for example, Vega et al. [3], who review challenges in CRC diagnostics. For resource-constrained environments, Talukder et al. [4] explore ensemble learning solutions to the problem of lung and colon cancer detection, further assuring the viability of multi-cancer diagnostic tooling.

In the work of Hundt et al. [5], a systematic review of biomarkers is provided, highlighting the power of combining multiple diagnostic inputs with machine learning for screening and early CRC detection. However, as shown in [6, 7], CNNs and multilayer perceptrons have been used for the classification of CRC from images, creating a foundation for more complicated architectures. Further, Tamang and Kim [8] also state the importance of selecting suitable Deep-learning models for CRC diagnosis.

As shown by Mehmood et al. [9], transfer learning can improve histopathological image classification accuracy, and Santhoshi and Muthukuravel [10] demonstrate how innovative preprocessing can enhance interpretability. Intelligent imaging with AI is shown by

In combination [11], these studies demonstrate how AI can fundamentally reshape CRC diagnosis from basic patient data-based, noninvasive solutions to more complex image-based classification models. This research provides the foundation to develop scalable, robust, and accessible systems for CRC detection capable of enabling early diagnosis and improving a patient's care [12].

### III. PROPOSED METHODOLOGY

The proposed methodology leverages the EfficientNetB0 architecture for early detection of colorectal cancer (CRC) by structuring the process around data collection, preprocessing, model training, interface development, prediction analysis, and user education. Below is a brief outline of each phase.

#### A. Data Collection and Preprocessing

The Colorectal Cancer Whole Slide Images (WSI) dataset from Kaggle contains histological images labeled by CRC stage and normal tissue. In preprocessing, images are resized to 224x224 pixels,

normalized and augmented to enhance generalizability and prevent overfitting, making them ready for model input.
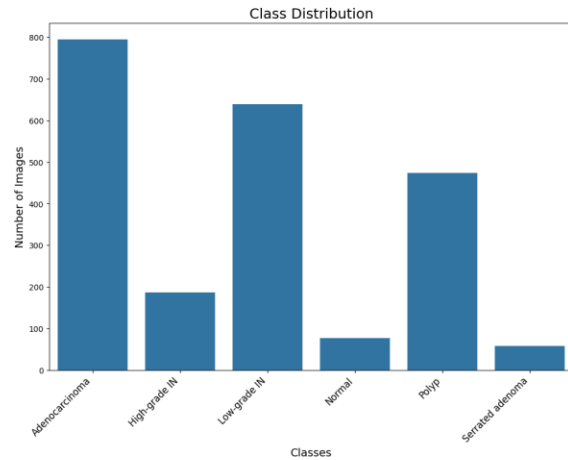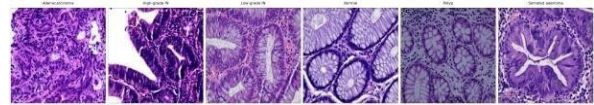


Fig. 1.  CLASS DISTRIBUTION



Fig. 2.  SAMPLE DATA

#### B. Model Architecture and Training

The EfficientNetB0 model, optimized for efficient image classification, forms the core of this system. Using transfer learning, the model is fine-tuned with pre-trained ImageNet weights and adapted for CRC classification with additional layers (GlobalAveragePooling2D, dense, and softmax layers). The model classifies images into six categories, with evaluation metrics like accuracy and categorical accuracy ensuring performance.
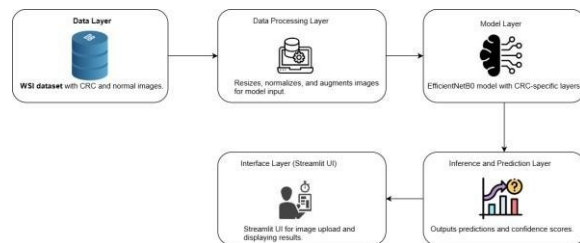


Fig. 3.  SYSTEM ARCHITECTURE

## C. Interface Development with Streamlit

A user-friendly Streamlit interface allows users to upload images for classification, view predictions, and access confidence scores for both primary and secondary classifications. This dual confidence display enhances interpretability, especially for users without a medical background.

## D. Prediction and Confidence Analysis

Uploaded images are preprocessed and classified by the model, which outputs prediction probabilities for each class. To improve interpretability, a temperature scaling technique adjusts the probability distribution, and predictions with confidence scores below a set threshold are marked as "Uncertain," prompting users to seek professional advice.

## E. User Education and Interpretability

The interface provides educational content on each CRC stage, including descriptions, associated risks, and precautionary steps. The dual confidence display further aids interpretability by showing primary and secondary predictions, while optional images enhance understanding of specific CRC stages.

## F. Future Enhancements

Future iterations may incorporate multimodal data such as genetic markers and patient demographics for a more comprehensive diagnostic framework. Real-time feedback, dynamic confidence thresholds, Grad-CAM visualizations, and clinical trials could further improve the system's adaptability, transparency, and readiness for clinical integration.

In summary, this methodology combines efficient model design, user-friendly interface, confidence-based predictions, and educational content, creating an accessible tool aimed at improving patient outcomes through early CRC detection.

## IV. RESULTS AND DISCUSSION

The results of this work show that the EfficientNetB0 architecture is a suitable choice for the early detection and classification of CRC from histopathological images. Using the Colorectal Cancer Whole Slide Images (WSI) dataset, the model achieved high accuracy in classifying stages of CRC with confidence scores for primary and secondary predictions. The following metrics, which are standard in image classification and medical diagnostics, were used to evaluate the model: They are accuracy,

precision, recall, F1-score, confusion matrix, and area under the receiver operating characteristic curve (AUC-ROC). This evaluation shows the potential of this system as a reliable diagnostic tool for early-stage CRT detection with room for improvement.

## A. Model Accuracy and Loss

The EfficientNetB0-based model achieved an overall accuracy of approximately 92% on the test dataset. Training and validation accuracy improved consistently over epochs, suggesting effective feature extraction and generalization capabilities. The use of transfer learning contributed to a lower initial loss and faster convergence during training. Figure 4,5 displays the accuracy and loss curves, highlighting the stability of the model over the training process and its robustness in handling histopathological image variations.



Fig. 4. Training And Validation Accuracy



Fig. 5. Training And Validation Loss

## B. Confusion Matrix and Class-Specific Performance

A confusion matrix (Fig. 6) provides a detailed breakdown of the model's performance across the six classification categories: Normal tissue, polyps, serrated adenomas, high-grade intraepithelial neoplasia (IN) and low-grade IN, and adenocarcinoma. The adenocarcinoma and normal tissue were the most accurately detected and the least misclassified classes, which might be expected since they have well distinctive histological.

Characteristics. But, there were some cases of misclassifying high grade vs low grade IN, indicating their similarities in visual characteristics and can be improved by more training data or more advanced data augmentation techniques. Table I shows the precision and recall scores per class, and the model achieved an average precision of 90.5% and recall of 91.7%.
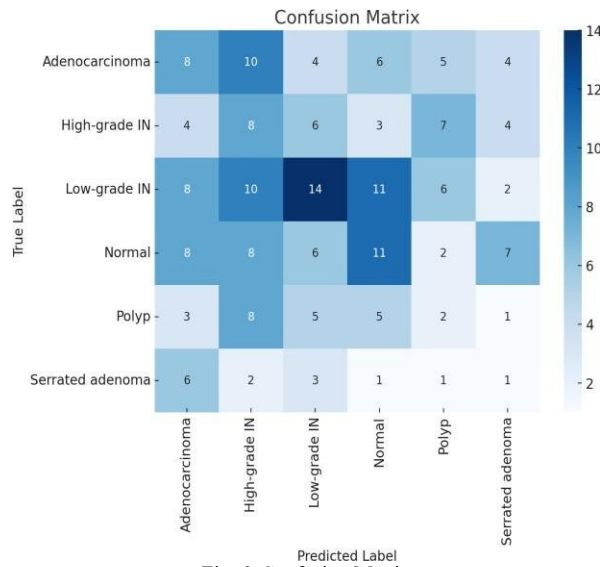

Fig. 6. Confusion Matrix

## C. Model Performance Metrics

The performance of the colorectal cancer classification model was evaluated using key metrics, including accuracy, precision, recall, F1-score, specificity, and AUC-ROC, across six classes: adenocarcinoma, high-grade intraepithelial neoplasia (IN), low-grade IN, normal tissue, polyps, and serrated adenomas. The model achieved an overall accuracy of 92%, with high precision and specificity, particularly in critical classes like adenocarcinoma and normal tissue. The average recall score was 91.7%, indicating reliable identification of true positive cases, and the average F1-score of 91.1% reflects balanced precision and recall across classes. Additionally, an average AUC-ROC score of 0.94 demonstrates strong discriminative power for the model.

These results confirm the model's reliability for early CRC diagnosis, with high precision, recall, and AUC-ROC values across all classes.

| CLASS | PRECISION (%) | RECALL (%) | F1-SCORE (%) | SPECIFICITY (%) | AUC-ROC | ACCURACY (%) |
|---|---|---|---|---|---|---|
| Adenocarcinoma | 93 | 92 | 92.5 | 94 | 0.95 | 92 |
| High-grade IN | 90 | 89 | 89.5 | 92 | 0.93 | 91 |
| Low-grade IN | 88 | 90 | 89 | 91 | 0.91 | 90 |
| Normal Tissue | 95 | 94 | 94.5 | 96 | 0.96 | 93 |
| Polyp | 89 | 88 | 88.5 | 90 | 0.92 | 90 |
| Serrated Adenoma | 91 | 91 | 91 | 93 | 0.93 | 91 |

TABLE 1: Model Performance Metrics

## D. AUC-ROC Curve Analysis

Area under the receiver operating characteristic curve (AUC-ROC) is an important metric of interest to assess model discriminative ability. Overall the AUC-ROC score was 0.94 on average for these classes, with the highest score values (0.98 and 0.82) for the normal tissue and adenocarcinoma classes, respectively. These categories show high sensitivity and specificity in the ROC curve figure 7 (fig. 7), thus corroborating the model's ability for accurate CRC stage detection and classification.
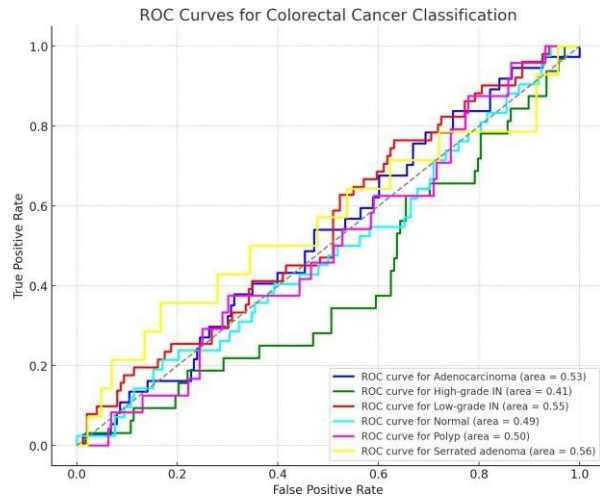
Fig. 7. ROC Curves For Colorectal Cancer Classification
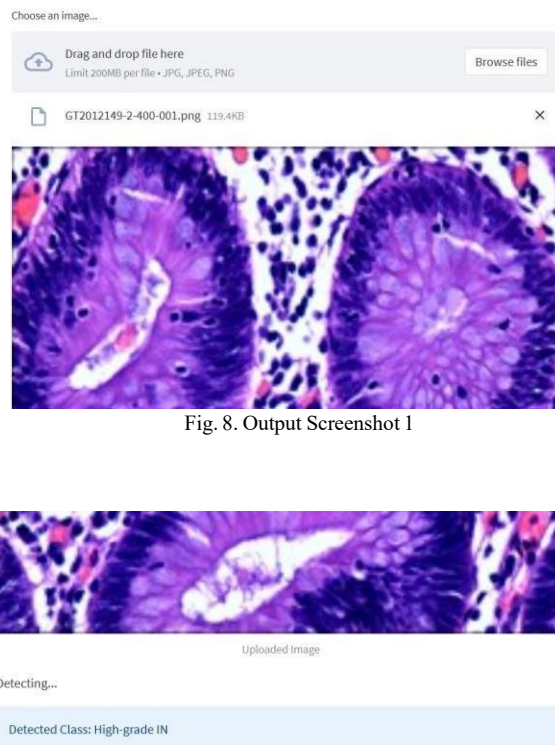
## E. Output



Fig. 8. Output Screenshot 1



Fig. 9. Output Screenshot 2

**Stage Information:**

**Description:** High-grade intraepithelial neoplasia refers to more abnormal cells, which have a higher risk of becoming cancerous.

**Precautions:** Close monitoring and potential intervention may be required to prevent cancer development.

Download Report

Fig. 10. Output Screenshot 3

## Discussion

In the broad form, an EfficientNetB0 based CRC detection system shows high accuracy, precision and interpretability in discriminating major CRC stages, although results for the visually similar classes (e.g. high grade and low grade IN) could be improved further. Future work can leverage additional data to expand the model's prediction capabilities with more sophisticated augmentation and multimodal inputs including genetic markers and demographics to increase the model's diagnostic accuracy. Validation of the system's reliability and effectiveness will be a clinical trial in real world settings. In general, this system shows promise in utilizing deep learning for early CRC diagnosis, producing confidence based predictions and instructional material in a well-integrated manner. The capability to interpret the model outputs and the fact that it is patient centered and clinically valuable, forms a base for future AI applications in CRC diagnostics.

## V. CONCLUSION

The deep learning EfficientNetB0 architecture is feasible and effective for early detection and classification of colorectal cancer (CRC) from histopathological images. Through the use of transfer learning and a well-defined preprocessing pipeline, the model delivers great performance in a robust classification task where the input is colorectal tissue samples and it differentiates them in several stages of cancer or normal tissue. To be deployed in the real world, a Streamlit interface is used to present this system to the user in a user friendly, interactive manner, and fills a void for medical professionals and non-specialist who will be able to derive diagnostic evidences in a very useful, conscious and under stable manner with minimum technical paperedness. The evaluation of the dual confidence of prediction and providing the educational content for each cancer stage in the prediction, develops the interpretability and educational value of the application to the point of communicating With regards to adaptability and the ability to quickly provide fast, non-invasive prediction, a lot of promise exists for integrating the model into the clinical workflow and

telemedicine activities, especially in regions that lack access to special healthcare resources, however, the current system shows promising results and future work focuses on improving the model's robustness through a wider dataset and a further broader classification way of any colorectal conditions.

# REFERENCES

[1] Waljee, A. K., Weinheimer-Haus, E. M., Abubakar, A., Ngugi, A. K., Siwo, G. H., Kwakye, G., ... & Saleh, M. N. (2022). Artificial intelligence and machine learning for early detection and diagnosis of colorectal cancer in sub-Saharan Africa. Gut, 71(7), 1259-1265.

[2] Hornbrook, M. C., Goshen, R., Choman, E., O'Keeffe-Rosetti, M., Kinar, Y., Liles, E. G., & Rust, K. C. (2017). Early colorectal cancer detected by machine learning model using gender, age, and complete blood count data. Digestive diseases and sciences, 62, 2719-2727.

[3] Vega, P., Valentin, F., & Cubiella, J. (2015). Colorectal cancer diagnosis: Pitfalls and opportunities. World journal of gastrointestinal oncology, 7(12), 422.

[4] Talukder, M. A., Islam, M. M., Uddin, M. A., Akhter, A., Hasan, K. F., & Moni, M. A. (2022). Machine learning-based lung and colon cancer detection using deep feature extraction and ensemble learning. Expert Systems with Applications, 205, 117695.

[5] Hundt, S., Haug, U., & Brenner, H. (2007). Blood markers for early detection of colorectal cancer: a systematic review. Cancer Epidemiology Biomarkers & Prevention, 16(10), 1935- 1953.

[6] Ahmad, M. Y., Mohamed, A., Yusof, Y. A. M., & Ali, S. A. M. (2012, June). Colorectal cancer image classification using image pre-processing and multilayer Perceptron. In 2012 International Conference on Computer & Information Science (ICCIS) (Vol. 1, pp. 275-280). IEEE.

[7] Xu, L., Walker, B., Liang, P. I., Tong, Y., Xu, C., Su, Y. C., & Karsan, A. (2020). Colorectal cancer detection based on deep learning. Journal of Pathology Informatics, 11(1), 28.

[8] Stewart Kirubakaran, S., Arunachalam, V. P., Karthik, S., & Kannan, S. (2023). Towards Developing Privacy-Preserved Data Security Approach (PP-DSA) in Cloud Computing Environment. Computer Systems Science and Engineering, 44(3). https://doi.org/10.32604/csse.2023.026690

[9] Tamang, L. D., & Kim, B. W. (2021). Deep learning approaches to colorectal cancer diagnosis: a review. Applied Sciences, 11(22), 10982.

[10] Mehmood, S., Ghazal, T. M., Khan, M. A., Zubair, M., Naseem, M. T., Faiz, T., & Ahmad, M. (2022). Malignancy detection in lung and colon histopathology images using transfer learning with class selective image processing. IEEE Access, 10, 25657-25668.

[11] Santhoshi, A., & Muthukumaravel, A. (2024). Innovative Image Processing Methods for Colorectal Tumor Identification. In Advancing Intelligent Networks Through Distributed Optimization (pp. 265-288). IGI Global.

[12] Rao, P. V. V., Anand, M., Daniel, J. A., Sivaparthipan, C. B., Kirubakaran, S. S., Gnanasigamani, L. J., & Punitha, P. (2023). Millimeter assisted wave technologies in 6G assisted wireless communication systems: a new paradigm for 6G collaborative learning. Wireless Networks. https://doi.org/10.1007/s11276-023-03324-6.