# Deep Learning-Based Colon Cancer Tumor Prediction Using Histopathological Images

1st Rahul Deb Mohalder
*Computer Science and Engineering Discipline*
*Khulna University*
Khulna, Bangladesh
rahul@ictcell.ku.ac.bd

2nd Ferdous Bin Ali
*Statistics Department*
*Jahangirnagar University*
Dhaka, Bangladesh
hridoyferdous@yahoo.com

3rd Laboni Paul
*Computer Science and Engineering Discipline*
*Khulna University*
Khulna, Bangladesh
laboni1124@cseku.ac.bd

4th Kamrul Hasan Talukder
*Computer Science and Engineering Discipline*
*Khulna University*
Khulna, Bangladesh
khtalukder@cse.ku.ac.bd

*Abstract*—Colorectal cancer is one of the deadliest diseases and one of the most difficult diseases to diagnose. A big reason for this is that it takes a long time to identify at an early stage. For treatment, a rapid and precise diagnosis of nodules is very crucial. In order to identify cancer in its early stages, a variety of techniques have been employed. Deep learning approaches were used in this work in order to identify Colorectal cancer tumors. In our this research, we used a dataset of same dimension of colon cancer tissues histopathological images. We proposed a deep learning model for predicting CRC tumors from histopathological images. CNN technique used for analyzing complex data. By CNN technique we analyzed our complex tumor images for identifying abnormal or suspicious tumor patterns. We made a five-layer deep neural network model. It consists of the input layer, four hidden layers, and the output layer. We used Rectified linear unit (ReLU) activation function in the hidden layer and the Softmax function in the output layer. We obtained an accuracy 99.70% from our deep learning model and our model loss was 0.0160. We calculate precision, recall, and F-score for the performance evaluation of our method. It is evident from our experiment that our proposed model produces a better result than some related works.

*Index Terms*—deep learning, colon cancer, histopathological image, medical image analysis.

## I. INTRODUCTION

Cancer is one of the most common causes of death and a major barrier to improving life expectancy across countries. According to WHO, colon cancer is the second most common cancer and third leading cause of cancer death [1]. In all, a rapid increase in the global burden of incidence and killing of cancer is reflecting the aging and growth of both population and changes, many of these are related to socioeconomic development in the prevalence and distribution of primary cancer risk factors.

Colorectal cancer is a major public health concern, both in terms of the number of individuals afflicted and the financial expenses connected with it. Colonoscopy is an essential screening procedure that improves the survival rate of colorectal cancer patients [2]. Researchers are presently focusing on the relationship between colonoscopy and computer-aided diagnostic tools since some approaches have already been presented and show tremendous promise for better disease management.

Colonoscopy is a popular procedure for detecting colorectal cancer early. According to the WHO, a colonoscopy test should be performed every three years for patients with colon cancer. During a colonoscopy examination, a long, flexible tube (colonoscope) is placed in the rectum. The doctor can scan a polyp and see the interior of the collar on the screen using a small camera at the end of the colonoscope. If no abnormalities are identified, the polyp is surgically removed. On the other side, if an anomaly is discovered, the next step is to obtain biopsy tissue samples [3]. This manual anomaly detection approach is not optimal since it takes time and is vulnerable to observer prejudices. It also relies solely on the competency or experience of medical practitioners and technologists [4].

Researchers are trying to bridge the gap between expert and manual decisions. As a result, the effective removal of pre-cancerous polyps using automated computer-assisted procedures may give an advantage in lowering colorectal cancer death rates. Rapid technology advancements have to lead to a variety of low-cost, fast diagnostic computational-based approaches in the field of image processing and machine learning. Pattern recognition-based algorithms are used in traditional approaches for rapid and automated cancer detection. This entails collecting a predefined collection of hand-crafted features from histology pictures based on, for example, texture and morphological qualities, then training a classifier to classify/detect malignant cells using these features [5].

In this work, we proposed a CNN model to predict colon cancer tumor from histopathological data. Our main objective of this research was to identify and predict tumors perfectly. We organized the rest of this paper as Section II about previous work, Section III about proposed methodology, Section IV about result analysis and visualization, Section V about

comparisons, and Section VI about conclusion and future work of this work.

## II. LITERATURE REVIEW

To identify CRC tumors from histopathological data different types of detection or prediction and feature extraction techniques are used. But it was not possible to completely overcome all the limitations of tumor identification.. It has surprisingly emerged from an area completely different from medicine and healthcare that a viable resolution to this dilemma. In the previous half-century, computer science has arguably advanced the greatest compared to other scientific and technology fields. A broad variety of applications in Machine Learning (ML) vary from the identification of sickness to smart systems, which may prescribe traditional medications in accordance with patients' symptoms [6].

Computational procedures to support pathologists in the interpretation of microscopic images have been developed in recent years. These techniques of image analysis mostly focus on fundamental segmentation (e.g. nucleosections) and extraction of functions (e.g., orientation, shape, and texture). These extracted or hand constructed functions can be utilized in some techniques, such as a vectors for support or an automatic forester for tissue categorization and illness grading, as an input to conventional machine classification frameworks.

Likewise, Mehedi et al. have created a new classification approach based on a revolutionary neural network (CNN) to discriminate between five distinct types of lung and colon tissues using a number of histological pictures. The results demonstrate that the model can classify with high confidence the related lung and colon cancer species. The next parts described this cancer diagnostic technique by providing necessary diagrams, tables, charts and other visuals for simple understanding [7].

Many recent research have proposed the machine to categorize, locate and segment tumor regions in histology pictures using machine learning approaches [8]. Deep neural networking (DNNs), used extensively in the extraction and learning of topics, may successfully categorize or identify tumours; nevertheless, very few research utilize CRCs that were particularly tailored to image data termed convolutional neural networking (CNNs) [9], [10]. Łukasz et al. [11] established a precise, reliable and active (ARA) learning system for classifying histological colorectal cancer images. In this context, they have developed a novel CNN model (named ARA-CNN), which is built on the Kather et al. [12] dataset to classify colorectal cancer tissue.

## III. METHODOLOGY

In the Fig. 1 we illustrated our proposed workflow diagram. For our this work we collected a dataset of colon cancers images from Borkowski et al.

### A. Dataset Description

We have used the LC25000 dataset created and curated by Borkowski et al. [13]. This dataset originally contains histopathological images with 5 classes which are benign lung tissue, lung adenocarcinomas, squamous cell carcinomas, benign colon tissue and colon adenocarcinomas. Each class have 5000 images of $768 \times 768$ pixels. We have followed the following notation throughout the paper to represent the classes.

1) Colon_aca = Colon Adenocarcinoma
2) Colon_n = Benign Colonic Tissue

### B. Exploratory Data Analysis

To better understand the low level features of the images, we have done histogram oriented color channel intensity analysis as part of our experimental data analysis process. It's visible from the analysis that the target classes have significant differences in image intensity so a good feature extractor will be able to perform more accurately.

### C. Image Pre-Processing

There are 3 steps in image pre-processing. These are

*1) Image Resize:* Given input images have size are RGB with $768 \times 768$ pixels. This large size is very computationally expensive to handle. So we have reduced the image size to $224 \times 224$ pixels.

*2) Image Normalization:* Image normalization is an image processing technique to change the intensity level of pixels to ensure better contrast. As these input images are RGB, the pixel values of each channel can vary from 0 to 256 where each of the numbers represent a distinct color code. To use this image for training a Deep Neural Network, these computations of high numeric values may lead to complexity. So we normalize those images values in between 0 to 1.

$$
\begin{aligned}
R' &= R/255 \\
G' &= G/255 \\
B' &= B/255
\end{aligned} \tag{1}
$$

In Equation 1, R', G' and B' represent the normalized values for each pixel of each color channel.

*3) Color Channel Conversion:* As noted here, the OpenCV library in python by default reads an image in BGR format. Therefore, it needs to be converted into an RGB coloring format for our further operation. The RGB to HSV conversions by the following mathematical options note in equation 2, 3, 4, 5 [14].

$$
\begin{aligned}
C_{max} &= \max(R', G', B') \\
C_{min} &= \min(R', G', B') \\
\triangle &= C_{max} - C_{min}
\end{aligned} \tag{2}
$$

Hue Calculation:

$$
H = \begin{cases}
0° & , \triangle = 0 \\
60° \times (\frac{G'-B'}{\triangle} \, mod \, 6) & , C_{max} = R' \\
60° \times (\frac{B'-R'}{\triangle} + 2) & , C_{max} = G' \\
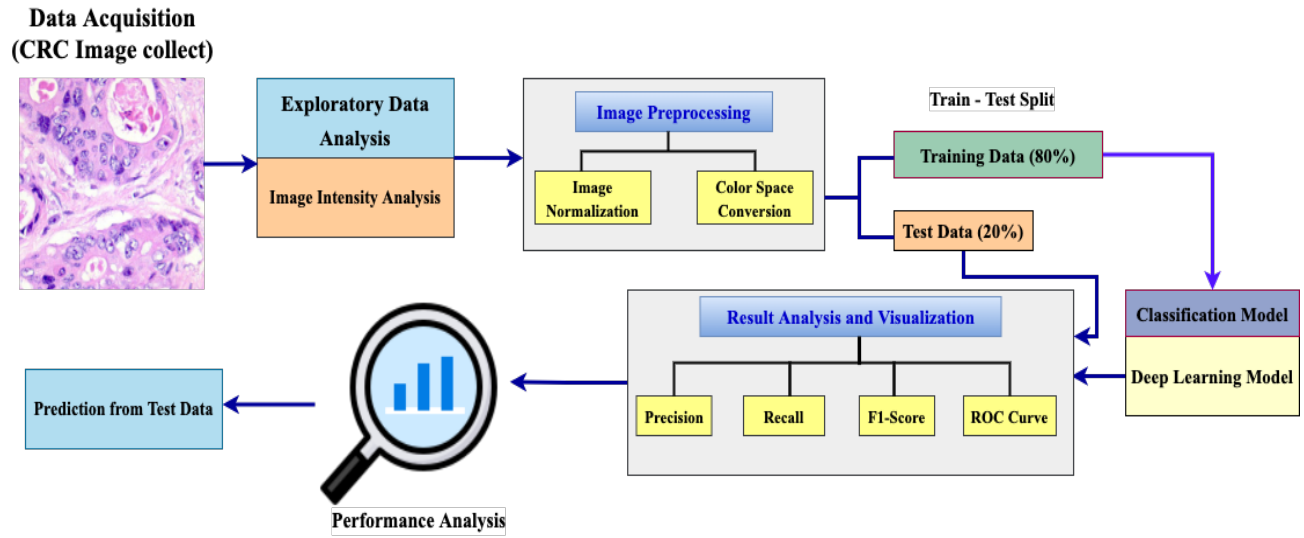60° \times (\frac{R'-G'}{\triangle} + 4) & , C_{max} = B'
\end{cases} \tag{3}
$$
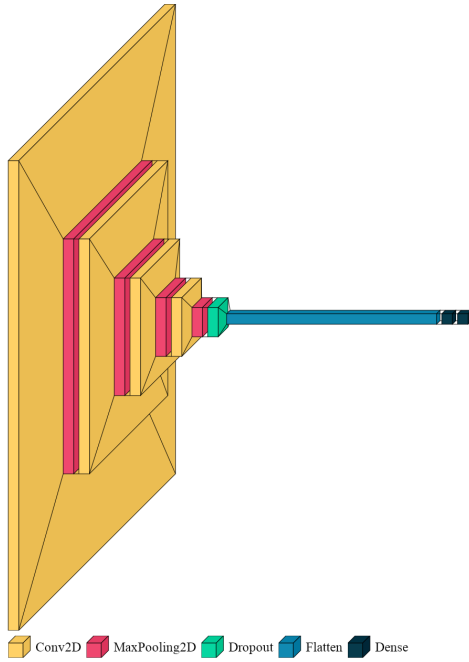
Page 630

Fig. 1: Workflow Diagram of our Proposed System



Fig. 2: Proposed Deep Learning Model

Saturation Calculation:

$$S = \begin{cases} 0 & , \mathbf{C}_{max} = 0 \\ \frac{\triangle}{C_{max}} & , \mathbf{C}_{max} \neq 0 \end{cases} \quad (4)$$

Value Calculation:

$$V = C_{max} \quad (5)$$

### D. Train-Test Split

We have splitted our dataset into train and test data randomly at 80:20 ratio. Our proposed model is trained on the 80% of data and evaluated on the remaining 20% of the data.

### E. Model Description

We have developed a Deep Convolutional Neural Network to address our classification problem. Convolutional Neural Networks are neural networks having convolution filters with learnable filter weights. A convolutional layer i can be represented as a function $Y_i = F_i(X_i)$, where $X_i$ is the input tensor, $Y_i$ is the output tensor and $F_i$ is the convolution operator. It is also to be noted here that input tensor $X_i$ has a tensor shape $(H_i, W_i, C_i)$ where $H_i$ and $W_i$ are considered as spatial dimensions and $C_i$ is the channel dimension.

There are 4 convolution blocks in our proposed network. In every block, there is a convolution filter layer and Max Pooling layer with filter size $2 \times 2$ whereas there is a dropout operation in the fourth block. First and second convolution block consists of 32 $3 \times 3$ filters with the same padding and relu as the activation function. The purpose of using Rectified Linear Unit(ReLU) as an activation function due to its sparse activation nature, computational efficiency and ability to avoid vanishing gradient problems. Mathematically ReLU is defined as in equation 6. Third convolution block consists of 1 convolution layer with 64 3x3 convolution filters keeping other parameters the same followed by another MaxPooling layer with filter size $2 \times 2$. In the last convolution block, the convolution layer has 128 $3 \times 3$ convolution with MaxPooling layer with filter size $2 \times 2$. These four convolution blocks act as a feature extractor. The output of the final block is converted into a single long linear vector of length 25088 and this flattened matrix is fed into the fully connected layer to classify images. There are 2 layers in the fully connected layers, a hidden layer with 128 nodes and output layer with 2 nodes having a softmax (equation 6) as the activation function as it converts the values to a probability distribution [15].
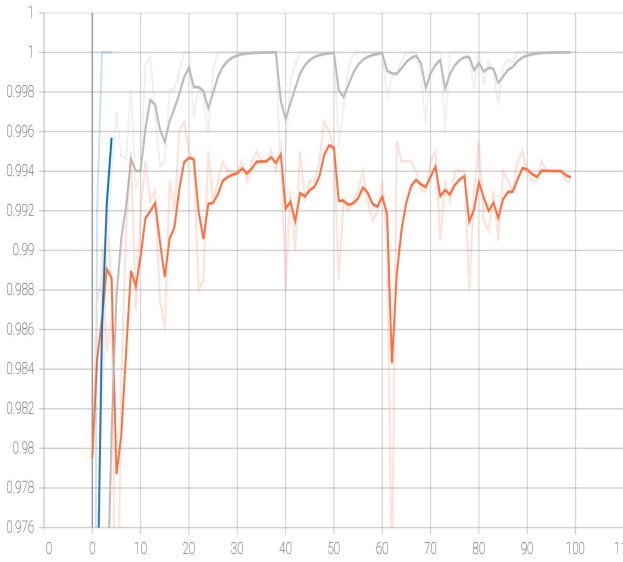
Page 631

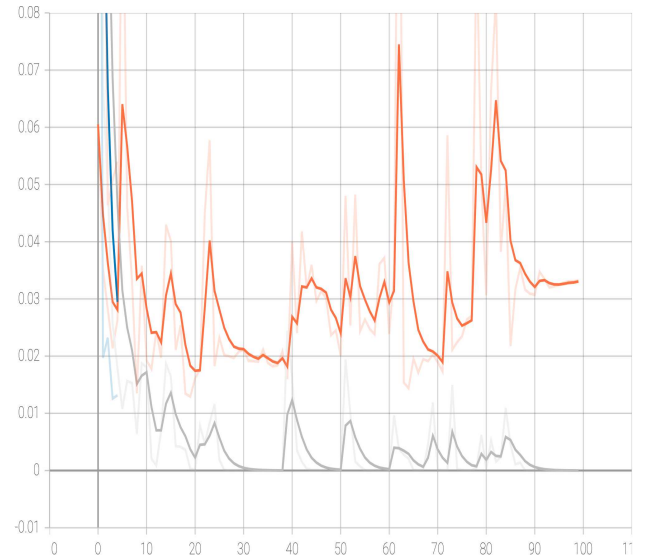Fig. 3: Training and Test accuracy per epoch


Fig. 4: Training and Test loss per epoch

$$f(x) = max(0, x)$$

$$ReLU(x) = \begin{cases} 0 & , \text{if } x < 0 \\ 1 & , \text{if } x > 0 \end{cases} \qquad (6)$$

$$Softmax : \sigma(x_j) = \frac{e_j^x}{\sum_i e_i^x}$$

### F. Training Procedure

We have carefully designed our training procedure in porter to avoid overfitting. The dataset was divided into train and test data randomly at 80:20 ratio. 80% of the training dataset is used to train the model with k-fold cross validation strategy where the value of k=10.

We have used Binary Cross Entropy (equation 7) as a loss function [16]. Adam [17] was used as the optimizer with learning rate 0.001, exponential decay rate for the 1st moment was 0.9 and 2nd moment is 0.999. Backpropagation algorithm was used to train the whole network [18] with batch size 32.

$$BCE = -\frac{1}{n}\sum_{i=1}^{n}(y_i.\log(\hat{y_i}) + (1 - y_i).\log(1 - \hat{y_i})) \quad (7)$$

### G. Experiment Setup

The experiment was implemented using Python programming language of version 3.7.5, Tensorflow version 2.8.1 and Tensorboard for logging and monitoring model training. All experiments were done in Intel Core i7-7700 with 16GB of RAM and NVIDIA GTX 1060 6GB gpu.

### H. Analysis and Visualization

The defined test set that contains 20% of the data is used for testing our model's performance. For analyzing result and measuring accuracy or performance of our model we used Precision, Recall and F1-score technique.
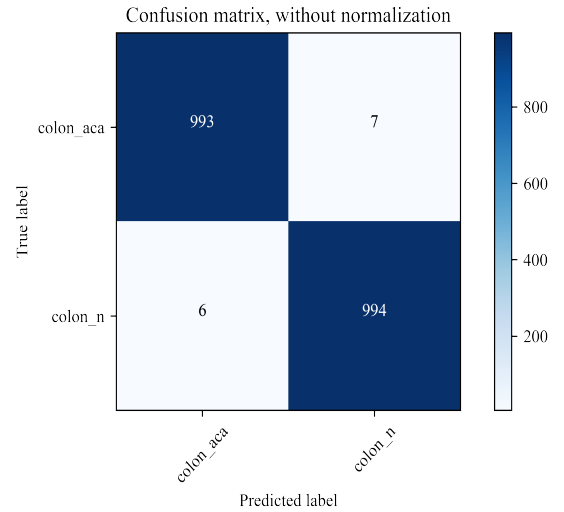

Confusion matrix, without normalization

Fig. 5: Confusion Matrix (CM)

TABLE I: Performance Measures of our DL Model

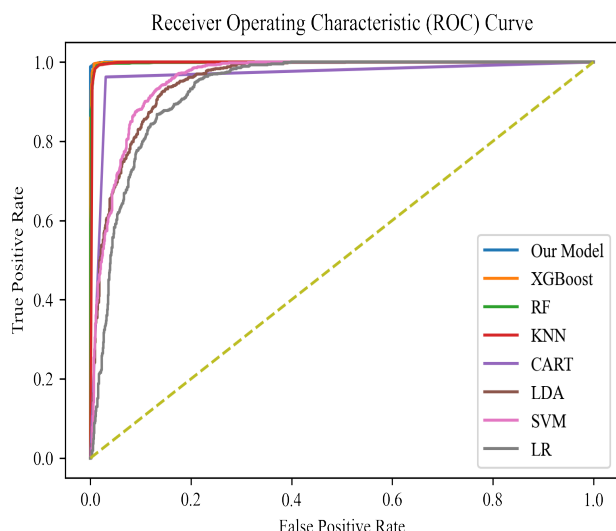|  | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0 (colon_aca) | 0.99 | 0.99 | 0.99 | 1000 |
| 1 (colon_n) | 0.99 | 0.99 | 0.99 | 1000 |
|  |  |  |  |  |
| accuracy |  |  | 0.99 | 2000 |
| macro avg | 0.99 | 0.99 | 0.99 | 2000 |
| weighted avg | 0.99 | 0.99 | 0.99 | 2000 |

Fig. 6: ROC Curve

TABLE II: Accuracy and Loss Comparisons with other Algorithms or Classifiers

| Model, Algorithms or Classifiers | Accuracy | Loss |
|---|---|---|
| XGBoost | 99.06% | 0.007731 |
| RF | 98.54% | 0.009237 |
| KNN | 98.54% | 0.010504 |
| CART | 96.23% | 0.015153 |
| SVM | 91.13% | 0.023542 |
| LDA | 89.61% | 0.022392 |
| LR | 88.56% | 0.026317 |
| Our Model | 99.23% | 0.004361 |

## IV. RESULT ANALYSIS AND VISUALIZATION

For analyzing each epoch's accuracy and loss outcome we used Tensorboard [23] toolkit. In Fig. 3 shows the training and test accuracy per epoch, Fig. 4 training and test loss per epoch.

We compare our results with the five user summaries to measure the performance. We use precision, recall, and f-1score values to determine the performance rate. Those performance outcomes were precision 0.99, recall 0.99, and f-1score 0.99. The TABLE I illustrates our DL model performance measure output. The visualization of the performance of our DL model over test data was done by a confusion matrix shown in Fig. 5 using our proposed deep learning model.

## V. COMPARISONS

### A. Comparison with other Machine Learning Algorithms

We applied the K-Fold (K=10) Cross-Validation process to validate our DL model with other machine learning algorithms. By the receiver operating characteristic (ROC) we measured the classification outcomes which we obtained from our proposed deep learning classifier. Fig. 6 shows the ROC curve. We also compared our proposed model performance with

the ensemble machine learning algorithm, random forest, k-nearest neighbors, decision trees, linear discriminant analysis, support vector machine, logistic regression algorithms. The TABLE II shows the cross validation result. In the Fig. 7 illustrates the machine learning algorithms comparisons.

### B. Comparison with Related Works

We compared our results with other researchers' work. TABLE III shows the comparison with different researchers' used techniques for colon cancer prediction or detection and their proposed models performance with our proposed model.

## VI. CONCLUSION

Classifying Colorectal Cancer Tumors from histopathological images is one of the most challenging and rigorous tasks in the field of biomedical imaging and machine learning. Currently, most of the effective and convenient systems are not automated but need humans in the loop for diagnosis purposes. Our goal is to reduce human involvement to make diagnosis easier. Our proposed deep learning model is able to classify CRC Tumor from images with 99.35% accuracy reaching human-level accuracy. Model performance decreases when image quality decreases are the only limitation of our system.

In the future, we want to achieve human-level accuracy from low-resolution images so that our model is invariant to input image quality. We will also address the explainability and fairness of the model in the future.

## REFERENCES

[1] H. Sung, J. Ferlay, R. L. Siegel, M. Laversanne, I. Soerjomataram, A. Jemal, and F. Bray, "Global cancer statistics 2020: Globocan estimates of incidence and mortality worldwide for 36 cancers in 185 countries," *CA: a cancer journal for clinicians*, vol. 71, no. 3, pp. 209–249, 2021.

[2] J. M. Walsh and J. P. Terdiman, "Colorectal cancer screening: scientific review," *Jama*, vol. 289, no. 10, pp. 1288–1296, 2003.

[3] L. F. Sánchez-Peralta, L. Bote-Curiel, A. Picón, F. M. Sánchez-Margallo, and J. B. Pagador, "Deep learning to find colorectal polyps in colonoscopy: A systematic literature review," *Artificial intelligence in medicine*, p. 101923, 2020.

[4] D. A. Lieberman, D. G. Weiss, J. H. Bond, D. J. Ahnen, H. Garewal, W. V. Harford, D. Provenzale, S. Sontag, T. Schnell, T. E. Durbin *et al.*, "Use of colonoscopy to screen asymptomatic adults for colorectal cancer," *New England Journal of Medicine*, vol. 343, no. 3, pp. 162–168, 2000.

[5] F. Stracci, M. Zorzi, and G. Grazzini, "Colorectal cancer screening: tests, strategies, and perspectives," *Frontiers in public health*, vol. 2, p. 210, 2014.

[6] S. Das, S. Biswas, A. Paul, and A. Dey, "Ai doctor: an intelligent approach for medical diagnosis," in *Industry interactive innovations in science, engineering and technology*. Springer, 2018, pp. 173–183.

[7] M. Masud, N. Sikder, A.-A. Nahid, A. K. Bairagi, and M. A. AlZain, "A machine learning approach to diagnosing lung and colon cancer using a deep learning-based classification framework," *Sensors*, vol. 21, no. 3, p. 748, 2021.

[8] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2010, pp. 249–256.

[9] A. Teramoto, T. Tsukamoto, Y. Kiriyama, and H. Fujita, "Automated classification of lung cancer types from cytological images using deep convolutional neural networks," *BioMed research international*, vol. 2017, 2017.
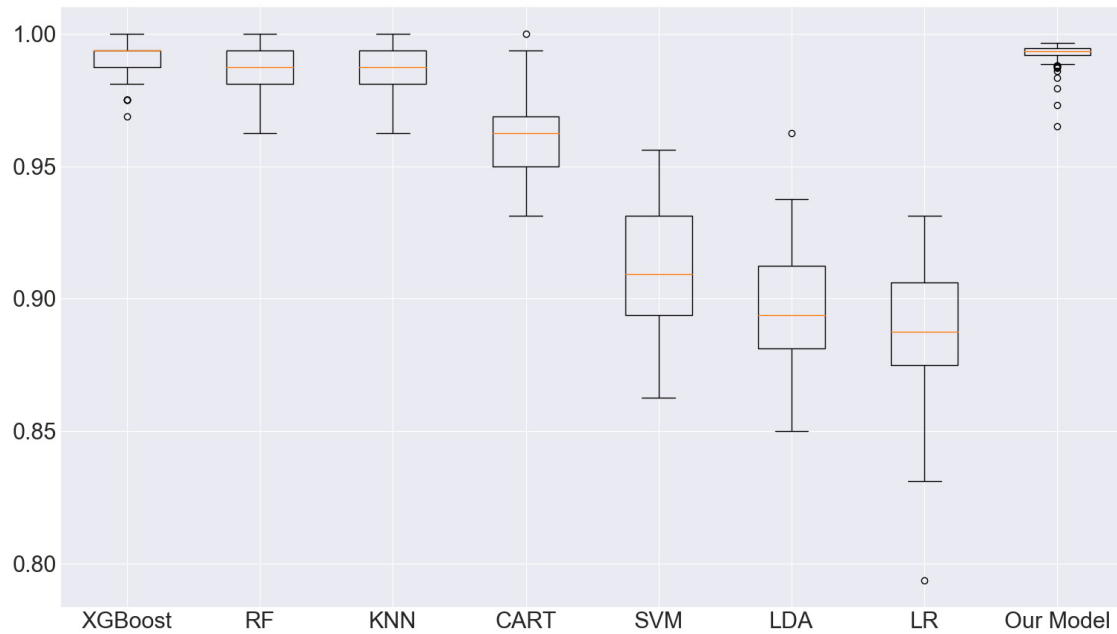
## Machine Learning Algorithm Comparison



Fig. 7: Machine Learning Algorithm Comparison

TABLE III: Comparison with CRC Prediction and Detection Models and Techniques Performance

| References | Models and Techniques | Results |
|---|---|---|
| Masud et al. [7] | Deep Learning-Based Classification | Accuracy-96.33%. |
| Yoon et al. [19] | CNNs, VGG | Accuracy-94.29%, Loss-0.365. |
| Bukhari et al. [20] | ResNet-18, ResNet-30, and ResNet- 50 | Best accuracy 93.91%. |
| Reis et al. [21] | DenseNet169 model | Accuracy-95%. |
| Sakr et al. [22] | Deep CNN model | Accuracy-99.5%. |
| Our Proposed Model | Deep Learning Model | Accuracy-99.7% & Loss-0.016. |

[10] R. Deb Mohalder and K. H. Talukder, "Deep learning based colorectal cancer (crc) tumors prediction," in *2021 12th International Conference on Computing Communication and Networking Technologies (ICCCNT)*, 2021, pp. 01–06.

[11] Ł. Raczkowski, M. Możejko, J. Zambonelli, and E. Szczurek, "Ara: accurate, reliable and active histopathological image classification framework with bayesian deep learning," *Scientific reports*, vol. 9, no. 1, pp. 1–12, 2019.

[12] J. N. Kather, C.-A. Weis, F. Bianconi, S. M. Melchers, L. R. Schad, T. Gaiser, A. Marx, and F. G. Zöllner, "Multi-class texture analysis in colorectal cancer histology," *Scientific reports*, vol. 6, no. 1, pp. 1–11, 2016.

[13] A. A. Borkowski, M. M. Bui, L. B. Thomas, C. P. Wilson, L. A. DeLand, and S. M. Mastorides, "Lung and colon cancer histopathological image dataset (lc25000)," Dec 2019. [Online]. Available: https://arxiv.org/abs/1912.12142

[14] R. C. Gonzalez, *Digital image processing*. Pearson education india, 2009.

[15] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.

[16] F. Topsøe, "Bounds for entropy and divergence for distributions over a two-element set," *J. Ineq. Pure & Appl. Math*, vol. 2, no. 2, p. 300, 2001.

[17] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[18] Z. Zhang, "Derivation of backpropagation in convolutional neural network (cnn)," *University of Tennessee, Knoxville, TN*, vol. 22, p. 23, 2016.

[19] H. Yoon, J. Lee, J. E. Oh, H. R. Kim, S. Lee, H. J. Chang, and D. K. Sohn, "Tumor identification in colorectal histology images using a convolutional neural network," *Journal of digital imaging*, vol. 32, no. 1, pp. 131–140, 2019.

[20] S. U. K. Bukhari, S. Asmara, S. K. A. Bokhari, S. S. Hussain, S. U. Armaghan, and S. S. H. Shah, "The histological diagnosis of colonic adenocarcinoma by applying partial self supervised learning," *medRxiv*, 2020.

[21] H. C. Reis and V. Turk, "Transfer learning approach and nucleus segmentation with medclnet colon cancer database," *Journal of Digital Imaging*, pp. 1–20, 2022.

[22] A. S. Sakr, N. F. Soliman, M. S. Al-Gaashani, P. Pławiak, A. A. Ateya, and M. Hammad, "An efficient deep learning approach for colon cancer detection," *Applied Sciences*, vol. 12, no. 17, p. 8450, 2022.

[23] TensorBoard, TensorFlow's visualization toolkit. [Online]. Available: https://www.tensorflow.org/tensorboard