

Deep Learning - Based Prediction Models for Colon and Rectum Cancer Disease: Enhancing the Precision Diagnosis

1st Sunil Kumar Suman

Department of Computer Science and Engineering
Supaul College of Engineering
Supaul, India.
Sunilkumarsuman12@gmail.com

2nd Dr. D. Kumaraswamy

Electronics and Communication Engineering
Ellenki College of Engineering and Technology
Sangareddy, Telangana, India.
kumaraswamy.btec@gmail.com

3rd Sunil Kumar Sahu

Department of Computer Science and Engineering
Supaul College of Engineering
Supaul, India.
Sunilsahu847@gmail.com

4th Prasanna Lakshmi Akella

Electronics and Communication Engineering
Koneru Lakshmaiah Education Foundation (Deemed to be University) Hyderabad, India.
prasannalakshmiakella@gmail.com

5th Dr. S. P. Santhoshkumar

Department of Computer Science and Engineering
School of Computing,
Vel Tech Rangarajan Dr.
Sagunthala R&D Institute of Science and Technology
Chennai, India.
0000-0001-8531-759

6th R. Karthika

School of Computing
Department of Information Technology Rathinam Technical Campus Coimbatore, India.
karthika.cse@rathinam.in

Abstract—Colon and rectal cancer are among the leading causes of cancer-related illness and death worldwide. The detection and treatment of colon cancer are viewed as social and economic issues due to the high fatality rates. Every year, around half a million people worldwide develop cancer, including colon cancer. Although improving patient outcomes requires an accurate and timely diagnosis, traditional diagnostic methods usually suffer from subjectivity and variability problems. An extensive summary of colon cancer diagnosis is given here. Deep learning, a branch of artificial intelligence, has become a powerful tool for increasing cancer detection accuracy by utilizing large datasets and complex algorithms. This study combines genomic data, medical imaging, histopathological analysis, and clinical decision-making tools to better comprehend the creation and application of deep learning-based prediction models in the diagnosis of colon and rectal cancer. Through the use of Recurrent Neural Networks (RNNs), Convolutional Neural Networks (CNNs), and hybrid models, deep learning can improve the accuracy of cancer detection, anticipate the course of a disease, and assist with customized treatment planning. By automating and enhancing diagnostic precision, these models have the potential to revolutionize healthcare workflows, reduce diagnostic errors, and ultimately increase patient survival rates. This research also discusses the challenges and possible advantages of using deep learning to clinical practice for colon and rectal cancer. Last but not least, the problems that exist today and possible directions for future research are examined.

Keywords—Deep Learning, Colorectal Cancer, Precision Diagnosis, Colon and Rectum Cancer, prediction Model.

I. INTRODUCTION

CRC, sometimes referred to as colon and rectal cancer, is one of the most common and deadly illnesses in the world. Since CRC is becoming more common, it is becoming a serious public health issue, particularly in wealthy nations. Since early-stage malignancies are more curable and have better prognoses, early diagnosis and prompt treatments are essential for lowering death rates [1]. Though useful, conventional methods of diagnosing colorectal cancer

(CRC), such as colonoscopy, histological analysis, and imaging techniques, are frequently constrained by subjectivity, human error, and the availability of qualified practitioners. Because of these drawbacks, cutting-edge instruments that support early, accurate, and automatic CRC detection must be developed.

One of the organs most frequently impacted by cancer is the colon. Specifically, colon cancer ranks third among cancers that affect both men and women [2]. Despite the availability of sophisticated screening technologies, With an expected 2 million new cases identified each year and over 1 million fatalities from the disease as of 2020, CRC is the second most common cause of cancer-related deaths globally [1], [6]. Unfortunately, the World Health Organization and the International Agency for Research on Cancer estimate that the annual number of deaths associated with colon and rectal cancer will rise by roughly 69% between 2020 and 2040, and the overall burden of the disease will rise by about 56% [7]. Additionally, current research indicates that colorectal cancer (CRC) is becoming more common worldwide, even in sub-Saharan Africa, despite the fact that it was previously believed to affect only those nations that had embraced a lifestyle similar to that of the West [11]. Therefore, the prevention of CRC-related fatalities requires early diagnosis and appropriate treatment. Among the newest instruments for accomplishing this aim is AI, and more specifically, the subtype of AI known as "deep learning (DL)" [12].

AI has advanced significantly in the healthcare sector in current years. One area of AI that has shown particular promise is deep learning, which is used in medical diagnostics. Deep learning models have proven to be exceptionally capable of processing huge and complicated datasets, including genetic data, histopathology slides, and medical pictures. In particular, CNNs and RNNs constitute the foundation of these models. By learning from patterns in these datasets, deep learning algorithms can assist in identifying malignant lesions and predict the progression of

the disease, the chance of a recurrence, and each patient's reaction to a certain treatment plan.

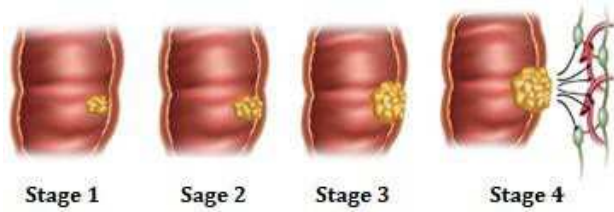


Fig. 1. The many phases of colon cancer

When cancers have not yet developed on the organ wall, the mucosa, or lining, of the colon or rectum is regarded as being in the initial stage. In the second stage, no lymph nodes or adjacent tissues are yet affected, but tumors begin to form on the colon or rectum walls [22]. The tumor is considered to be at the third stage when it has only spread to the lymphatic tissues and has not yet spread to any other part of the body. In the fourth stage, the tumor extends to other organs, including the lungs [23].

II. LITERATURE REVIEW

Dimitrios Bousis et al. (2023), Colon cancer is a serious public health issue that is affecting more and more people globally [1]. A fast and precise diagnosis of colon cancer is the first step toward effective therapy and/or avoiding a future recurrence of the disease. An overview of the full spectrum of deep learning applications in colon cancer diagnostic processes, such as endoscopy, histologic analysis, medical imaging, and screening serologic testing, is the aim of this report.

Kyriaki Katsaounou et al. (2022) highlighted the CRC, one of the most prevalent cancers in people, has a complex etiology that includes both hereditary and environmental factors [2]. Unlike tumors with known environmental, heritable, or sex-linked origins, sporadic colorectal cancer is hard to forecast and lacks clinically meaningful genetic biomarkers of risk. One out of every twenty cases of colorectal cancer has a known heritable component. Partially unknown dietary, behavioral, microbiological, regenerative, genetic, and epigenetic factors are responsible for the remaining cases, which happen seldom. The colonoscopy technique, which is always recommended beyond a certain age, needs to be improved to include an assessment of biomarkers that indicate an individual's risk of colorectal cancer in order to address this problem. Ideally, these indicators would be the source of the disease and could be altered by dietary or therapeutic adjustments. transcriptional analysis among other omics. There is an urgent demand for epigenetic, metagenomic, and metabolomic profiles to provide data for risk evaluations. In order to identify much-needed predictive biomarkers, this article aims to provide an overview of the multifactorial disruption of homeostasis that causes CRC, which can be examined utilizing multi-omics and gut-on-chip analysis.

Hyun-Jong Jang et al. (2020), Since unique mutational patterns can be highly instructive in determining the best course of treatment, it has become more and more crucial to identify genetic alterations in cancer patients. Recent research has demonstrated that typical hematoxylin and eosin (H&E) sections from a variety of tumors, including CRCs, can be used directly for deep learning-based molecular

cancer subtyping. Since H&E-stained tissue slides are widely available, predicting mutations from tumor pathology images can be a simple and inexpensive supplement to customized treatment [25].

Andrew M D Wolf et al. (2018), According to one analysis, colorectal cancer is the fourth most common disease diagnosed in adults and the second leading cause of cancer-related fatalities. These included a fresh evaluation of the age at which screening should begin by race and sex, as well as updated modeling that accounts for variations in the incidence of colorectal cancer in the United States. By identifying and removing adenomatous polyps and other precancerous lesions, screening with any of several modalities is linked to a considerable diminish in the incidence of CRC and a decrease in mortality from early identification and incidence reduction [3].

Kainz P et al. (2017), showcase the quantitative and qualitative segmentation outcomes on the newly made public Warwick-QU colon adenocarcinoma dataset linked to the Glas@MICCAI2015 challenge [24]. The findings demonstrate that deep learning techniques can produce extremely precise and repeatable outcomes for biomedical image analysis, which could greatly enhance the caliber and speed of medical diagnosis.

III. OTHER DIAGNOSTIC TESTS FOR CRC USING APPLICATIONS OF DL

The results of using DL procedures in various tests meant to identify CRC will be presented in this section.

A. Virtual colonoscopy

Virtual colonoscopy is a different and becoming more and more common technique for finding polyps in the intestine lumen. It is especially helpful for people who, for a variety of reasons, including age, medical issues, or contraindications, cannot endure or comply with a traditional colonoscopy. This surgery involves thorough bowel scanning followed by a precise reconstruction of the gut lumen utilizing a specific CT scan protocol. CT colonography is less invasive than traditional colonoscopy because it doesn't involve inserting a scope into the colon. Rather, it offers a comprehensive, three-dimensional picture of the colon and rectum, making it possible to identify any anomalies, such as tumors or polyps, that might point to colorectal cancer [13].

Furthermore, patients typically find CT colonography to be more pleasant and speedier. Its limitations, like those of any diagnostic tool, include its inability to identify flat lesions or smaller polyps, and in the event that a polyp is found, a traditional colonoscopy is still necessary for biopsy or excision. However, CT colonography is becoming a more important tool for colorectal cancer screening, providing a good alternative for those who cannot undergo traditional colonoscopy.

The result is X-ray-derived pictures that are used to detect polyps, although this method is subject to the same patient preparation constraints as traditional colonoscopy.

B. CT Scan

CT scans are not always the initial diagnostic method used to diagnose CRC, as endoscopy and pathologic inspection are. Instead, they are the most reliable method for

figuring out how much disease staging is necessary at the time of early diagnosis [20]. The need to extract as much diagnostic information as possible from a CT scan has led several researchers to look into the potential use of deep learning in CT scan data extraction, as a CT scan is usually performed on a patient with colorectal cancer before any surgical procedures. Wang et al. increased the sensitivity of extra-peritoneal colorectal cancer diagnosis by about 95% in 2022 by employing 3D reconstruction of CT scan images of patients in stages II and III of the disease and correlating serum levels of Carcinoembryonic Antigen using a DL protocol. Future study on the search for more diagnostic information from CT scan images is highly desirable [18].

C. OCT (Optical Coherence Tomography)

Although it has long been used to identify diseases of the retina, recent research suggests that it may also be applied to other tissues and organs by differentiating between layers of normal and cancerous tissue [13]. For instance, CRC is diagnosed using OCT. According to published results from 2020, the CNN exhibited 100% sensitivity and an area under the receiver operating characteristic curve of about 99%, suggesting that it could become a more powerful tool for optical diagnosis of colorectal cancer [19].

IV. EVALUATION METRICS FOR PERFORMANCE

Machine learning models' performance is measured using evaluation measures. The efficacy of the DL algorithm for training on fresh data can be evaluated using these assessment metrics. A variety of assessment indicators can be used to test a model. Using several metrics to evaluate the quality of a trained model can yield more accurate results because each model run using one evaluation metric is different from the same model run using another evaluation metric. The next part presents the formulas and provides an explanation of the evaluation measures employed by academic papers. Accuracy, which can be calculated as follows, measures the proportion of true observations to all samples measured:

$$\text{Precision or Accuracy} = \text{Sum of TN, TP} / \text{Sum of FP, FN, TN, TP}$$

The percentage of incorrect observations relative to the total number of measured samples is displayed by the Rate of Error, which can be computed as follows:

$$\text{Error Rate} = \text{Sum of FN and FP} / \text{Sum of FP, FN, TN and TP}$$

The accuracy, which can be computed as follows, quantifies the true categorized positive estimations of all categorized estimates in a valid category.

$$\text{Accuracy} = \text{TP} / \text{Sum of FP, TP}$$

The ratio of accurate estimations to accurately predicted estimates is measured using the recall. The calculation for this is

$$\text{Recall} = \text{TP} / \text{Sum of FN, TP}$$

The specificity can be computed as follows and is used to measure the positive observations rate of incorrect samples:

$$\text{Specificity} = \text{TN} / \text{Sum of FN and TN}$$

The sensitivity, which may be calculated as follows, is the number of true samples that are considered accurate.

$$\text{Sensitivity} = \text{TN} / \text{Sum of FN and TN}$$

The following method can be used to calculate the ROC curve: $\text{TPR} = \text{TP} / \text{TP} + \text{FN}$ displays the ratio of false positives to TPs by displaying the outcomes of the various threshold values that were employed.

$$\text{FPR} = \text{FP} / \text{Sum of TN and FP}$$

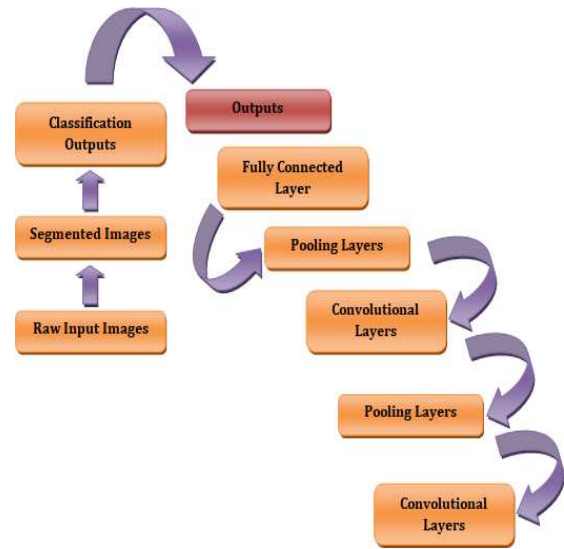


Fig. 1. Deep learning is used in the segmentation and classification process

V. CLASSIFICATION MODEL FOR LUNG CANCER HISTOPATHOLOGY IMAGES

The process of building a CNN for the categorization of lung cancer histology images using Tensor Flow and Keras. The dataset contains images of lung tissue in three separate categories, and train a model to distinguish between them. Preprocessing methods, model construction, training, and evaluation through visualization techniques are some of the salient features.

A. Libraries and Configuration of the Environment

Importing the required libraries first. Here's a quick rundown of each:

- *OS and warnings*: Specifically used to silence TensorFlow messages for a cleaner output, this function manages system warnings and environment setup. Seaborn, matplotlib, and numpy are crucial libraries for data visualization and numerical calculations.
- *Tensorflow*: The primary library for deep learning model construction and training. Performance reports like the confusion matrix and classification report are produced using sklearn.metrics.

B. Preprocessing and Data Loading

Using the `image_dataset_from_directory` function, which generates TensorFlow datasets for training and validation by reading photos straight from the directory, we load the dataset. 90% of the dataset is used for training, while 10% is used for validation. Images are normalized by rescaling pixel values and scaled to (8, 8). To gather the validation data into arrays `x_val` and `y_val` and apply the rescaling transformation to the training and validation datasets.

C. Seeing Examples of Pictures

To have a better understanding of the data, we show a few sample photographs from the training set. The class name that corresponds to each image is labeled.

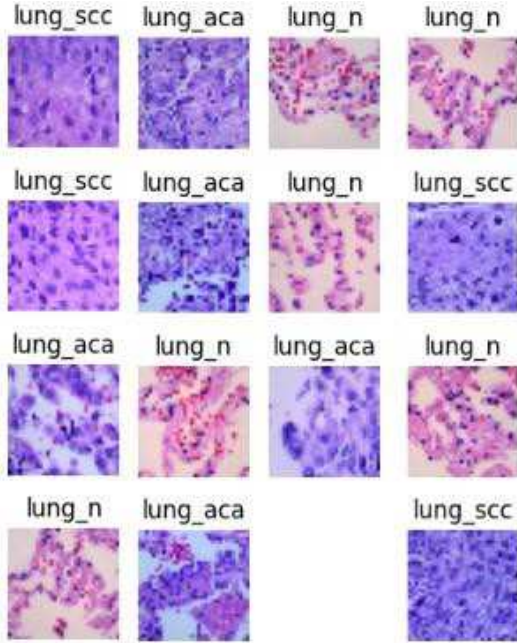


Fig. 2. Sample Images with Corresponding Class

The datasets for effective training and one-hot encode the labels to make them compatible with categorical cross-entropy loss. We also optimize data loading by using the prefetch function.

D. Architecture Model

Use the Sequential API to define a CNN model. This architecture consists of using convolutional layers, features may be extracted from pictures.

- Batch Normalization - To enhance training by normalizing the activations.
- Dropout - To stop overfitting, neurons are sporadically deactivated.
- To further lessen overfitting, L2 Regularization is added to each convolution layer.
- To downsample the feature maps, use MaxPooling. For output including multi-class categorization, use a dense layer with Soft-max Activation.

E. Model Compilation:

To use the Adam optimizer to build the model with a 0.0001 learning rate. Accuracy and AUC for multi-label classification are the metrics monitored, and categorical cross-entropy is the loss function. Three essential callbacks should be used in order to improve training:

- *Early Stopping*: Terminates training after 15 epochs if no progress is observed.
- *Reduce LR On Plateau*: If validation loss reaches a plateau, this feature cuts the learning rate in half.
- *Model Checkpoint*: This feature saves the optimal

Training of Models Utilize the callbacks to control the learning process after 300 epochs of training the model. Assessment and Visualization Loss and Accuracy.

TABLE I. CNN MODEL SEQUENTIAL API- SOFTMAXACTIVATION

| Layer | Output contour | Parameter |
|-----------------------|-----------------|-----------|
| conv2d | (0, 8, 8, 16) | 1,792 |
| batch_normalization | (0, 8, 8, 16) | 64 |
| conv2d_1 | (0, 8, 8, 16) | 36,928 |
| batch_normalization_1 | (0, 8, 8, 16) | 64 |
| max_pooling2d | (0, 16, 16, 16) | 0 |
| dropout | (0, 16, 16, 16) | 0 |
| conv2d_2 | (0, 16, 16, 32) | 73,856 |
| batch_normalization_2 | (0, 16, 16, 32) | 512 |
| conv2d_3 | (0, 16, 16, 32) | 1,47,584 |
| batch_normalization_3 | (0, 16, 16, 32) | 512 |
| max_pooling2d_1 | (0, 8, 8, 32) | 0 |
| dropout_1 | (0, 8, 8, 32) | 0 |
| conv2d_4 | (0, 8, 8, 64) | 2,95,168 |
| batch_normalization_4 | (0, 8, 8, 64) | 1,024 |
| conv2d_5 | (0, 8, 8, 64) | 5,90,080 |
| batch_normalization_5 | (0, 8, 8, 64) | 1,024 |
| max_pooling2d_2 | (0, 4, 4, 64) | 0 |
| dropout_2 | (0, 4, 4, 64) | 0 |
| conv2d_6 | (0, 4, 4, 512) | 11,80,160 |
| batch_normalization_6 | (0, 4, 4, 512) | 2,048 |
| conv2d_7 | (0, 4, 4, 512) | 23,59,808 |
| batch_normalization_7 | (0, 4, 4, 512) | 2,048 |
| max_pooling2d_3 | (0, 2, 2, 512) | 0 |
| dropout_3 (Dropout) | (0, 2, 2, 512) | 0 |
| flatten | (0, 2048) | 0 |
| dense | (0, 3) | 6,147 |

BN – BatchNormalization, MP2D - MaxPooling2D, Total parameter: 4,699,203 (17.93 MB). Trainable parameter: 4,695,363 (17.91 MB), Non-trainable parameter: 3,840 (15.00 KB)

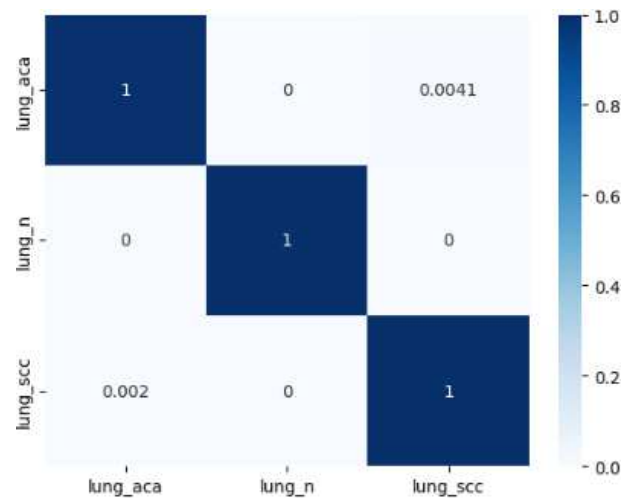


Fig. 3. ClassificationReport and Confusion Matrix

F. Confusion Matrix and Classification Report

After that, the best model is loaded, predictions are formed, and the validation data's classification report and confusion matrix are generated.

TABLE II CLASSIFICATION REPORT AND CONFUSION MATRIX

| | Accuracy | Recall | F1-score | Support |
|------------------|----------|--------|----------|---------|
| 0 | One | One | One | 493 |
| 1 | One | One | One | 515 |
| 2 | One | One | One | 492 |
| | | | | |
| Precision | | | One | 1500 |
| Macro Average | One | One | One | 1500 |
| Weighted Average | One | One | One | 1500 |

VI. CONCLUSION

Determining the cancer's stage is typically the most crucial step after determining its presence, as this knowledge determines the best treatment option and length. Analyzing the tissue region's structure is the primary method of determining the cancer stage or grade. This is done utilizing a variety of screening techniques is used to look for polyps or colorectal cancer in pictures. A thorough argument of the diagnosis procedure was provided in this publication. To illustrate the many imaging modalities used in the analysis process. HIs, which are visible under a microscope, are the most often used technique for colon cancer. The state-of-the-art methods, which are divided into ML and DL procedures, were examined because they aid in identifying cancer in its early stages, which lowers the death rate and allows for early treatment. By removing nonmalignant cells following an early diagnosis of colorectal cancer through screening tests, these methods can also help slow the disease's progression. Also, provided a number of prospective study techniques that would be examined in relation to the scientific endeavors on this research issue. In the future, a set of benchmark datasets and predetermined assessment metrics will be used to compare and examine the most widely used ML and DL algorithms in a single environment.

REFERENCES

- [1] Bousis, D., Verras, G., Bouchagier, K., Antzoulas, A., Panagiotopoulos, I., & Katinioti, A. et al., (2023). The role of deep learning in diagnosing colorectal cancer. *Gastroenterology Review/Przegląd Gastroenterologiczny*, volume 18, issue 3, pp. 266-273. <https://doi.org/10.5114/pg.2023.129494>
- [2] Katsaounou K, Nicolaou E, Vogazianos P, et al., (2022). "Colon cancer: from epidemiology to prevention", *Metabolites*, volume 12: no. 499.
- [3] Wolf AMD, Fonham E, Church T, et al., (2018). Colorectal cancer screening for average-risk adults: 2018 guideline update from the American Cancer Society. *CA Cancer J Clin*, volume 68: pp. 250-81.
- [4] Mitsala A, Tsalikidis C, Pitiakoudis M, et al., (2021) "Artificial intelligence in colorectal cancer screening, diagnosis and treatment", A new era. *CurrOncol*, volume 28, pp. 1581-607.
- [5] Mulita F, Lotfollahzadeh S., (2022). "Intestinal stoma", *StatPearls Publishing, Treasure Island, FL*.
- [6] Hossain MS, Karuniawati H, Jairoun AA, et al., (2022). "Colorectal cancer: a review of carcinogenesis, global epidemiology, current challenges, risk factors, preventive and treatment strategies", *Cancers* volume 14, no. 178.
- [7] "Colorectal Cancer Awareness Month 2022" – IARC. <https://www.iarc.who.int/featured-news/colorectal-cancer-awareness-month-2022/>.
- [8] Wong MCS, Huang J, Huang JLW, et al., (2020), "Global prevalence of colorectal neoplasia: a systematic review and meta-analysis", *Clin Gastro enterol Hepatol*, volume 18, pp. 553-61.
- [9] Waljee AK, Weinheimer-Haus EM, Abudakar A, et al., (2022). "Artificial intelligence and machine learning for early detection and diagnosis of colorectal cancer in sub-Saharan Africa", *Gut*, volume 71, pp. 1259-65.
- [10] Ben Hamida A, Devanne M, Weber J, et al., (2021), "Deep learning for colon cancer histopathological images analysis", *ComputBiol Med* volume 136, pp.104730.
- [11] Ho C, Zhao Z, Chen XF, et al., (2022). "A promising deep learning-assistive algorithm for histopathological screening of colorectal cancer", *Sci Rep.*, volume 12, no. 2222.
- [12] Pacal, I, Karaboga D, Basturk A, et al., (2022). "A comprehensive review of deep learning in colon cancer", *ComputBiol Med.*, volume 126, no. 104003.
- [13] Mulita F, Tepetes K, Verras GI, et al., (2022). "Perineal colostomy: advantages and disadvantages", *Gastroenterology Rev.*, volume 17, pp. 89-95.
- [14] Triantafyllidis JK, Vagianos C, Malgarinos G., (2015), "Colonoscopy in colorectal cancer screening: current aspects", *Indian J SurgOncol*, volume 6, pp. 237-50.
- [15] Wesp P, Grosu S, Graser A, et al., (2022). "Deep learning in CT colonography: differentiating premalignant from benign colorectal polyps", *EurRadiol*, volume 8, pp. 4749-57.
- [16] Hegde N, Shishir M, Shashank S, et al., (2021). "A survey on machine learning and deep learning-based computer-aided methods for detection of polyps in CT colonography", *Curr Med Imaging*, volume 17, pp. 3-15.
- [17] Cao W, Pomeroy MJ, Gao Y, et al., (2019). "Multi-scale characterizations of colon polyps via computed tomographic colonography", *Vis ComputInd Biomed Art*, volume 2, no. 25.
- [18] Wang X, Gu C, Zha Y, et al., (2022), "Diagnosis of nonperitonealized colorectal cancer with computerized tomography image features under deep learning", *Contrast Media Mol Imaging*, no.1881606.
- [19] Zeng Y, Xu S, Chapman Jr WC, et al., (2020). "Real-time colorectal cancer diagnosis using PR-OCT with deep learning", *Theranostics* volume 10, pp. 2587-96.
- [20] Dighe S, Swift I, Brown G., (2008). "CT staging of colon cancer", *ClinRadiol* volume 63, pp. 1372-9.
- [21] Tharwat M, Sakr NA, El-Sappagh S, Soliman H, Kwak K-S., (2022). "Elmogy M. Colon Cancer Diagnosis Based on Machine Learning and Deep Learning: Modalities and Analysis Techniques", *Sensors*, volume 22, issue 23, no.9250. <https://doi.org/10.3390/s22239250>
- [22] Baxter, N.N.; Goldwasser, M.A.; Paszat, L.F.; Saskin, R.; Urbach, D.R.; Rabeneck, L., (2009). "Association of colonoscopy and death from colorectal cancer", *Ann. Intern. Med.*, 150, pp 1–8.
- [23] Bera, K.; Schalper, K.A.; Rimm, D.L.; Velcheti, V.; Madabhushi, A., (2019), "Artificial intelligence in digital pathology—New tools for diagnosis and precision oncology", *Nat. Rev. Clin. Oncol.*, volume 16, pp 703–715.
- [24] Kainz, P.; Pfeiffer, M.; Urschler, M., (2017). "Segmentation and classification of colon glands with deep convolutional neural networks and total variation regularization", *PeerJ.*, volume 5, e3874.
- [25] Amelie Echle et al., (2021). "Deep learning for the detection of microsatellite instability from histology images in colorectal cancer: A systematic literature review", *Immunoinformatics*, volumes 3–4, no. 100008 <https://doi.org/10.1016/j.immuno.2021.100008>