

A Deep Convolutional Neural Network for the Detection of Polyps in Colonoscopy Images.

Tariq Rahim, Syed Ali Hassan, Soo Young Shin*

Kumoh National Institute of Technology, Gumi, Gyeongbuk 39177, Republic of Korea

tariqrahim@ieee.org, Syedali@kumoh.ac.kr, wdragon@kumoh.ac.kr

Abstract

Computerized detection of colonic polyps remains an unsolved issue because of the wide variation in the appearance, texture, color, size, and presence of the multiple polyp-like imitators during colonoscopy. In this paper, we propose a deep convolutional neural network based model for the computerized detection of polyps within colonoscopy images. The proposed model comprises 16 convolutional layers with 2 fully connected layers, and a Softmax layer, where we implement a unique approach using different convolutional kernels within the same hidden layer for deeper feature extraction. We applied two different activation functions, MISH and rectified linear unit activation functions for deeper propagation of information and self regularized smooth non-monotonicity. Furthermore, we used a generalized intersection of union, thus overcoming issues such as scale invariance, rotation, and shape. Data augmentation techniques such as photometric and geometric distortions are adapted to overcome the obstacles faced in polyp detection. Detailed benchmarked results are provided, showing better performance in terms of precision, sensitivity, F1- score, F2- score, and dice-coefficient, thus proving the efficacy of the proposed model.

Keywords: Colonoscopy, convolutional neural network, MISH, polyp

[☆]Fully documented templates are available in the elsarticle package on CTAN.

Soo Young Shin

Email address: wdragon@kumoh.ac.kr (tariqrahim@ieee.org, Syedali@kumoh.ac.kr, wdragon@kumoh.ac.kr)

detection, precision, rectified linear unit, sensitivity.

1. Introduction

Diagnosing of distinct diseases within the small intestine is a time-consuming and hectic process for physicians. This has led to the introduction of technologies such as colonoscopy and wireless capsule endoscopy [1]. Colorectal cancer (CRC) is the second-highest cause of death by cancer worldwide with 880,792 deaths and a mortality rate of 47.60% in 2018 reported by American Cancer Society [2] 95% of CRC cases start with the appearance of a growth on the inner lining of the rectum or colon, called a polyp. Various types of polyps exist including, adenoma polyps, which can worsen into CRC. CRC is curable in 90% of cases assuming early detection [3]. Colonoscopy has emerged as minimally invasive and additional tool for investigating polyps by examining the gastrointestinal tract [3]. Colonoscopy relies on highly skillful endoscopists, and recent clinical investigations have shown colonoscopy misses 22% 28% of polyps. This false negatives can lead to late diagnosis of colon cancer, resulting in a survival rate as low 10% [4].

Deep learning (DL) is a subtype of machine learning concerned with the structure and function of brain-like systems known as artificial neural networks [5]. DL plays an important role in many areas, including text recognition tasks, self-driving cars, image recognition, and healthcare. Computer vision and machine learning-based methods have revolved over several decades to automatically detect polyps [6, 7, 8]. Such systems have generally examined, hand-crafted features, such as texture, histograms of oriented gradients, color wavelets, Haar, and local binary patterns [9, 10]. More advanced algorithms have been suggested to evaluate polyp appearance based on factors such as context information [11] and edge shape [12]. However, the decrease in detection performance is mainly due to the similar appearance of polyp-like and polyp structures.

Convolutional neural networks (CNNs) present promising outcomes in polyp detection and segmentation. CNN features outperformed hand-crafted features

in the MICCAI 2015 polyp detection challenge [6]. The Region-based CNN approaches, such as *R-CNN* [13], *Fast R-CNN* [14], and *Faster R-CNN* [15] have shown promising results for object detection in natural images. Work has also been done on regression-based object detection models such as You Only Look Once (YOLO) [16] and single shot multibox detector (SSD) [17]. However, recent investigations have shown that deep neural networks (DNNs), including CNNs, are extremely vulnerable to noise and perturbations [18]. Even one single-pixel addition increases the miss detection vulnerability of current DNNs including CNNs [19]. Even though computer-aided detection techniques can effectively classify frames from a colonoscopy, detection of polyps remains challenging due to significant size, appearance, and intensity variations between frames. This is a serious issue, because polyps and polyp-like objects have similar appearances in consecutive frames, leading to miss-detection even when implementing powerful models such as CNNs. Furthermore, the performance of DL approaches is highly correlated with the amount of data available for training. The lack of availability of labeled polyp images for training makes the detection and segmentation of the polyp a difficult task [20].

This work presents a new-CNN based detection model of polyp in colonoscopy images. The proposed CNN model employs fewer hidden layers, making the model lighter and less time-consuming during training. The proposed CNN model employs MISH as an activation function in some of the hidden layers for better deep propagation of information within the CNN [21]. Data augmentation such as photometric and geometric distortions is performed due to the scarcity of annotated polyp images generated from the colonoscopy process. The rest of the paper is categorized as follows: Section 2 presents recent related work done on polyp detection in colonoscopy images using DL. In Section 3, the proposed CNN model for polyp detection in colonoscopy images is explained in detail. In Section 4, the experimental results are described in detail, along with the dataset acquisition and augmentation process. Finally, in Section 5 the paper is reviewed and concluded and future work is presented.

2. Related Work

As CNN techniques for single and multiple object detection advance, they increasingly outperform previous conventional image processing techniques [22]. For multiple object detection, a region-based CNN combined with a deformable part-based model has been proposed to handle feature extraction and occlusion [23]. Recently, with the progress of DL in multiple image processing applications, a CNN-based method has been introduced for polyp detection [24, 25]. Due to this and related progress, CNN features outperformed hand-crafted features in the MICCAI 2015 polyp detection challenge [6]. A regression-based CNN model using ResYOLO combined with efficient convolution operators has been shown to successfully track and detect polyps in colonoscopy videos [7]. To avoid miss-detection of polyp between neighboring frames, a two-stage detector including a CNN-based object detector and a false-positive reduction unit can be applied [18]. Automatic detection of hyperplastic and adenomatous colorectal polyps in colonoscopy images has been performed using sequentially connected encoder-decoder based CNN [26]. Furthermore, automatic polyp detection in colonoscopy videos can be conducted via ensemble CNN, which learns a variety of polyp features such as texture, color, shape, and temporal information [27].

To overcome the lack of sufficient training samples for the use of pre-trained CNN on large-scale natural images, transfer learning systems have been proposed. This has been successfully implemented in various medical applications, such as automatic interleaving between radiology reports and diagnostic CT [28], MRI imaging, and ultrasound imaging [29]. Furthermore, the performance of various CNN architectures based on transfer learning, such as AlexNet and GoogLeNet has been evaluated for classification of interstitial lung disease and detection of thoracic-abnormal lymph nodes [30]. Similarly, a transfer learning-based method using the deep-CNN model Inception Resnet has been used to detect polyps in colonoscopy images [31]. Questions of whether a CNN with adequate fine-tuning can overcome the full training of the model from scratch have been answered in detail by examination of four different medical imag-

ing applications in three different specialties: gastroenterology, radiology, and cardiology for the purpose of classification, detection, and segmentation [32].

CNN has been used for decades in the field of computer vision for various applications. However, training a deep CNN model from scratch a complicated task [32]. Deep CNN models require a large amount of labeled training data. This becomes a difficult requirement when large-scale annotated medical data set are unavailable. Training the models is tedious and computationally time-consuming, and becomes even more so when facing complications such as overfitting and convergence. To overcome these issues, this work presents a form of CNN-based detection of polyps in colonoscopy images. The proposed model employs fewer hidden layers, making the model lighter and less time-consuming during training. The model uses MISH as an activation function in some of the hidden layers for better deep propagation of information within the CNN [21]. Data augmentation methods such as photometric and geometric distortions are used due to the scarcity of annotated polyp images generated from the colonoscopy process.

3. Proposed Deep Convolutional Neural Network (CNN) Architecture

Initially, the input image is divided into a grid during the training phase. Then the image is labeled using the RectLabel tool, generating a bounding box “ B ” consisting of five features. The horizontal and vertical components are labeled “ x ,” and “ y ,” respectively. Height and width are labeled “ h ” and “ w ,” respectively. Finally, a confidence score “ C_s ” is defined for each defined grid cell. The objective function of bounding box “ B ” is a bag of freebies using mean square error (MSE) to perform regression on the center coordinate points, height, and width of the box “ B ”. The intersection over union (IoU) is a vital indicator for estimating the distance between the predicted truth and the ground truth “ B ” and is given as generalized form as

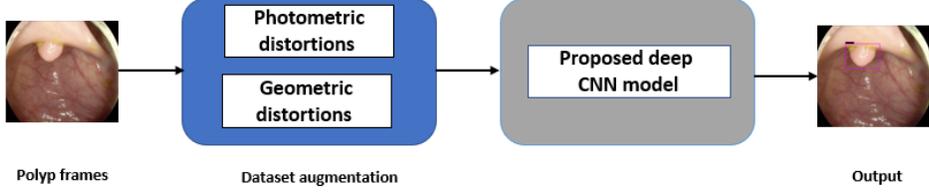


Figure 1: Flow of the proposed deep CNN model for polyp detection.

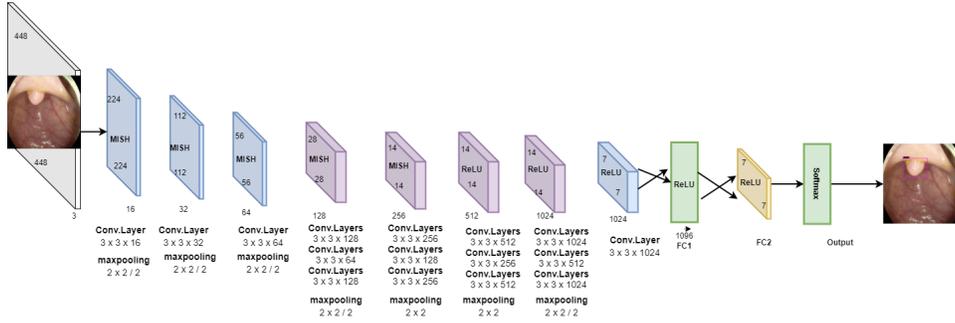


Figure 2: Architecture of the proposed deep CNN for polyp detection.

$$IoU = \frac{|E \cap F|}{|E \cup F|} = \frac{|I|}{|U|} \quad (1)$$

where “ E ” and “ F ” represent the predicted truth and ground truth, respectively. Here, the IoU distance $\mathcal{L}_{IoU} = 1 - IoU$ fulfills all properties of a metric, including the identity of indiscernibles, non-negativity, triangle inequality, and symmetry, but has a scale-invariant issue. The cost or loss function for object detection use l_1 norm and l_2 norm for x, y, w, h , but due to the scale-invariant property of IoU, there is an increase in loss with respect to scaling. In the proposed deep CNN approach to polyp detection, we have implemented generalized (GIoU) [33] as a new loss to optimize the non-overlapping “ B ” in consideration of the shape and orientation of the object in “ B ”. The GIoU finds the smallest convex shape $C \subseteq \mathbb{S} \in \mathbb{R}^n$ for two arbitrary convex shapes $E, F \subseteq \mathbb{S} \in \mathbb{R}^n$ followed by the calculation of the ratio between the area occupied by C minus “ E ” and “ F ”, divided by the total area occupied by “ C ”. Details for the algorithm

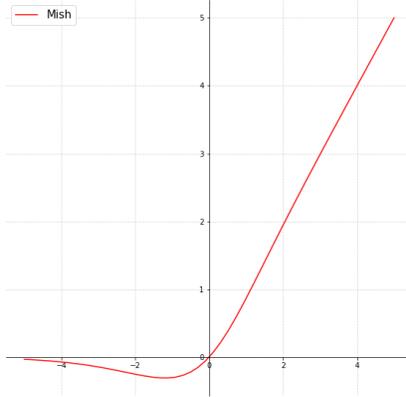


Figure 3: Mish activation function

and formulation can be found in [33], where a GIoU is expressed in simple mathematical form as, $GIoU = IoU - \frac{|C \setminus (E \cup F)|}{|C|}$. Furthermore, we have applied a non-maximum suppression algorithm involving “ C_s ” to avoid multiple and overlapping GIoUs.

Fig. 1 shows the general flow of the proposed approach for polyp detection in colonoscopy images. As shown in Fig. 2, the proposed deep CNN consists of 16 convolutional layers, two fully connected layers, and a softmax layer. To lessen computational complexity and improve hierarchical image features, *maxpooling* is used for the first 15 convolutional layers. For better image feature extraction, different sizes of *convolution kernels* are employed, with a stride of 2. In the proposed model, we have implemented *Mish* [21], which is a self-regularized smooth non-monotonic activation function, in the first 15 convolutional layers. This implementation was done after extensive trials to find the best matching position of the activation function. As observed in Fig. 3, *Mish* is an unbounded above result in avoiding saturation due to capping. This may normally lead to slow training, i.e., near-zero gradients. A better gradient flow and smooth propagation of information across deeper layers are achieved by the infinite order of continuity and a small allowance of negative values, in comparison to a strictly bounded rectified linear unit (ReLU) as an activation function. MISH

can be expressed mathematically as:

$$f(x) = x \cdot \tanh(\zeta(x)) \quad (2)$$

where, $\zeta(x) = \ln(1 + e^x)$ is softplus activation [21].

In the last layers, ReLU is used as an activation function to reduce the likelihood of gradient vanishing and achieve the sparsity. Flattening is done by two fully connected layers to yield a single continuous linear vector followed by *softmax* or the regression layer to generate the required output. The approach of using *Mish* and ReLU as an activation function results in smooth propagation of information across deeper layers. MISH helps to avoid capping, and ReLU prevents the gradient from vanishing.

The proposed deep CNN is trained with a similar concept of multi-layer perceptions i.e., a back-propagation algorithm which minimizes the cost function concerning the unknown weights “ W ”:

$$\mathcal{L} = -\frac{1}{|L|} \sum_i^{|L|} \ln(p(m^i|L^i)) \quad (3)$$

where $|L|$ represents the number of training images, $p(m^i|L^i)$ represents the probability that L^i is accurately classified, and L^i represents the i^{th} training image with the associated label m^i . We have applied stochastic gradient descent (SGD) as an optimizer, which minimizes the cost function over the whole training data set along with the cost over mini-batches of data. If W_j^t represents the weights in the j^{th} convolutional layer at t iteration, and $\hat{\mathcal{L}}$ represents the cost over a mini-batch of size M , then in the next iteration the updated weights are calculated as given below:

$$\begin{aligned} \gamma^t &= \gamma^{\lfloor tM/|X| \rfloor} \\ V_j^{t+1} &= \mu V_j^t - \gamma^t \eta_j \frac{\partial \hat{\mathcal{L}}}{\partial W_j} \\ W_j^{t+1} &= W_j^t + V_j^{t+1} \end{aligned} \quad (4)$$

where η_j is the learning rate of the j^{th} , μ is the momentum indicating the previously updated weight contribution in the current iteration, and γ represent

the scheduling rate which after each epoch decreases the learning rate η . The simulation parameters used for the proposed deep CNN are given in Table. 1.

Table 1: Simulation parameters used for deep CNN model

Network parameters	Configuration values
Input image dimension	448 X 448
Learning rate	0.0001
Optimizer	Stochastic gradient descent (SGD)
Momentum	0.9
Bath size	32
Iterations (t)	10,000

4. Experimental Results and Discussion

This section detail the data set specifications and experimental results generated by implementing the proposed deep CNN for the detection of polyps in colonoscopy images.

4.1. Dataset Specifications and Augmentation

The study used a publicly available dataset of polyp-frames obtained from the ETIS-Larib database [34], containing 196 polyp images. These images were obtained from 34 different colonoscopy videos of 44 different polyps with various appearances and sizes, having a resolution of 1225×966 pixels. The ground truth of polyp areas for polyp datasets is determined by expert video endoscopists. A CNN model trained with such a small amount of data is likely to be meaningless and unstable, so data augmentation was performed on the polyp dataset. Data augmentation had to be performed on the colonoscopy images by considering vivid variations. Otherwise over-fitting would have occurred. In a colonoscopy imagery, polyps exhibits large variations in location, color, and scale. Moreover, variations in brightness and definition also occur due varying the view-point of the camera. Therefore, in addition to photometric distortions and geometric

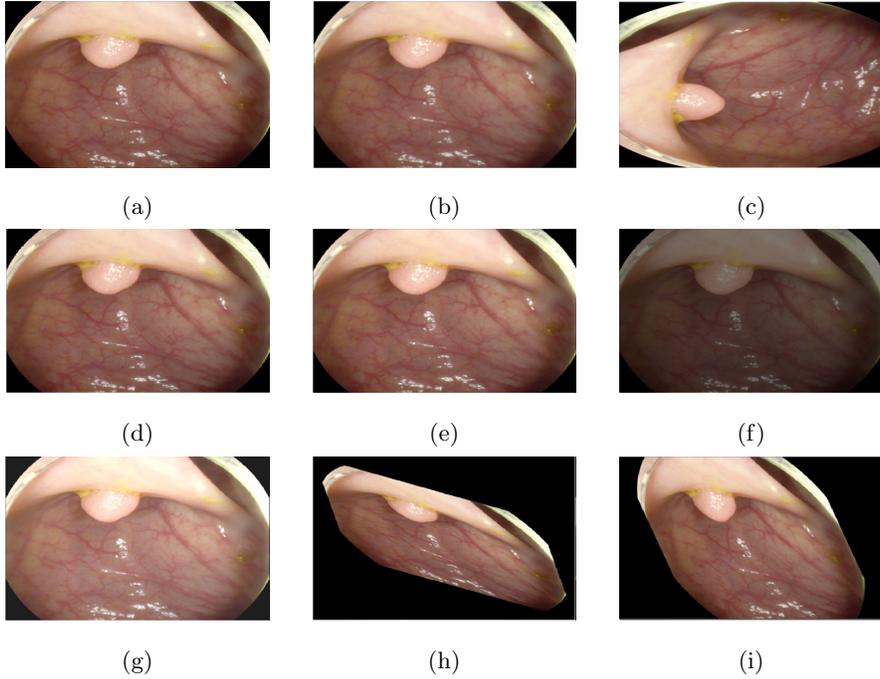


Figure 4: Image from the dataset with photometric and geometric augmentation. (a) original frame of polyp, (b) noisy polyp frame with $\sigma = 1.1$, (c) rotated polyp frame with 90° , (d) polyp frame with 15.00% zoom in, (e) polyp frame with 15.00% zoom out, (f) dark polyp frame, (g) bright polyp frame, (h) sheared polyp frame by y-axis, (i) sheared polyp frame by x-axis.

distortions, we also have considered zooming, shearing, and altering brightness as strategies for data augmentation.

For photometric distortions, we controlled brightness and contrast as an enhancement, while blurring by adding noise with a standard deviation (σ) of 1.0. Similarly, for geometric distortions, clock-wise rotation of the polyp images with angles of 90° , 180° , and 270° were performed. Zoom-in and zoom-out with zooming parameters such as 30.00% and 10.00% were performed to obtain different scales of polyp images. Lastly, shearing for both the x-axis and the y-axis was performed to shear the images from left to right and top to bottom, respectively. Fig. 4 shows photometric and geometric forms of image augmentation. In this way, we augmented the data set of the ETIS-Larib

database from 196 polyp images to 2,156 images, which is more suitable for training the proposed deep CNN model.

4.2. Performance Metrics for Evaluations

The metrics used to evaluate the detection of polyps within colonoscopy frames in this work are the same as those used in the MICCAI 2015 challenge [6]. The output obtained using the proposed model has rectangular shaped coordinates (x, y, w, h) . The following parameters are defined as follows:

True Positive (TP): True output detection if the detected centroid falls within the polyp ground truth. For multiple true output detection within the same frame and of the same polyp, TP is counted as one.

True Negative (TN): True detection, i.e., negative frames (frames without polyps) yielding no detection output.

False Positive (FP): False detection output where the detected centroid falls outside the polyp ground truth.

False Negative (FN): False detection output, i.e, polyp is missed in a frame having a polyp.

Employing the above parameters, we can compute the following performance metrics to efficiently evaluate the performance of the proposed deep CNN model.

Precision: This metric computes how precisely the model is detecting a polyp within an image

$$\text{Precision (Pre)} = \frac{TP}{TP + FP} \times 100 \quad (5)$$

Sensitivity: This metric is also called recall or True Positive Rate and computes the proportion of the actual polyps that were detected correctly

$$\text{Sensitivity (Sen)} = \frac{TP}{TP + FN} \times 100 \quad (6)$$

F1- score and F2- score: F1 and F2- score is simply the harmonic mean between precision and sensitivity , in a range of $[0, 1]$. Both scores are recognized to balance the precision and sensitivity. The F1- score is given as:

$$F1 - \text{score} = \frac{2 \times \text{Sen} \times \text{Pre}}{\text{Sen} + \text{Pre}} \times 100 \quad (7)$$

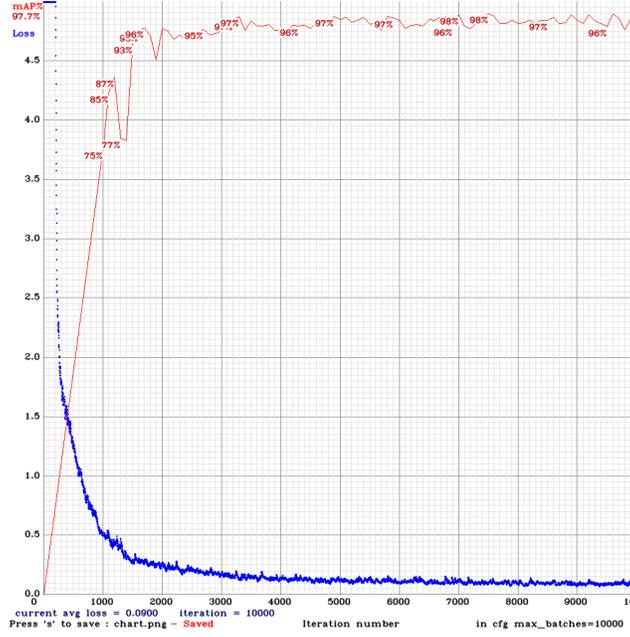


Figure 5: Training phase of the proposed deep CNN model.

while the F2- score can be calculated as:

$$F2 - score = \frac{5 \times Pre \times Sen}{4 \times Pre + Sen} \quad (8)$$

Dice Coefficient: This metric is used for pixel-wise result comparison between ground truth and predicted detection that ranges $[0, 1]$, and is given as:

$$Dice\ coefficient\ (E, F) = \frac{2 \times |E \cap F|}{|E| + |F|} = \frac{2 \times TP}{2 \times TP + FP + FN} \quad (9)$$

4.3. Polyp Frames Evaluation

This section reports the polyp detection performance of the proposed CNN model. For implementation of the model, 80.00% and 20.00% of the 2,156 augmented polyp frames were used for training and testing, respectively. Fig. 5 shows the real-time training phase of the proposed model, where 10,000 iterations were run to achieve the best weights. The model was trained using the

simulation parameters as given in Table. 1, both for the non-augmented ETIS-Larib database [34] containing 196 poly images, and the augmented data set, for fair performance comparison. A high mean average precision of 97.70% with an MSE of 0.900 was obtained in the early iterations, resulting in the best weights for testing purposes.

The results listed in Table. 2 using the proposed deep CNN model, show better performance, with high values for precision, the F1-score, the F2-score, and the dice coefficients. Note that the low sensitivity or recall values is an indication of better polyp detection performance in the proposed model. As observed in Table. 2, the proposed model is compared to other works [31] that showed better performance for both the non-augmented and augmented case. For the non-augmented data set of ETIS-Larib, the generated TP, FP, and FN values were 90, 35, and 51, respectively. Similarly, 20.00% of the augmented data set was employed for testing purposes, generating TP, FP, and FN values of 340, 20, and 70, respectively.

Table 2: Detection performance comparison of the proposed deep CNN model on ETIS-Larib database without(w/0) and with augmentation strategies.

Data set	Performance metrics			
	[31] (%)		Propose deep CNN model (%)	
Non-augmented ETIS LARIB database (196)	Pre	48.00	Pre	72.00
	Sen	39.40	Sen	63.82
	F1- score	43.30	F1- score	67.66
	F2- score	40.90	F2- score	65.30
	Dice-coefficient	NA	Dice-coefficient	0.676
Augmented ETIS LARIB database (2,156)	Pre	91.40	Pre	94.44
	Sen	71.20	Sen	82.92
	F1- score	80.00	F1- score	88.30
	F2- score	74.50	F2- score	85.00
	Dice-coefficient	NA	Dice-coefficient	0.88

The results shown in Fig. 6 are generated using the proposed deep CNN model on the augmented data set. It can be observed that the proposed model

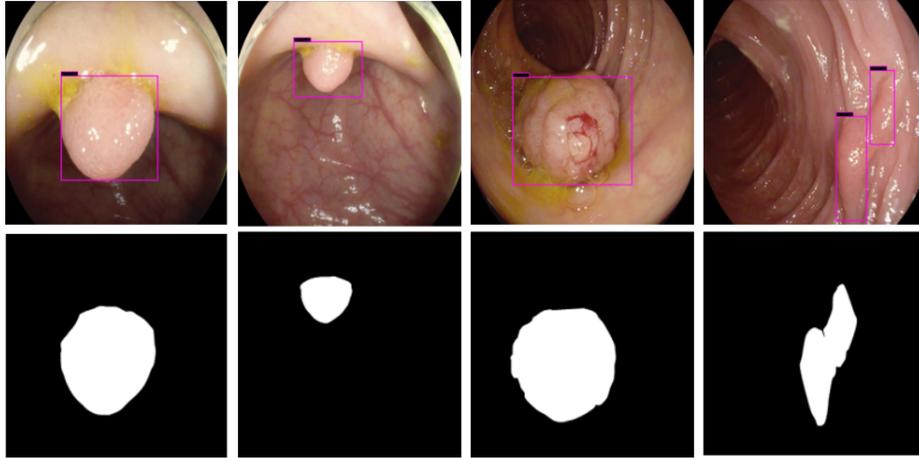


Figure 6: Example of accurate detection along with the correct ground truth using deep CNN model. The first row shows the detection results for different polyp from data augmentation process. The second shows the ground truth images of test images.

shows better polyp detection performance. As illustrated in Fig. 6, polyps within a frame can be identified at multiple positions, and as noted above in this case, the TP for detection is considered to be 1. The proposed deep CNN model performed better than other benchmark results in terms of the performance metrics listed above, as shown in Table. 2 and Fig. 6.

For single and deep layer of the proposed model, we have shown channel activation representing the convolutional kernels accurately detected the polyp. Fig. 7 shows different bright and dark parts corresponding to the spatial property of the object within the test images for single and deep layers. The *top left* is the test polyp image followed by *top right* detection output generated by proposed deep CNN model. The *bottom left* shows the single layer activation channel while *bottom right* shows the deep layer for deeper feature analysis represented by green rectangular boxes. It can be observed in Fig. 7, that both single and deep layers are extracting polyp features with a high score, resulting in high polyp detection.

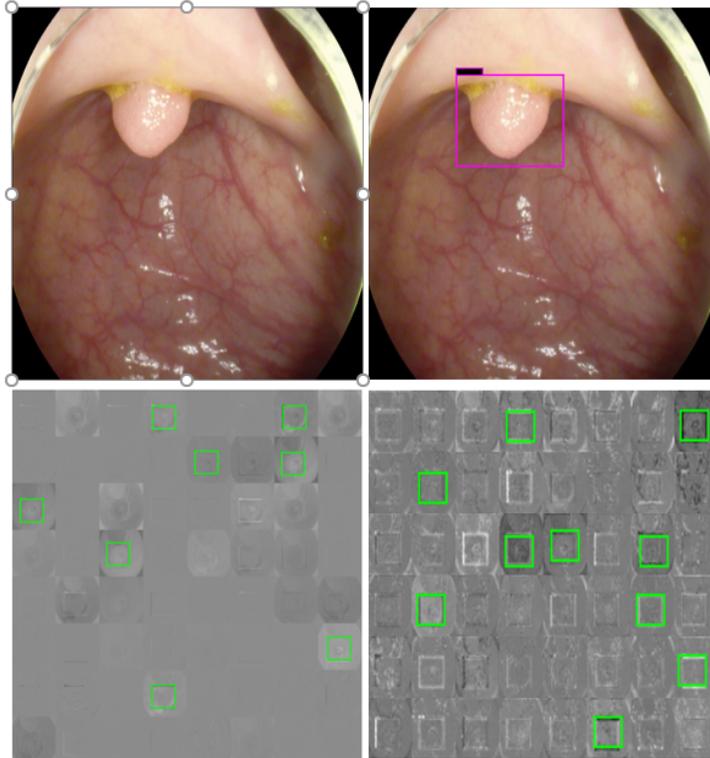


Figure 7: Test polyp channel activation visualization of CNN after training of the proposed deep CNN model. Top left: test polyp images, Top right: detection output of the test image, Bottom left: activation channel for single layer, Bottom right: activation channel for deep layer (both shown by green rectangular box).

4.4. Benchmark Performance with Other Approaches

The performance of the proposed deep CNN model was benchmarked against the 2015 MICCAI challenge [6], as the dataset used was the same. The top three experimental results from each team in the challenge UNS-UCLAN, OUS, and CUMED were selected for benchmarking. These results were selected because CNN has been used for learning end-to-end detection of the polyp. The UNS-UCLAN team [6] used three CNNs for the extraction of features on multiple spatial scales, followed by a classification approach with a multi-layer perceptron network. AlexNet, a CNN-based model was adopted with a conventional sliding window to perform patch-based classification [31]. The CUMED team used a

Table 3: Performance comparison of the proposed deep CNN model on ETIS-Larib database with other methods.

Implemented methods	Performance metrics				
	Pre (%)	Sen (%)	F1- score (%)	F2- score (%)	Dice-coefficient
UNS-ULCAN	32.70	52.80	40.40	47.10	NA
OUS	69.70	63.00	66.10	64.20	NA
CUMED	72.30	69.20	70.70	69.80	NA
Proposed deep CNN model	94.44	82.92	88.30	85.00	0.88

segmentation approach based on CNN [34], where classification was conducted pixel-wise along with a ground truth mask.

As shown in Table. 3, the generated results from the proposed model using the augmented dataset outperform the other team’s methods on several metrics, including precision, sensitivity, F1- score, F2- score, and dice-coefficient. As DL-based methods employ different computers with different specifications, it is hard to benchmark detection processing directly. In our work, the dataset was trained and tested on NVIDIA Titan RTX GPUs to reduce processing time. Compared to the other studies listed in Table. 3 the mean detection processing time 0.6 sec per frame. This is slightly greater than that in competing models, but the increased processing time comes with better performance.

Conclusion

In this paper, we presented a computerized DL-based detection model for colonic polyps. A deep CNN model consisting of 16 convolutional layers with two full-connected layers, and a Softmax layer, was implemented with different kernel sizes in the same hidden layer being employed. Moreover, two different activation functions MISH and ReLU were implemented for the first time to provide deeper propagation of information, better self-regularization, and better capping avoidance. The scale invariance issue related to IoU was addressed by adopting a GIoU that is robust under rotation and shape variation. Furthermore, photometric and geometric strategies were used for data augmentation,

thus overcoming the image scarcity issue. We provided a detailed benchmark performance comparison of our detection output, which outperforms the other approaches in performance metrics such as precision, sensitivity, F1- score, F2-score, and dice-coefficient.

Acknowledgments

This work was supported by Priority Research Centers Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education, Science and Technology(2018R1A6A1A03024003).

References

- [1] T. Rahim, M. A. Usman, S. Y. Shin, A survey on contemporary computer-aided tumor, polyp, and ulcer detection methods in wireless capsule endoscopy imaging, arXiv preprint arXiv:1910.00265.
- [2] R. Segal, K. Miller, A. Jemal, Cancer statistics, 2018, *Ca Cancer J Clin* 68 (1) (2018) 7–30.
- [3] O. Chuquimia, A. Pinna, X. Dray, B. Granado, Polyp follow-up in an intelligent wireless capsule endoscopy, in: 2019 IEEE Biomedical Circuits and Systems Conference (BioCAS), IEEE, 2019, pp. 1–4.
- [4] A. Leufkens, M. Van Oijen, F. Vleggaar, P. Siersema, Factors influencing the miss rate of polyps in a back-to-back colonoscopy study, *Endoscopy* 44 (05) (2012) 470–475.
- [5] S. A. Hassan, T. Rahim, S. Y. Shin, Real-time uav detection based on deep learning network, in: 2019 International Conference on Information and Communication Technology Convergence (ICTC), IEEE, 2019, pp. 630–632.
- [6] J. Bernal, N. Tajkbaksh, F. J. Sánchez, B. J. Matuszewski, H. Chen, L. Yu, Q. Angermann, O. Romain, B. Rustad, I. Balasingham, et al., Comparative

validation of polyp detection methods in video colonoscopy: results from the miccai 2015 endoscopic vision challenge, *IEEE transactions on medical imaging* 36 (6) (2017) 1231–1249.

- [7] R. Zhang, Y. Zheng, C. C. Poon, D. Shen, J. Y. Lau, Polyp detection during colonoscopy using a regression-based convolutional neural network with a tracker, *Pattern recognition* 83 (2018) 209–219.
- [8] N. Tajbakhsh, S. R. Gurudu, J. Liang, System and methods for automatic polyp detection using convolutional neural networks, uS Patent 10,055,843 (Aug. 21 2018).
- [9] J. Bernal, J. Sánchez, F. Vilarino, Towards automatic polyp detection with a polyp appearance model, *Pattern Recognition* 45 (9) (2012) 3166–3182.
- [10] S. Y. Park, D. Sargent, I. Spofford, K. G. Vosburgh, A. Yousif, et al., A colon video analysis framework for polyp detection, *IEEE Transactions on Biomedical Engineering* 59 (5) (2012) 1408–1418.
- [11] N. Tajbakhsh, S. R. Gurudu, J. Liang, Automated polyp detection in colonoscopy videos using shape and context information, *IEEE transactions on medical imaging* 35 (2) (2015) 630–644.
- [12] J. Bernal, F. J. Sánchez, G. Fernández-Esparrach, D. Gil, C. Rodríguez, F. Vilariño, Wm-dova maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from physicians, *Computerized Medical Imaging and Graphics* 43 (2015) 99–111.
- [13] R. Girshick, J. Donahue, T. Darrell, J. Malik, Rich feature hierarchies for accurate object detection and semantic segmentation, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 580–587.
- [14] R. Girshick, Fast r-cnn, in: *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1440–1448.

- [15] S. Ren, K. He, R. Girshick, J. Sun, Faster r-cnn: Towards real-time object detection with region proposal networks, in: *Advances in neural information processing systems*, 2015, pp. 91–99.
- [16] J. Redmon, S. Divvala, R. Girshick, A. Farhadi, You only look once: Unified, real-time object detection, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.
- [17] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, A. C. Berg, Ssd: Single shot multibox detector, in: *European conference on computer vision*, Springer, 2016, pp. 21–37.
- [18] H. A. Qadir, I. Balasingham, J. Solhusvik, J. Bergsland, L. Aabakken, Y. Shin, Improving automatic polyp detection using cnn by exploiting temporal dependency in colonoscopy video, *IEEE Journal of Biomedical and Health Informatics*.
- [19] J. Su, D. V. Vargas, K. Sakurai, One pixel attack for fooling deep neural networks, *IEEE Transactions on Evolutionary Computation* 23 (5) (2019) 828–841.
- [20] W.-L. Chao, H. Manickavasagan, S. G. Krishna, Application of artificial intelligence in the detection and differentiation of colon polyps: a technical review for physicians, *Diagnostics* 9 (3) (2019) 99.
- [21] D. Misra, Mish: A self regularized non-monotonic neural activation function, arXiv preprint arXiv:1908.08681.
- [22] Q. Zhang, N. Huang, L. Yao, D. Zhang, C. Shan, J. Han, Rgb-t salient object detection via fusing multi-level cnn features, *IEEE Transactions on Image Processing*.
- [23] J. Li, H.-C. Wong, S.-L. Lo, Y. Xin, Multiple object detection by a deformable part-based model and an r-cnn, *IEEE Signal Processing Letters* 25 (2) (2018) 288–292.

- [24] S. Y. Park, D. Sargent, Colonoscopic polyp detection using convolutional neural networks, in: *Medical Imaging 2016: Computer-Aided Diagnosis*, Vol. 9785, International Society for Optics and Photonics, 2016, p. 978528.
- [25] S. Park, M. Lee, N. Kwak, Polyp detection in colonoscopy videos using deeply-learned hierarchical features, Seoul National University.
- [26] D. Bravo, J. Ruano, M. Gómez, E. Romero, Automatic polyp detection and localization during colonoscopy using convolutional neural networks, in: *Medical Imaging 2020: Computer-Aided Diagnosis*, Vol. 11314, International Society for Optics and Photonics, 2020, p. 113143A.
- [27] N. Tajbakhsh, S. R. Gurudu, J. Liang, Automatic polyp detection in colonoscopy videos using an ensemble of convolutional neural networks, in: *2015 IEEE 12th International Symposium on Biomedical Imaging (ISBI)*, IEEE, 2015, pp. 79–83.
- [28] H.-C. Shin, L. Lu, L. Kim, A. Seff, J. Yao, R. M. Summers, Interleaved text/image deep mining on a very large-scale radiology database, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1090–1099.
- [29] H. Chen, D. Ni, J. Qin, S. Li, X. Yang, T. Wang, P. A. Heng, Standard plane localization in fetal ultrasound via domain transferred deep neural networks, *IEEE journal of biomedical and health informatics* 19 (5) (2015) 1627–1636.
- [30] H.-C. Shin, H. R. Roth, M. Gao, L. Lu, Z. Xu, I. Nogues, J. Yao, D. Mollura, R. M. Summers, Deep convolutional neural networks for computer-aided detection: Cnn architectures, dataset characteristics and transfer learning, *IEEE transactions on medical imaging* 35 (5) (2016) 1285–1298.
- [31] Y. Shin, H. A. Qadir, L. Aabakken, J. Bergsland, I. Balasingham, Automatic colon polyp detection using region based deep cnn and post learning approaches, *IEEE Access* 6 (2018) 40950–40962.

- [32] N. Tajbakhsh, J. Y. Shin, S. R. Gurudu, R. T. Hurst, C. B. Kendall, M. B. Gotway, J. Liang, Convolutional neural networks for medical image analysis: Full training or fine tuning?, *IEEE transactions on medical imaging* 35 (5) (2016) 1299–1312.
- [33] H. Rezatofghi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, S. Savarese, Generalized intersection over union: A metric and a loss for bounding box regression, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 658–666.
- [34] J. Silva, A. Histace, O. Romain, X. Dray, B. Granado, Toward embedded detection of polyps in wce images for early diagnosis of colorectal cancer, *International Journal of Computer Assisted Radiology and Surgery* 9 (2) (2014) 283–293.