

Enhancing Endoscopic Image Classification with Symptom Localization and Data Augmentation

Trung-Hieu Hoang
hthieu@selab.hcmus.edu.vn
University of Science, VNU-HCM
Vietnam

Thanh-An Nguyen
ntan@selab.hcmus.edu.vn
University of Science, VNU-HCM
Vietnam

Hai-Dang Nguyen
nguyenhd@eurecom.fr
Eurecom
France

Vinh-Tiep Nguyen
tiepvn@uit.edu.vn
University of Information Technology,
VNU-HCM
Vietnam

Viet-Anh Nguyen
nvanh160013@ump.edu.vn
University of Medicine and Pharmacy,
Ho Chi Minh city
VietNam

Minh-Triet Tran
tmtriet@fit.hcmus.edu.vn
University of Science, VNU-HCM
Vietnam

ABSTRACT

Inspired by recent advances in computer vision and deep learning, we proposed new enhancements to tackle problems appearing in endoscopic image analysis, especially abnormalities findings and anatomical landmarks detection. In details, a combination of **Residual Neural Network** and **Faster R-CNN** are applied jointly in order to take all of their advantages and improve the overall performance. Nevertheless, novel data augmentation has been designed and adapted to corresponding domains. Our approaches prove their competitive results in term of not only the accuracy but also the inference time in Medico: The 2018 Multimedia for Medicine Task and The Biomedia ACM MM Grand Challenge 2019. These results show the great potential of the collaborating between deep learning models and data augmentation in medical image analysis applications.

KEYWORDS

endoscopic image, symptoms localization, anatomical landmarks detection, object detection, image classification

1 INTRODUCTION

The abnormalities finding and landmark detection in endoscopy images challenge aim to bring new achievements in computer vision, image processing and machine learning to the next level of computer and multimedia assisted diagnosis. In order to encourage the research community to tackle the problems with endoscopy images, besides proposing new endoscopic images dataset [8], several challenges are conducted, such as the Medico: Multimedia Task at MediaEval 2018 [7] and The Biomedia ACM MM Grand Challenge 2019 [13]. The goal of these challenges is encouraging research communities to establish an assisted diagnosis systems that can detect abnormalities and anatomy landmarks in human gastrointestinal tract automatically in an efficient way with as less training data as possible.

In our approach, we introduce a stacked model consisting of two deep networks, a Residual Neural Network (Resnet) [2] followed by a Faster Region-based Convolutional Neural Network (Faster R-CNN) [10]. As our observation, the Resnet mostly focuses on deep global features of image, it fails to classify images that symptoms

of *abnormal symptoms* or *instruments* appear as small objects on diversity backgrounds. Therefore, we aim to apply the Faster R-CNN to re-classify the images of some classes that Resnet usually mis-classify.

Besides, due to the limitation and the imbalance between classes in the training samples and using extra source of endoscopic data is infeasible, we proposed a novel data augmentation mechanism which can enhance the performance of both models. Nevertheless, a multi-tasks classifiers is introduced to reduce the confusion level of classifier module in cases that multiple symptoms of diseases appeared in the same image.

2 RELATED WORK

Early approaches that tackle with endoscopic image analysis usually work on the bleeding detection problem. By using color threshold and applying some handcrafted features, Shah et al.[6] and Jung et al.[5] propose several solutions using color domain and region segmentation. Later, supervised learning machine techniques are used jointly with superpixel saliency to identify bleeding regions of by recognizing blood color patterns [3].

Since introduced, deep neural networks have been used in order to solve several problems in the field of analyzing endoscopic images of the gastrointestinal (GI) tract. Particularly, to localize and identify polyps within real-time constraint, deep CNNs has recently shown an impressive potential when achieving up to 96.4% accuracy - published in 2018 by Urban G et al. [14]. Another interesting article of Satoki Shichijo et al. [12] also applies multiple deep CNNs to diagnose Helicobacter pylori gastritis based on endoscopic images. Further, gastrointestinal bleeding detection using deep CNNs on endoscopic images has been successfully done and published by Xiao Jia et al. [4].

3 APPROACH

3.1 Overview of proposed methods

As mentioned before, a stacked model consists of a Residual Neural Network (ResNet) and a Faster Region-based Convolutional Neural Network (Faster R-CNN). In order to training the Faster R-CNN model, addition information regarding to the location of abnormal symptom need to be annotated which is described in details in



Figure 1: Our proposed symptoms region localization, annotated images with bounding boxes and classes name.

section 3.2. Additionally, unbalancing between classes can make deep models bias to some classes than others, which reduce the overall accuracy. In the context of these challenges that require using as less as training data as possible, therefore, using extra data from other dataset is infeasible. In Section 3.3, we also provide some augmentation mechanism on the given training dataset in order to provide a better version for the training step of both modules.

Nevertheless, inference time must be taken into account, while the Faster R-CNN needs longer time to process than that of ResNet. Instead of feeding all images through the ResNet module, we should have different strategies in Section 3.5 that reduce the number of images need to go through the Faster R-CNN module and as result, balancing between the accuracy and inference time.

Last but not least, an other improvement of our method is using a multi-task classification architecture that can predict multiple classes at the same time, which can reduce the confusing level of the image classifier model in these difficult cases. Further information about this architecture is presented in Section 3.6.

3.2 From Classification to Symptoms Region Localization

An object detection module can be really useful to detect small abnormal symptoms and diseases, such as *polyps*, *instruments*, *dyed-lifted-polyps* and *dyed-resection-margins*. Besides, in some cases that multi-symptoms appear in the same image, identify all of them is necessary to draw the final conclusion. Nevertheless, system that can not only predict the abnormalities but also propose the corresponding positions of that is more reliable and more convenient for endoscopists.

Due to the limitation of the given training dataset, there is no extra information about positions of abnormalities corresponding to each endoscopic image. Therefore, we decide to annotate all the abnormal symptoms in every images of the following classes: *dyed-resection-margins*, *dyed-lifted-polyps*, *instruments* and *polyps*. Examples of our proposed bounding box can be seen in Figure 1.

Totally, 5241 images belonged to the mentioned classes, are annotated with 5715 bounding boxes from the Kvasir dataset [8] and Medico 2018 development dataset. Although, the number of training samples we annotated is much larger than the number of actual training samples used in training phase, it is still useful for future works.

3.3 Instruments, polyps dataset Augmentation

Noticeably, *Instruments* - the second highest priority class has only 36 training samples. In order to maintain the balancing between all of these classes and also improve the diversity of the *instruments* images, we aim to augment the given dataset by generate more

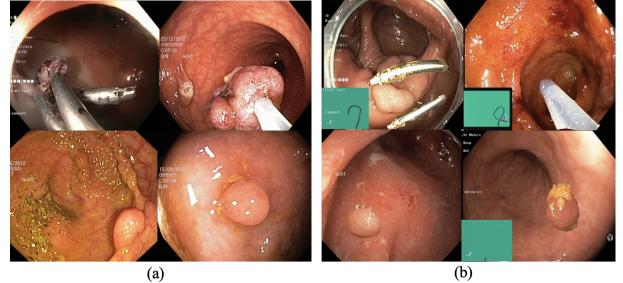


Figure 2: Dataset augmentation result. (a) original samples (b) augmented samples generated by our method

images for the *instruments* based on the current given development set.

From a given image, we carefully crop the region that contains symptom of diseases or *instruments* along their edges. Then, we randomly select images from other classes and use them as the background of the cropped instruments. Random affine transformation are also applied. With this strategy, we are able to generate more than 800 images for the target classes. As can be seen in Figure 2, there is not significant visual differences between original and augmented samples.

The augmented dataset can be used to enhance the robustness of the classification model. On the other hand, we can easily get the position of the foreground object on generated images, which can be utilized to augment the training dataset of the objects detection module.

The augmented dataset can be used for both sub-task. By adding more positive examples to the training phase of the classification model, the performance of this model can be more robust. On the other hand, we can easily get the position of the foreground object on generated images, which can be utilized to augment the training dataset of the objects detection module (Faster R-CNN).

3.4 Fine-tuning Deep Neural Network on Endoscopic images

Besides high computational cost, one of the main drawback of deep learning architecture is that it requires a large amount of training data. Moreover, labeled medical data for supervised learning is limited and manual labelling of medical images is a difficult task. Therefore, it requires a lot of effort and time to train the network, which would depend on the size of training data used. However, there is a possible solution to deal with these limitations is using transfer learning, where a pre-trained network on a large dataset (such as ImageNet [11]) is used.

In our approach, both Residual Network with 101 layers and Faster R-CNN [1] are both share a same features encoder. Therefore, it is necessary to propose a features encoder that is specialized on endoscopic images. This is the reason that we fine-tuned our deep neural network models (pre-trained on ImageNet) by using our modified development dataset. After training the whole neural network and then we freeze several first layers in its architecture and fine-tune the remains with small learning-rate. We also tried to train the network from scratch and all of our experiments point

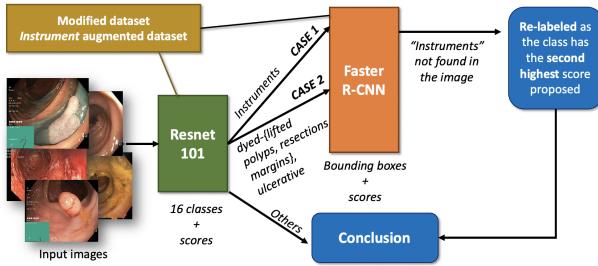


Figure 3: Configuration of the Third scenario (best performance).

out that in term of using convolution neural network for medical images, knowledge transferring from natural images to medical images is possible, even though there is a large difference between the source and the target databases.

This idea is also mentioned in [9], it is especially useful in the case of small dataset of images provided. Fine-tuning on the ImageNet pre-trained model significantly improves the efficient of deep learning model on medical domains.

3.5 Configuration of Conditional Scenarios

As mentioned before, there is a trade-off between the inference time and the accuracy. The following section describes the pipeline for each of our scenarios that we proposed in order to evaluate the performance in term of the accuracy and inference time.

First scenario: instruments, polyps double-checked. Residual network with 101 layers model are fine-tuned on the original development set provided by the task organizers along with our instruments increased dataset. After passed through ResNet101, output images classified as special classes become the input of Faster R-CNN network, which is trained for detecting instruments in images.

- First case: Images predicted as *instruments* by Resnet101 are double-checked. In case instruments are not detected by Faster R-CNN in those images, they are re-labeled as the class of their second highest score proposed by Resnet101.
- Second case: Images predicted as *dyed-lifted-polyps*, *dyed-resection-margins*, *ulcerative colitis* by ResNet101 are fed forward through Faster R-CNN network to detect *instruments*. They are classified as *instruments* if detected or keep the original prediction otherwise.

Second scenario: instruments double-checked. Feeding forward a large number of images in the three classes through Faster R-CNN causes a bottle-neck of inference time, as Faster R-CNN has high time complexity. Therefore, in this second , we limited the images passed through Faster R-CNN by only performing the first case of the first scenario.

Third scenario: instruments double-checked and data augmentation. The configuration of the third scenario is as same as the second scenario which is illustrated in Figure 3. Instead of using the original training set mentioned in the first scenario, we train our model on the re-labeled development set combined with the augmented instrument set.

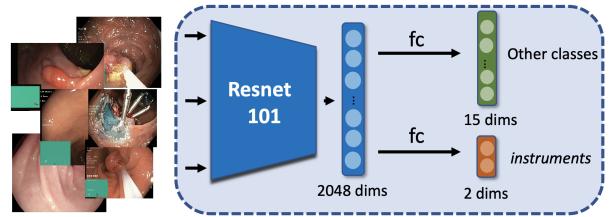


Figure 4: Overview of the multi-classes classifier network.

Forth scenario: 75% of training set. In this scenario, we reduce the number of images used for training by selecting randomly 75% images of each class in the same training set as the third scenario. Other processing steps are also configured in the same way.

Fifth scenario: esophagitis priority. Throughout our experiments, *normal-z-line* and *esophagitis* are the top most confusing classes not only for Resnet101 but also for human to distinguish them. In the priority list, *esophagitis* has a higher rank than *normal-z-line*'s. Thus, after several times evaluating our model on the development dataset, we propose a condition for these two classes when they are predicted by ResNet101. As ResNet101 provides a probability distribution over the 16 classes for each image, whenever the *normal-z-line* appears to be the highest class, we add a small bias 0.3 to the probability of the *esophagitis*. Hence, the model is more likely to emit the *esophagitis* class. This intuitively means that our model prefers *esophagitis* to *normal-z-line* when it is confused between these classes.

3.6 Solving Multi-classes Problem with Multi-tasks Classifier

After Medico 2018, another improvement of our work is to reduce the confusion level of the deep neural network model in cases that various type of abnormalities appeared in a same image. Instead of using only one classifier and forcing the deep neural network model to follow the priority list in these cases by feeding a number of positive and negative samples, we narrow down this job to multiple classifiers. For instance, given an image that both *esophagitis* and *instruments* appear simultaneously, the job to determine whether or not *instruments* appear inside the image is then left for a 2-classes classifier. This classifier can output the probability that the given image contains *instruments*. Meanwhile, the second classifier works independently, which can output the probability for other classes, except *instruments*.

Architecture. In order to reduce the inference time, we decide to share the weights of backbone ResNet 101 for both classifier. The overview of this module can be seen in Figure 4. In general, the proposed multi-task classification model consists of a ResNet 101 architecture except the last fully connected layer, working as a features extractor module that can output a 2048 dimensions vector for each input image. There are two fully connected branches on top of the output of that features extractor in order to get the prediction of *instruments* class and other 15 classes. The number of classifiers is extendable in the future.

Table 1: Official evaluation result of Medico: : The 2018 Multimedia for Medicine Task for both sub-tasks (provided by the organizers) and speed (fps) on Tesla K80 GPU

RunID	PREC	REC	ACC	F1	MCC	RK	FPS
Run01	94.245	94.245	99.281	94.245	93.861	93.590	6.589
Run02	93.959	93.959	99.245	93.959	93.556	93.273	23.191
Run03	94.600	94.600	99.325	94.600	94.240	93.987	23.148
Run04	93.043	93.043	99.130	93.043	92.579	92.257	22.654
Run05	94.508	94.508	99.314	94.508	94.142	93.884	21.413

Loss function. Two classifiers are trained simultaneously with the overall loss function is a weighted sum of each of their loss, given as follow

$$\mathcal{L}(p_i, p_o, p_i^*, p_o^*) = \lambda \sum_i \mathcal{L}_{instr}(p_i, p_i^*) + (1-\lambda) \sum_i \mathcal{L}_{others}(p_o, p_o^*) \quad (1)$$

where (p_i, p_i^*) and (p_o, p_o^*) denote the prediction and the ground-truth of *instr* class and other classes, respectively. \mathcal{L} is Cross Entropy Loss function. λ is a combination weight.

Final prediction. Since, there are two output vectors from the model, the final prediction can be determined by

$$y_f = \begin{cases} argmax(p_o) & p_i < 0.5 \\ y_{instruments} & p_i \geq 0.5 \end{cases} \quad (2)$$

where y_f stands for the final prediction of input image, p_o is a 15 dimensions vector indicates the probability that image is likely to belong to. p_i is the probability that *instruments* appear inside the image and $y_{instruments}$ is the corresponding label of *instruments* class.

4 EXPERIMENTAL SETUP

Medico: The 2018 Multimedia for Medicine Task

The configuration of *Run01* to *Run05* corresponding to five scenarios described in section 3.5.

Biomedia ACM MM Grand Challenge 2019

We continue to tackle existing problems that we did not solve efficiently in Medico 2018, which are the confusing between *normal-z-line* and *esophagitis*; multi-classes problem with *instruments* classes. Besides increasing the size of training data with our augmentation strategy, there are two major improvements in this challenge.

- (1) With the multi-classes problem, **we applied the Multi-tasks Classifier (MUL)** which is introduced in Section 3.6.
- (2) With the *normal-z-line* and *esophagitis*, with a help of a medical expert, we re-annotate the labels of the original dataset and train our models on the modified version of the dataset (*RE_LBL*).

5 RESULTS

As illustrated in Table 1, there is a trade-off between speed and accuracy when comparing the result of *Run01* and *Run02*. When reducing a large number of images passing through Faster R-CNN for the sake of time, so its performance seems to be relatively worse than *Run01*'s.

Table 2: Official evaluation result of The Biomedia ACM MM Grand Challenge 2019 for both sub-tasks (provided by the organizers) and speed (fps) on GTX 1080 Ti GPU

RunID	PREC	REC	F1	MCC	FPS
SIN	0.8565	0.8503	0.8458	0.9126	3.5461
MUL	0.8763	0.8771	0.8746	0.9406	3.6101
SIN+RE_LBL	0.8667	0.8567	0.8483	0.9168	3.5842
MUL+RE_LBL	0.8698	0.8573	0.8491	0.9202	3.5842

As we mentioned earlier, training samples takes an important role in building a deep-neural network model. Through our experiments, in the case of less training data, the augmented dataset helps us improve the performance of deep-neural network model. *Run03* and *Run05* show impressive results comparing to the first two runs. This implies that training on our re-labeled development set provides better models.

On the other hand, using the Residual neural network cannot classify efficiently the two classes *esophagitis* and *normal-z-line*. The same problem also occurs between the *dyed-resection-margins* and *dyed-lifted-polyps* classes.

Additionally, the configuration of *Run05* intuitively prefers *esophagitis* to *normal-z-line*, which may leads to an increasing of the false-positive cases in the result.

By comparison to the others, *Run04* has the lowest precision since it uses 75% of training data. Decreasing the amount of training samples of course affects the performance in deep-learning models. Nevertheless, the result is still acceptable when it decreases only a few percentages and its configuration is as same as *Run03*.

Regarding to the evaluation results in Table 2, the performance of the Multi-tasks classifiers has successfully proved to have better performance than the single-task classifier used in previous experiments.

Although, having competitive results on our own validation set, re-labeling the development dataset on *esophagitis* and *normal-z-line* approach seem to have worse performance on the official test set. Distinguishing these classes is an challenging problem, especially for computers since there are still disagreements between experts in some cases.

6 CONCLUSION AND FUTURE WORKS

Endoscopic image classification is a challenging problem because of the fine-grained images, less training data and require high accuracy. In our approach, we focus on enhancing the performance of image classification model, such as Residual Neural Network by using object detection model, such as Faster R-CNN that can localize small symptoms of diseases, which are useful evidences. Besides, data augmentation strategy can be applied to solve the limitation of training samples which is commonly occurred in medical datasets. Accuracy and inference time that we reach is acceptable and appropriate for real-time constraint. However, for future works, additional medical testing must be taken into account besides visual information to create a more robust approach to exploit the distinction between easy-confused classes.

ACKNOWLEDGMENTS

We would like to express our appreciation to Honors Program of Computer Science, Software-engineering Laboratory, University of Science, VNU-HCM

REFERENCES

- [1] Xinlei Chen and Abhinav Gupta. An implementation of faster r-cnn with study for region sampling. *arXiv preprint arXiv:1702.02138*, 2017.
- [2] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- [3] Dimitris Iakovidis, Dimitris Chatzis, Panos Chrysanthopoulos, and Anastasios Koulaouzidis. Blood detection in wireless capsule endoscope images based on salient superpixels. *volume 2015*, 08 2015.
- [4] Xiao Jia and Max Q.-H. Meng. A deep convolutional neural network for bleeding detection in wireless capsule endoscopy images. *2016 38th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, 2016.
- [5] Y. S. Jung, Y. H. Kim, D. H. Lee, and J. H. Kim. Active blood detection in a high resolution capsule endoscopy using color spectrum transformation. In *2008 International Conference on BioMedical Engineering and Informatics*, volume 1, pages 859–862, May 2008.
- [6] Subodh K Shah, Pragya P Rajauria, Jeongkyu Lee, and M. Emre Celebi. Classification of bleeding images in wireless capsule endoscopy using hsi color domain and region segmentation. 07 2019.
- [7] PÈŽal Halvorsen Thomas de Lange Kristin Ranheim Randel Duc-Tien Dang-Nguyen Mathias Lux-Olga Ostroukhova Konstantin Pogorelov, Michael Riegler. Medico multimedia task at mediaeval 2018. *Media Eval' 2018*, 2018.
- [8] Konstantin Pogorelov, Kristin Ranheim Randel, Carsten Griwodz, Sigrun Losada Eskeland, Thomas de Lange, Dag Johansen, Concetto Spampinato, Duc-Tien Dang-Nguyen, Mathias Lux, Peter Thelin Schmidt, Michael Riegler, and Pál Halvorsen. Kvasir: A multi-class image dataset for computer aided gastrointestinal disease detection. In *Proceedings of the 8th ACM on Multimedia Systems Conference*, MMSys'17, pages 164–169, New York, NY, USA, 2017. ACM.
- [9] Adnan Qayyum, Syed Anwar, Muhammad Majid, Muhammad Awais, and Majdi Alnowami. Medical image analysis using convolutional neural networks: A review. 42, 09 2017.
- [10] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 91–99. Curran Associates, Inc., 2015.
- [11] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Fei-Fei Li. Imagenet large scale visual recognition challenge. *CoRR*, abs/1409.0575, 2014.
- [12] Aoyama Kazuharu Nishikawa Yoshitaka Miura Motoi Shinagawa Takahide Takiyama Hirotoshi Tanimoto Tetsuya Ishihara Soichiro-Matsu Keigo-Tada Tomohiro Shichijo Satoki, Nomura Shuhei. Application of convolutional neural networks in the diagnosis of helicobacter pylori infection based on endoscopic images. *EBioMedicine*, 25:106–111, Nov 2017.
- [13] Pia Smedsrød Trine B. Haugen Kristin Ranheim Randel Konstantin Pogorelov HÅékon Kvale Stensland Duc-Tien Dang-Nguyen Mathias Lux Andreas Petlund Thomas de Lange Peter Thelin Schmidt PÅel Halvorsen Steven Hicks, Michael Riegler. Acm mm biomedia 2019 grand challenge overview. 2019.
- [14] Gregor Urban, Priyam Tripathi, Talal Alkayali, Mohit Mittal, Farid Jalali, William Karnes, and Pierre Baldi. Deep learning localizes and identifies polyps in real time with 96% accuracy in screening colonoscopy. *Gastroenterology*, 155(4), 2018.