

Hierarchical Deep Convolutional Neural Networks for Multi-category Diagnosis of Gastrointestinal Disorders on Histopathological Images

Rasoul Sali¹, Sodiq Adewole¹, Lubaina Ehsan², Lee A. Denson⁴, Paul Kelly^{5,6}, Beatrice C. Amadi⁶, Lori Holtz⁷, Syed Asad Ali⁸, Sean R. Moore², Sana Syed^{2,*}, and Donald E. Brown^{1,3,*}

¹ Department of Systems and Information Engineering, University of Virginia, Charlottesville, VA, USA

² Department of Pediatrics, School of Medicine, University of Virginia, Charlottesville, VA, USA

³ School of Data Science, University of Virginia, Charlottesville, VA, USA

⁴ Division of Gastroenterology, Hepatology, and Nutrition, Cincinnati Children's Hospital Medical Center, Cincinnati, OH, USA

⁵ Blizard Institute, Barts and The London School of Medicine, Queen Mary University of London, London, United Kingdom

⁶ Tropical Gastroenterology and Nutrition group, University of Zambia, School of Medicine, Lusaka, Zambia

⁷ Department of Pediatrics, School of Medicine, Washington University, Saint Louis, MO, USA

⁸ Department of Pediatrics and Child Health, Aga Khan University, Karachi, Pakistan

*co-corresponding authors: {sana.syed, brown}@virginia.edu

Abstract— Deep convolutional neural networks (CNNs) have been successful for a wide range of computer vision tasks, including image classification. A specific area of the application lies in digital pathology for pattern recognition in the tissue-based diagnosis of gastrointestinal (GI) diseases. This domain can utilize CNNs to translate histopathological images into precise diagnostics. This is challenging since these complex biopsies are heterogeneous and require multiple levels of assessment. This is mainly due to structural similarities in different parts of the GI tract and shared features among different gut diseases. Addressing this problem with a flat model that assumes all classes (parts of the gut and their diseases) are equally difficult to distinguish leads to an inadequate assessment of each class. Since the hierarchical model restricts classification error to each sub-class, it leads to a more informative model than a flat model. In this paper, we propose to apply the hierarchical classification of biopsy images from different parts of the GI tract and the receptive diseases within each. We embedded a class hierarchy into the plain VGGNet to take advantage of its layers' hierarchical structure. The proposed model was evaluated using an independent set of image patches from 373 whole slide images. The results indicate that the hierarchical model can achieve better results than the flat model for multi-category diagnosis of GI disorders using histopathological images.

Index Terms—Hierarchical deep convolutional neural network, Gastrointestinal disorders, Multi-category diagnosis, Histopathological images, Coarse categories, Fine classes.

I. INTRODUCTION

Gastrointestinal (GI) diseases are ailments linked to the digestive system, including the esophagus, stomach, and the intestines. GI diseases account for substantial morbidity, mortality, and financial burden by affecting the GI tract and impacting digestion and overall health. The National Institute of Health reports that between 60 and 70 million Americans are affected by GI diseases each year [1].

A common approach to GI disease diagnosis lies in digital pathology for pattern recognition. However, a significant challenge of interpreting clinical biopsy images to diagnose disease is the often striking overlap in histopathology images between distinct but related conditions. There is a critical clinical need to develop new methods to allow clinicians to translate heterogeneous biomedical images into precise diagnostics [2].

Convolutional Neural Networks (CNNs) have shown superior performance for the automated extraction of quantitative morphologic phenotypes from GI biopsy images and associated diseases diagnosis [3], [4]. Despite this, as the number of disease classes associated with different parts of the gut becomes larger, one of the problems that may arise is that visual separability between classes becomes more challenging. This is mainly due to structural similarities in different gut parts and the shared features among gut diseases. Furthermore, for multi-class classification problems, some classes become harder to distinguish than others and require dedicated classifiers [5]. Regular flat models cannot address this issue because they assume that all classes are equally difficult to distinguish [6]. Hierarchical relationships are often identified that exist between classes, which can be used to deploy a hierarchical classification model. This model can be more informative since the classification error is restricted to subcategories compared to treating all classes as arranged in a flat structure.

In CNNs architecture, lower layers capture low-level features while higher layers are likely to extract more abstract features [7]. This CNN property can be combined with the hierarchical structure of classes to enforce the network to learn different levels of class hierarchy in different layers. In this way, coarse categories that are easier to classify are repre-

sented in lower (shallow) layers while higher (deeper) layers output fine subcategories simultaneously [5], [6]. In this paper, we propose a hierarchical deep convolutional neural network to take advantage of GI diseases' hierarchical structure for classification.

This paper is organized as follows: Section II provides an introduction to the diseases studied in this paper. In section III some related researches are reviewed. The methodology is explained in section IV. The data used in this study, data preparation steps and empirical results are elaborated in section V. Finally, section VII concludes the paper along with outlining future directions.

II. GASTROINTESTINAL DISORDERS

GI disorders refer to any abnormal condition or disease that occurs within the GI tract. While there is a wide variety of disorders associated with different parts of the GI tract, this paper focuses on certain disorders involving only the duodenum, esophagus, and ileum. In this section, we give an introduction to each of the considered disorders.

A. Duodenum

1) *Celiac Disease (CD)*: It is an inability to normally process dietary gluten (present in foods such as wheat, rye, and barley) and is present in 1% of the US population. Gluten exposure triggers an inflammatory cascade, which leads to a compromised intestinal barrier function. Gluten consumption by people with CD can cause diarrhea, abdominal pain, bloating, and weight loss. If unrecognized, it can lead to anemia, decreased bone density, and, in longstanding cases, intestinal cancer [8].

2) *Environmental Enteropathy (EE)*: It is an acquired small intestinal condition resulting from the continuous burden of immune stimulation by fecal-oral exposure to enteropathogens leading to a persistent acute phase response and chronic inflammation [9], [10]. EE can be characterized histologically by villus shortening, crypt hyperplasia, and a resultant decrease in the surface area of mature absorptive intestinal epithelial cells, leading to a markedly reduced nutrient absorption, under-nutrition, and stunting [11].

B. Esophagus

1) *Eosinophilic Esophagitis (EoE)*: It is a chronic, allergic inflammatory disease of the esophagus. It occurs when eosinophils, a normal type of white blood cells present in the digestive tract, build up in the lining of the esophagus. EoE is characterized by symptoms of esophageal dysfunction and eosinophilic infiltration of the esophageal mucosa in the absence of secondary causes of eosinophilia [12].

C. Ileum

1) *Crohn's Disease*: It is an inflammatory bowel disease that causes patchy disease constituting of chronic inflammation, ulcers, and mucosal damage anywhere in the GI tract, although the most common being the terminal ileum and colon. The interaction of genetic susceptibility, environmental

factors, and intestinal microflora is believed to be the major cause of Crohn's disease. This interaction results in abnormal mucosal immune response, which compromises epithelial barrier function [13].

III. RELATED WORK

Hierarchical CNN has demonstrated improved performance in image classification compared to flat CNN models across multiple domains [14]–[17]. These models exploit the hierarchical structure of object categories [18] to decompose the classification tasks into multiple steps. Hierarchical Deep CNNs (HD-CNN) proposed by Yan et al. [5] embeds CNN into a categorical hierarchy by separating easy classes using a coarse category classifier and difficult categories using a fine category classifier. This model can be implemented without increasing the complexity of the training process; however, it requires multi-step training of each CNN. Zhu and Bain proposed a branched variant (B-CNN) of the hierarchical deep CNN [6]. Since shallow layers of a CNN capture low-level features while deeper levels capture high-level features, B-CNN outputs multiple predictions ordered from coarse to fine along concatenated convolutional layers corresponding to the hierarchical structure of the target classes. The model branch training strategy adjusts parameters on the output layers, forcing the input to learn successively coarse to fine concepts along with the layer blocks. Hierarchical architecture has been applied to both image [15] and video classification [19] tasks with superior performance compared with conventional flat CNN models. While deep learning has seen a significant application in medical image classification tasks [16], hierarchical models remain a relatively less explored area in literature.

Ranjan et al. reported a CNN-based hierarchical model in medical image classification on histopathological images. [17] with superior performance than flat CNN. They proposed to classify cancer and its states (in situ, invasive, or normal) using multiple CNNs organized hierarchically. With one CNN in each of the two-level classification tasks, the first level CNN, a pre-trained AlexNet, is trained to discriminate the normal class from the rest of the classes' images. The second level of the hierarchy is a tree of CNN-based binary classifiers using majority voting to discriminate the other three classes; in situ, invasive and benign cells.

Krauß et al. applied a hierarchical CNN on cytopathology to classify cellular images as healthy or cancer-affected cells [20]. To the best of our knowledge, there are no previously published studies that have applied hierarchical deep convolutional neural networks to gastrointestinal disease classification using histopathological images.

IV. METHODOLOGY

A. Base Model

There are many different architectures of CNNs in the literature with associated advantages and drawbacks. In this paper, we used VGGNet [21] (proposed by Visual Geometry Group in the University of Oxford) as a base model that has shown excellent performance in image classification problems,

including medical image analysis [22]–[24]. VGGNet obtained the state-of-the-art results in the ILSVRC’14 competition with 7.3% error rate, which was among the top 5 errors and was a significant improvement over ZFNet [7], the winner of ILSVRC’13. Two main intertwined characteristics of VGGNet were the increased depth of the network and applying smaller filters. It uses 3×3 sized filters and 2×2 sized pooling from the beginning to the end of the network. Since smaller filters have few parameters, it is possible to increase the depth by stacking more of them with the same effective receptive fields when using larger filters. For instance, effective receptive fields of three stacked 3×3 filters with stride 1 are the same as a 7×7 filter. VGG16 and VGG19 have been released as two variants of VGGNet. VGG16 has 16 trainable layers, including 13 convolutional layers organized in 5 blocks followed by 3 fully-connected layers. The final layer is a softmax layer that outputs class probabilities. In this paper, VGG16 was applied and was trained from scratch on biopsy patches. The network is shown at the top in Figure 1 is a plain VGG16.

B. Hierarchical Convolutional Neural Network

To propose a hierarchical convolutional neural network, the architecture of Branch Convolutional Neural Network (B-CNN) [6] proposed by Zhu and Bain was applied to embed different levels of class hierarchical on VGG16 to propose H-VGGNet. There were seven classes of GI disorders: Duodenum-Celiac, Duodenum-EE, Duodenum-Normal, Esophagus-EoE, Esophagus-Normal, Ileum-Crohn’s and Ileum-Normal, as fine classes and each class belonged to a coarse category: Duodenum, Esophagus or Ileum (see the hierarchy of classes in Figure 2). Since there are two class levels in this hierarchy, in addition to output layer, which will output fine classes, one branch was added to VGG16 to output the coarse categories. The network is shown at the bottom in Figure 1 is H-VGGNet. Although new branches can consist of both convolutional and fully-connected layers, in this paper, the new added branch was composed of only fully connected layers. This model for each input image computed both coarse and fine level predictions. The rectified linear unit (ReLU) [25] was employed as the activation function. To reduce over-fitting, dropout regularization [26] was used after ReLU of each fully-connected layer with $p = 0.5$. Also, Batch Normalization [27] was applied after ReLU of every trainable layer.

The loss function for such a model is weighted summation of coarse and fine prediction losses (see equation 1).

$$\mathcal{L}_i = - \sum_{k=1}^K w_k \log \left(\frac{e^{s_{y_i p}^k}}{\sum_j e^{s_{y_i j}^k}} \right) \quad (1)$$

Where K is number of levels in hierarchy of classes, w_k is weight of level k th in class hierarchy and term $-\log(\cdot)$ is cross-entropy loss function for the i th instance in k th level of class hierarchy. $s_{y_i p}^k$ is the element in class score of instance i th in k th level of class hierarchy corresponding to positive

element in target label y_i and $s_{y_i j}^k$ is the j th element in class score of instance i th in k th level of class hierarchy.

V. EXPERIMENTAL SETUP

This section is devoted to presenting the experimental setting, including data description, data pre-processing steps, training details, and the evaluation criterion.

A. Data

1,150 Hematoxylin and Eosin (H&E) stained whole slide images (WSI) from 441 patients were obtained for this study. Some patients had 2 – 3 biopsies for each diagnosis. All biopsies obtained from the University of Virginia (UVA), VA, USA were retrospectively retrieved archival samples while the other biopsies were obtained as part of prospective cohorts studying growth faltering among children (except Crohn’s Disease biopsies). Images were obtained based on the gut disease state: 1) Celiac Disease in the Duodenum: UVA ($n = 239$), Cincinnati Children’s Medical Center (CCHMC), OH, USA and Washington University (WashU), MO, USA ($n = 43$ and $n = 54$, respectively); 2) Environmental Enteropathy in Duodenum: Aga Khan University, Karachi, Pakistan ($n = 34$) and Zambia School of Medicine, Lusaka, Zambia ($n = 19$); 3) histologically normal duodenum: UVA ($n = 174$), CCHMC ($n = 36$), WashU ($n = 11$); 4) EoE in Esophagus: UVA ($n = 349$); 5) histologically normal esophagus: UVA ($n = 155$); 6) Crohn’s disease in Ileum: CCHMC as part of RISK study sub-cohort ($n = 20$); and, 7) histologically normal ileum: UVA ($n = 16$).

We split our data into training, development, and test sets using 50 : 20 : 30 ratio. Since the model must generalize to other unseen patients data, we performed our split to ensure no overlap between the training, development, and test set for a particular patient.

B. Data Pre-processing

1) *Image Patching*: A sliding window method was applied to each high-resolution WSI to generate patches of size 1000×1000 pixels. Since some classes had more whole-slide images than others, we generated patches with different overlapping areas for each class. Patches were resized to 224×224 pixels to reduce the computational cost. After generating tissue tiles from each WSI, tiles’ labels were assumed to be the same with its associated WSI.

2) *Patch Clustering*: In this work, a two-step clustering process was applied to filter useless patches which had mostly been created from the WSIs’ background. A convolutional auto-encoder (CAE) was employed to map each patch into an embedding space through the first step. In the second step, k-means clustering was applied to cluster embedded features into two clusters: useful and useless. Table I summarizes distribution of WSIs and patches (after cleaning) in each class.

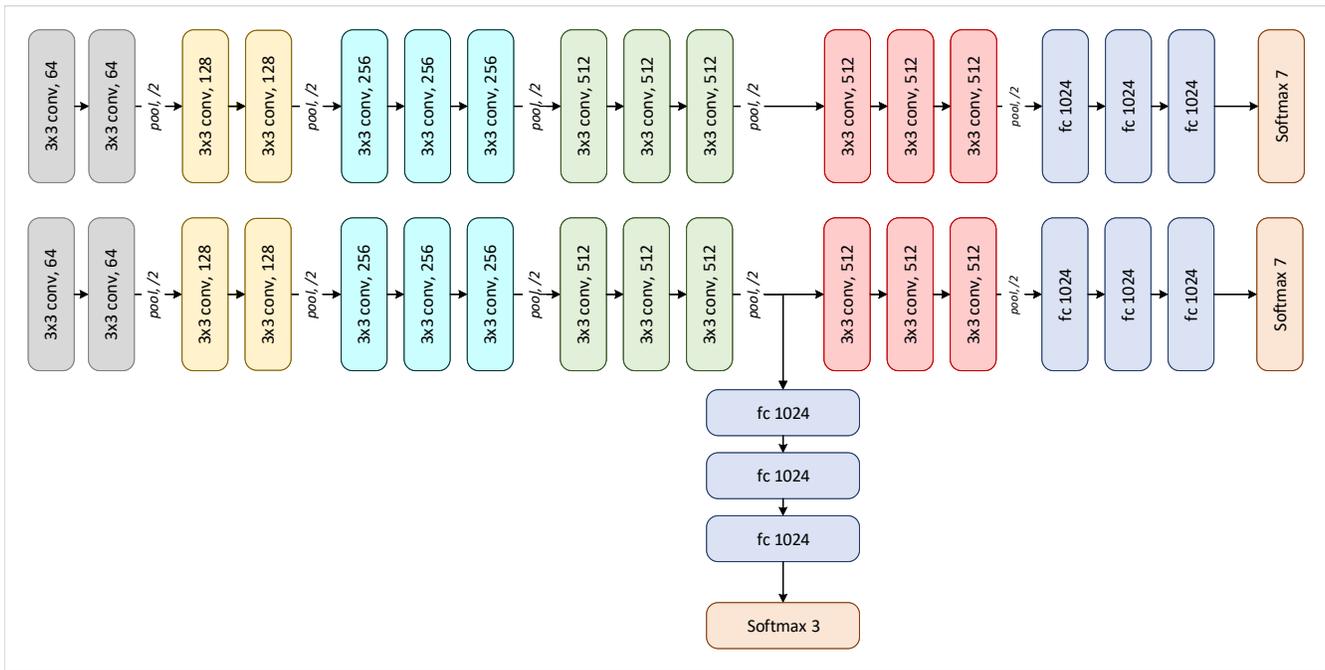


Fig. 1: Top: VGGNet architecture, Bottom: H-VGGNet architecture

TABLE I: Distribution of training, development, and test set data among different classes

Coarse category	Fine class	Train		Development		Test	
		Number of WSIs	Number of patches	Number of WSIs	Number of patches	Number of WSIs	Number of patches
Duodenum	Celiac	170	8521	62	1548	104	2738
	EE	28	8138	11	1058	14	1343
	Normal	120	8290	36	810	65	1257
Esophagus	EoE	140	12342	75	5389	134	10907
	Normal	80	9115	30	3170	45	5202
Ileum	Crohn's	11	5006	3	1300	6	1953
	Normal	8	3906	3	1312	5	2182

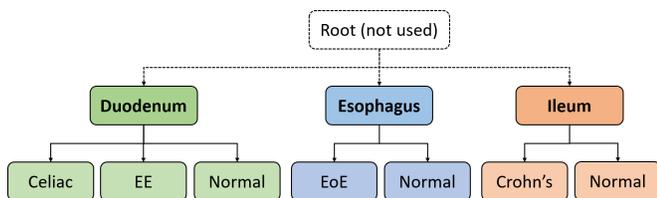


Fig. 2: Hierarchy of classes

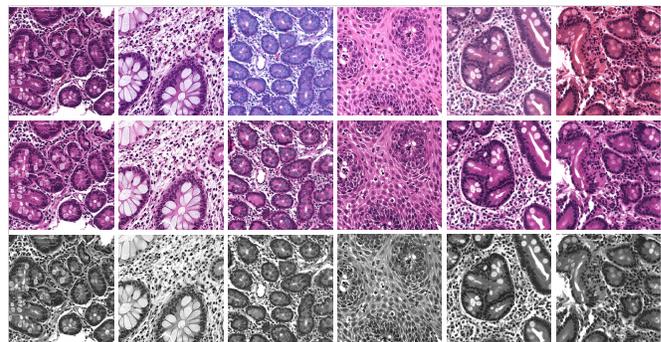


Fig. 3: Color normalization artifacts. The first row represents the original images, the second row is color normalized images using the method proposed by Vahadane et al. [28] and their associated gray-scale images are in the third row.

3) *Stain Color Normalization*: Histological images have substantial color variation that adds bias while training the model. This issue arises due to a wide variety of factors, such as differences in raw materials and manufacturing techniques of stain vendors, staining protocols of labs, and color responses of digital scanners [28]. Unwanted color variations should be addressed and resolved as an essential pre-processing step before any analyses to prevent any bias arising from this issue.

Various solutions such as color balancing [4], gray-scale, and stain normalization [3] have been proposed in the published literature to address the color variation issue. In this study, we used gray-scale images. However, before convert-

ing the RGB patches to gray-scale, the stain normalization approach proposed by Vahadane et al. [28] was applied to make sure that the effect of variation of color variation is significantly reduced. Figure 3 shows an example of the result of applying this process on representative biopsy patches.

C. Training Details

We conducted extensive experiments to compare the performance of the hierarchical model with a flat model. Both the base model and the hierarchical model were trained and tested ten times. Each time both the models were trained in 20 epochs. Optimization was performed using RMSprop [29] optimization with no momentum. The initial value of the learning rate is considered as 1×10^{-3} , it changed to 5×10^{-4} after the 10th epoch and to 1×10^{-4} after the 15th epoch. Different loss weights are applied to each level of the hierarchy to reflect the differences in each level of classes' importance. Since the low-level feature extraction is more important in initial epochs, more weights are assigned to it. As the model's training progresses, the weight of the coarse categories level decreases, and the weight of fine classes increases. The changes in loss weights follow [0.98, 0.02] in the first epoch, [0.30, 0.70] in the 5th epoch, [0.10, 0.90] in the 10th epoch, [0.00, 1.00] in the 15th epoch. This change in weights causes the algorithm to focus first on the optimization of the coarse category, and as the learning process progresses, this focus shifts to the fine level.

D. Evaluation Metrics

In order to assess the performance of models, accuracy, area under the ROC curve (AUC), Precision, Recall, and F1 score have been considered.

VI. RESULTS

Table II presents the performance comparison between two models in terms of the accuracy, AUC, Precision, Recall, and F1 score on the test set with 95% confidence intervals. As shown, the hierarchical model's performance for many classes was better than the flat model in terms of the mean of the criterion mentioned above. Also Table III presents the normalized confusion matrix of two models. The confusion between different coarse categories in the hierarchical model was less than the flat model.

VII. CONCLUSION

This paper proposes a hierarchical deep convolutional neural network for multi-category classification of GI disorders using histopathological biopsy images. Our proposed model was tested on 25,582 cropped images derived from an independent set of 373 WSIs. Our results showed that the hierarchical model had superior classification performance for a problem with an inherent hierarchical structure compared to a flat model, which assumes equal difficulty for classification. With the dataset collected from $n = 441$ patients and based on our training, development, and test set split, our model can be generalized to other patients that are not part of the training or development sets.

In CNN architecture, since lower layers capture low-level features while higher layers are likely to extract more abstract features, we utilized this property to build our model instead of employing separate models for different class hierarchy levels. The use of such a structure makes it possible to save

computational cost and benefit from shared information across the coarse levels in the training phase. Quantification of this synergy could be a possible avenue for future research.

ACKNOWLEDGEMENTS

The research reported in this manuscript was supported by the National Institute of Diabetes and Digestive and Kidney Diseases of the National Institutes of Health under award number K23DK117061-01A1 (SS), Bill and Melinda Gates Foundation under award numbers OPP1066203, OPP1066118, OPP1144149 and OPP1066153 and University of Virginia Translational Health Research Institute of Virginia (THRIV) Scholar Career Development Award (SS). The content is solely the authors' responsibility and does not necessarily represent the official views of the funding agencies.

REFERENCES

- [1] A. F. Peery, E. S. Dellon, J. Lund, S. D. Crockett, C. E. McGowan, W. J. Bulsiewicz, L. M. Gangarosa, M. T. Thiny, K. Stizenberg, D. R. Morgan *et al.*, "Burden of gastrointestinal disease in the united states: 2012 update," *Gastroenterology*, vol. 143, no. 5, pp. 1179–1187, 2012.
- [2] B. E. Bejnordi, M. Veta, P. J. Van Diest, B. Van Ginneken, N. Karssemeijer, G. Litjens, J. A. Van Der Laak, M. Hermsen, Q. F. Manson, M. Balkenhol *et al.*, "Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer," *Jama*, vol. 318, no. 22, pp. 2199–2210, 2017.
- [3] R. Sali, L. Ehsan, K. Kowsari, M. Khan, C. A. Moskaluk, S. Syed, and D. E. Brown, "Celiacnet: Celiac disease severity diagnosis on duodenal histopathological images using deep residual networks," *arXiv preprint arXiv:1910.03084*, 2019.
- [4] K. Kowsari, R. Sali, M. N. Khan, W. Adorno, S. A. Ali, S. R. Moore, B. C. Amadi, P. Kelly, S. Syed, and D. E. Brown, "Diagnosis of celiac disease and environmental enteropathy on biopsy images using color balancing on convolutional neural networks," in *Proceedings of the Future Technologies Conference*. Springer, 2019, pp. 750–765.
- [5] Z. Yan, H. Zhang, R. Piramuthu, V. Jagadeesh, D. DeCoste, W. Di, and Y. Yu, "Hd-cnn: hierarchical deep convolutional neural networks for large scale visual recognition," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2740–2748.
- [6] X. Zhu and M. Bain, "B-cnn: branch convolutional neural network for hierarchical classification," *arXiv preprint arXiv:1709.09890*, 2017.
- [7] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *European conference on computer vision*. Springer, 2014, pp. 818–833.
- [8] I. Parzanese, D. Qehajaj, F. Patricicola, M. Aralica, M. Chiriva-Internati, S. Stifter, L. Elli, and F. Grizzi, "Celiac disease: From pathophysiology to treatment," *World journal of gastrointestinal pathophysiology*, vol. 8, no. 2, p. 27, 2017.
- [9] D. Campbell, M. Elia, and P. Lunn, "Growth faltering in rural gambian infants is associated with impaired small intestinal barrier function, leading to endotoxemia and systemic inflammation," *The Journal of nutrition*, vol. 133, no. 5, pp. 1332–1338, 2003.
- [10] N. W. Solomons, "Environmental contamination and chronic inflammation influence human growth potential," *The Journal of nutrition*, vol. 133, no. 5, pp. 1237–1237, 2003.
- [11] S. Syed, A. Ali, and C. Duggan, "Environmental enteric dysfunction in children: a review," *Journal of pediatric gastroenterology and nutrition*, vol. 63, no. 1, p. 6, 2016.
- [12] E. S. Dellon and I. Hirano, "Epidemiology and natural history of eosinophilic esophagitis," *Gastroenterology*, vol. 154, no. 2, pp. 319–332, 2018.
- [13] J. Torres, S. Mehandru, J.-F. Colombel, and L. Peyrin-Biroulet, "Crohn's disease," *The Lancet*, vol. 389, no. 10080, pp. 1741–1755, 2017.
- [14] Z. Yan, R. Piramuthu, V. Jagadeesh, W. Di, and D. Decoste, "Hierarchical deep convolutional neural network for image classification," Aug. 20 2019, uS Patent 10,387,773.

TABLE II: Comparison of models' performance

Metric	Model	Class						
		Duodenum			Esophagus		Ileum	
		Celiac	EE	Normal	EoE	Normal	Crohn's	Normal
Accuracy	VGGNet	0.941 ± 0.011	0.986 ± 0.008	0.959 ± 0.009	0.917 ± 0.015	0.928 ± 0.013	0.960 ± 0.011	0.969 ± 0.012
	H-VGGNet	0.953 ± 0.007	0.987 ± 0.003	0.963 ± 0.008	0.937 ± 0.012	0.945 ± 0.013	0.973 ± 0.009	0.971 ± 0.010
AUC	VGGNet	0.842 ± 0.039	0.981 ± 0.011	0.843 ± 0.032	0.906 ± 0.016	0.923 ± 0.008	0.947 ± 0.024	0.863 ± 0.053
	H-VGGNet	0.870 ± 0.019	0.992 ± 0.003	0.872 ± 0.025	0.938 ± 0.005	0.931 ± 0.008	0.967 ± 0.008	0.874 ± 0.022
Precision	VGGNet	0.763 ± 0.052	0.837 ± 0.099	0.625 ± 0.070	0.965 ± 0.007	0.775 ± 0.046	0.710 ± 0.089	0.961 ± 0.024
	H-VGGNet	0.850 ± 0.020	0.915 ± 0.018	0.682 ± 0.053	0.942 ± 0.013	0.856 ± 0.025	0.798 ± 0.038	0.937 ± 0.013
Recall	VGGNet	0.712 ± 0.083	0.975 ± 0.025	0.711 ± 0.071	0.834 ± 0.035	0.915 ± 0.013	0.928 ± 0.049	0.729 ± 0.107
	H-VGGNet	0.756 ± 0.041	0.989 ± 0.005	0.764 ± 0.055	0.920 ± 0.017	0.901 ± 0.023	0.955 ± 0.019	0.752 ± 0.045
F1 score	VGGNet	0.728 ± 0.039	0.893 ± 0.058	0.653 ± 0.023	0.894 ± 0.019	0.838 ± 0.025	0.797 ± 0.056	0.820 ± 0.033
	H-VGGNet	0.799 ± 0.020	0.950 ± 0.008	0.714 ± 0.015	0.930 ± 0.006	0.877 ± 0.006	0.868 ± 0.021	0.833 ± 0.027

TABLE III: Normalized confusion matrix of flat and hierarchical model

True Label	Model	Predicted Label						
		Duodenum			Esophagus		Ileum	
		Celiac	EE	Normal	EoE	Normal	Crohn's	Normal
Celiac	VGGNet	0.712 ± 0.083	0.057 ± 0.046	0.162 ± 0.058	0.001 ± 0.001	0.001 ± 0.001	0.064 ± 0.025	0.004 ± 0.004
	H-VGGNet	0.756 ± 0.051	0.026 ± 0.008	0.157 ± 0.046	0.015 ± 0.015	0.004 ± 0.004	0.032 ± 0.008	0.011 ± 0.004
EE	VGGNet	0.020 ± 0.020	0.975 ± 0.025	0.001 ± 0.001	0.001 ± 0.001	0.001 ± 0.001	0.004 ± 0.004	0.000 ± 0.000
	H-VGGNet	0.009 ± 0.005	0.989 ± 0.005	0.001 ± 0.001	0.000 ± 0.000	0.001 ± 0.001	0.001 ± 0.001	0.000 ± 0.000
Normal (Duodenum)	VGGNet	0.218 ± 0.088	0.035 ± 0.035	0.711 ± 0.071	0.001 ± 0.001	0.001 ± 0.001	0.030 ± 0.022	0.007 ± 0.007
	H-VGGNet	0.197 ± 0.056	0.011 ± 0.006	0.764 ± 0.055	0.007 ± 0.006	0.004 ± 0.003	0.007 ± 0.002	0.011 ± 0.004
EoE	VGGNet	0.016 ± 0.006	0.004 ± 0.002	0.003 ± 0.003	0.834 ± 0.035	0.123 ± 0.033	0.017 ± 0.008	0.004 ± 0.002
	H-VGGNet	0.003 ± 0.001	0.003 ± 0.001	0.001 ± 0.001	0.920 ± 0.017	0.066 ± 0.015	0.003 ± 0.003	0.004 ± 0.001
Normal (Esophagus)	VGGNet	0.009 ± 0.003	0.002 ± 0.002	0.003 ± 0.003	0.053 ± 0.013	0.915 ± 0.013	0.018 ± 0.008	0.001 ± 0.001
	H-VGGNet	0.003 ± 0.001	0.000 ± 0.000	0.000 ± 0.000	0.086 ± 0.022	0.901 ± 0.023	0.007 ± 0.003	0.003 ± 0.001
Crohn's	VGGNet	0.042 ± 0.031	0.008 ± 0.008	0.017 ± 0.017	0.002 ± 0.001	0.001 ± 0.001	0.928 ± 0.049	0.004 ± 0.003
	H-VGGNet	0.019 ± 0.011	0.002 ± 0.001	0.009 ± 0.006	0.007 ± 0.003	0.003 ± 0.003	0.955 ± 0.019	0.006 ± 0.004
Normal (Ileum)	VGGNet	0.018 ± 0.007	0.019 ± 0.019	0.034 ± 0.031	0.025 ± 0.018	0.031 ± 0.015	0.144 ± 0.117	0.729 ± 0.093
	H-VGGNet	0.011 ± 0.004	0.004 ± 0.002	0.007 ± 0.005	0.057 ± 0.022	0.025 ± 0.007	0.144 ± 0.057	0.752 ± 0.045

- [15] Y. Seo and K.-s. Shin, "Hierarchical convolutional neural networks for fashion image classification," *Expert Systems with Applications*, vol. 116, pp. 328–339, 2019.
- [16] F. Milletari, N. Navab, and S.-A. Ahmadi, "V-net: Fully convolutional neural networks for volumetric medical image segmentation," in *2016 Fourth International Conference on 3D Vision (3DV)*. IEEE, 2016, pp. 565–571.
- [17] N. Ranjan, P. V. Machingal, S. S. D. Jammalmadka, V. Thenaknidiyoor, and A. Dileep, "Hierarchical approach for breast cancer histopathology images classification," 2018.
- [18] A.-M. Tousch, S. Herbin, and J.-Y. Audibert, "Semantic hierarchies for image annotation: A survey," *Pattern Recognition*, vol. 45, no. 1, pp. 333–345, 2012.
- [19] J. Fan, A. K. Elmagarmid, X. Zhu, W. G. Aref, and L. Wu, "Classview: hierarchical video shot classification, indexing, and accessing," *IEEE Transactions on Multimedia*, vol. 6, no. 1, pp. 70–86, 2004.
- [20] S. D. Krauß, R. Roy, H. K. Yosef, T. Lehtonen, S. F. El-Mashtoly, K. Gerwert, and A. Mosig, "Hierarchical deep convolutional neural networks combine spectral and spatial information for highly accurate raman-microscopy-based cytopathology," *Journal of biophotonics*, vol. 11, no. 10, p. e201800022, 2018.
- [21] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [22] A. Rakhlin, A. Shvets, V. Iglovikov, and A. A. Kalinin, "Deep convolutional neural networks for breast cancer histology image analysis," in *International Conference Image Analysis and Recognition*. Springer, 2018, pp. 737–744.
- [23] J. J. Gómez-Valverde, A. Antón, G. Fatti, B. Liefers, A. Herranz, A. Santos, C. I. Sánchez, and M. J. Ledesma-Carbayo, "Automatic glaucoma classification using color fundus images based on convolutional neural networks and transfer learning," *Biomedical optics express*, vol. 10, no. 2, pp. 892–913, 2019.
- [24] S. Y. Ko, J. H. Lee, J. H. Yoon, H. Na, E. Hong, K. Han, I. Jung, E.-K. Kim, H. J. Moon, V. Y. Park *et al.*, "Deep convolutional neural network for the diagnosis of thyroid nodules on ultrasound," *Head & neck*, vol. 41, no. 4, pp. 885–891, 2019.
- [25] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [26] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *The journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [27] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.
- [28] A. Vahadane, T. Peng, A. Sethi, S. Albarqouni, L. Wang, M. Baust, K. Steiger, A. M. Schlitter, I. Esposito, and N. Navab, "Structure-preserving color normalization and sparse stain separation for histological images," *IEEE transactions on medical imaging*, vol. 35, no. 8, pp. 1962–1971, 2016.
- [29] T. Tieleman and G. Hinton, "Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude," *COURSERA: Neural networks for machine learning*, vol. 4, no. 2, pp. 26–31, 2012.