

# Automated Histopathologic Assessment of Hirschsprung Disease Using a Multi-Stage Vision Transformer Framework

Youssef Megahed<sup>a,d,\*</sup>, Saleh Abou-Alwan<sup>a,\*</sup>, Anthony Fuller<sup>a</sup>, Dina El Demellawy<sup>e</sup>, Steven Hawken<sup>a,b,c,d,\*\*</sup> and Adrian D. C. Chan<sup>a,\*\*</sup>

<sup>a</sup>Department of Systems and Computer Engineering, Carleton University, Ottawa, Ontario, Canada

<sup>b</sup>Department of Clinical Science and Translational Medicine, University of Ottawa, Ottawa, Ontario, Canada

<sup>c</sup>School of Epidemiology and Public Health, University of Ottawa, Ottawa, Ontario, Canada

<sup>d</sup>Department of Methodological and Implementation Research, Ottawa Hospital Research Institute, Ottawa, Ontario, Canada

<sup>e</sup>Children's Hospital of Eastern Ontario (CHEO), Ottawa, Ontario, Canada

## ARTICLE INFO

### Keywords:

Hirschsprung Disease  
Vision Transformer  
Myenteric Plexus  
Ganglion Cell Detection  
Digital Pathology

## ABSTRACT

Hirschsprung Disease is characterized by the absence of ganglion cells in the myenteric plexus. Therefore, their correct identification is crucial for diagnosing Hirschsprung disease. We introduce a three-stage segmentation framework based on a Vision Transformer (ViT-B/16) that mimics the pathologist's diagnostic approach. The framework sequentially segments the muscularis propria, delineates the myenteric plexus, and identifies ganglion cells within anatomically valid regions. 30 whole-slide images of colon tissue were used, each containing expert manual annotations of muscularis, plexus, and ganglion cells at varying levels of certainty. A 5-fold cross-validation scheme was applied to each stage, along with resolution-specific tiling strategies and tailored postprocessing to ensure anatomical consistency. The proposed method achieved a Dice coefficient of 89.9% and a Plexus Inclusion Rate of 100% for muscularis segmentation. Plexus segmentation reached a recall of 94.8%, a precision of 84.2% and a Ganglia Inclusion Rate of 99.7%. For high-certainty ganglion cells, the model achieved 62.1% precision and 89.1% recall, while joint certainty scores yielded 67.0% precision. These results indicate that ViT-based models are effective at leveraging global tissue context and capturing cellular morphology at small scales, even within complex histological tissue structures. This multi-stage methodology has great potential to support digital pathology workflows by reducing inter-observer variability and assisting in the evaluation of Hirschsprung disease. The clinical impact will be evaluated in future work with larger multi-center datasets and additional expert annotations.

## 1. Introduction

Hirschsprung disease (HD), a congenital birth defect disorder defined by the absence of ganglionic cells within the myenteric plexus of the intestinal tract [7], rapidly manifesting after birth, affects 1/5000 neonatal patients worldwide [1]. Lack of ganglionic cells, caused by defective migration, proliferation, and differentiation of neural crest cells during developmental periods throughout gestation. Thus, leading to obstruction of bowel movements caused by malformation of colon nerve cells and therefore leading to symptoms of HD, including severe constipation and signs of intestinal obstruction [2]. Thus, accurate identification of ganglionic cells in the colon is vital for confirmation of HD.

After the onset of symptoms, traditional methods for diagnosing HD include a contrast enema radiograph, in which a barium contrast is inserted into the child's rectum and x-ray images of the colon are taken to visualize the structure of different colon segments, thereby confirming the presence of malformed intestines. A second popular method for diagnosing HD is a rectal biopsy, in which a section of tissue is extracted from the rectum and examined on a Whole Slide Image (WSI) by a trained pathologist. This

method involves quantitatively assessing the presence of ganglionic cells. Although both are typically considered the gold standard for confirming HD, they have inherent limitations, such as sampling error, inter- and intra-rater variability among pathologists [8, 10], and excessive time consumption. Highlighting the importance of automating the identification of ganglionic cells to confirm HD diagnosis.

The current advances in computational imaging and WSI have enabled Artificial Intelligence (AI) methods to automate the classification, detection, and segmentation of medical image data across various organ systems. For instance, several Deep Learning (DL) algorithms have been utilized for computer-aided diagnosis, attempting to classify breast cancer lesions [9] using WSIs [3] and human aortic and cardiac regions [5]. Other algorithms [4] have been deployed to segment images of bronchoscopes, further identifying lung-related diseases. Furthermore, object segmentation has also been widely applied to various medical images and diseases. Specifically, classifying the potential onset of autism spectrum disorder (ASD), including additional indicators of ASD and irregular brain structure [24, 25, 30]. Leveraging DL computer vision algorithms for identifying the presence of ganglion cells does seem promising. However, segmentation is required because ganglion cells are dispersed across the myenteric plexus layer, which lies within the muscularis

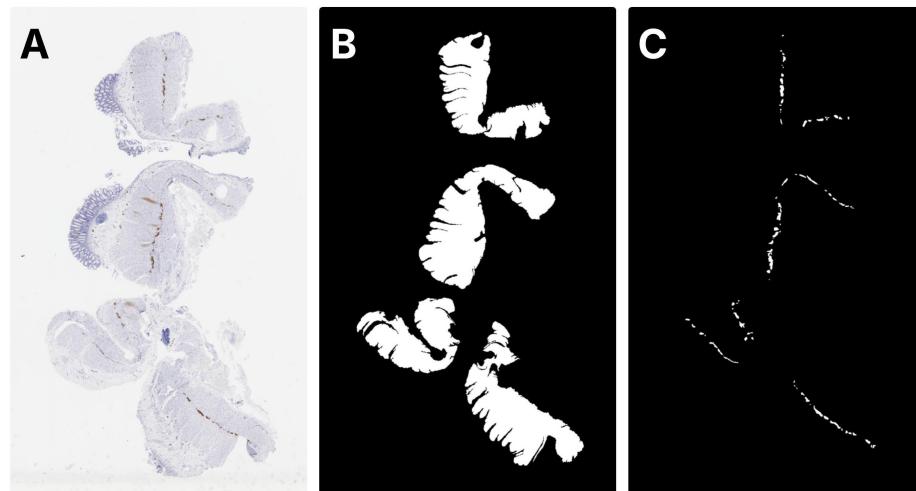
\*Co-first author

\*\*Co-advising author

✉ youssefmegahed@cmail.carleton.ca (Y. Megahed);

adrianchan@cunet.carleton.ca (A.D.C. Chan)

ORCID(s): 0009-0004-2595-5468 (Y. Megahed)



**Figure 1:** Reference annotations for tissue compartments in a whole slide section. **A:** Original WSI. **B:** Ground truth mask delineating the muscularis propria layer. **C:** Ground truth mask marking the myenteric plexus regions.

layer of the colon. Segmenting the different colon layers enables later identification of ganglion cells, or their absence, thereby further confirming the presence of HD.

Recent attempts [11, 12, 13, 14, 15] have been made to segment the muscularis layer, which encloses the myenteric plexus, first utilizing a shallow machine learning model, K-means clustering [14]. This was then followed by deploying a Convolutional Neural Network (CNN), which produced a far greater accuracy and plexus inclusion rate when segmenting the muscularis layer [13, 14]. In [14], a colour-thresholding algorithm was then utilized for segmenting the myenteric plexus region and effectively detecting Carletinin-positive ganglia inside plexus regions. Finally, ganglia were segmented from each plexus region using intensity thresholding on the previously segmented Calretinin signal. This was followed by deploying a Linear Discriminant Analysis (LDA) model to identify true-ganglia from false-ganglia [14]. Overall, the results suggest strong discriminative ability, offering a promising avenue for automated DL. Furthermore, other attempts have been made [6] that leverage the U-Net, a CNN algorithm, to similarly detect the presence of ganglion cells in the myenteric plexus, using WSIs from biopsy tissues. M. Duci et. al., in [6], demonstrated superior innate capabilities of the CNN architecture, especially given its ability to extract low-level and high-level features throughout different layers, in addition to capturing local relationships through convolutions. Overall, the algorithm is able to learn rich, essential semantic features from the colon layers in the WSI dataset. Although U-Net is appropriate for learning pixel-wise local representations, it cannot learn long-range, global contextual relationships across different regions of the image, which could be critical for detecting ganglion cells dispersed throughout the WSI.

Given its success in Natural Language Processing (NLP) tasks, Vision Transformer (ViT) [17] have achieved superior performance to CNNs in computer vision-related tasks, when utilized on general image datasets in addition

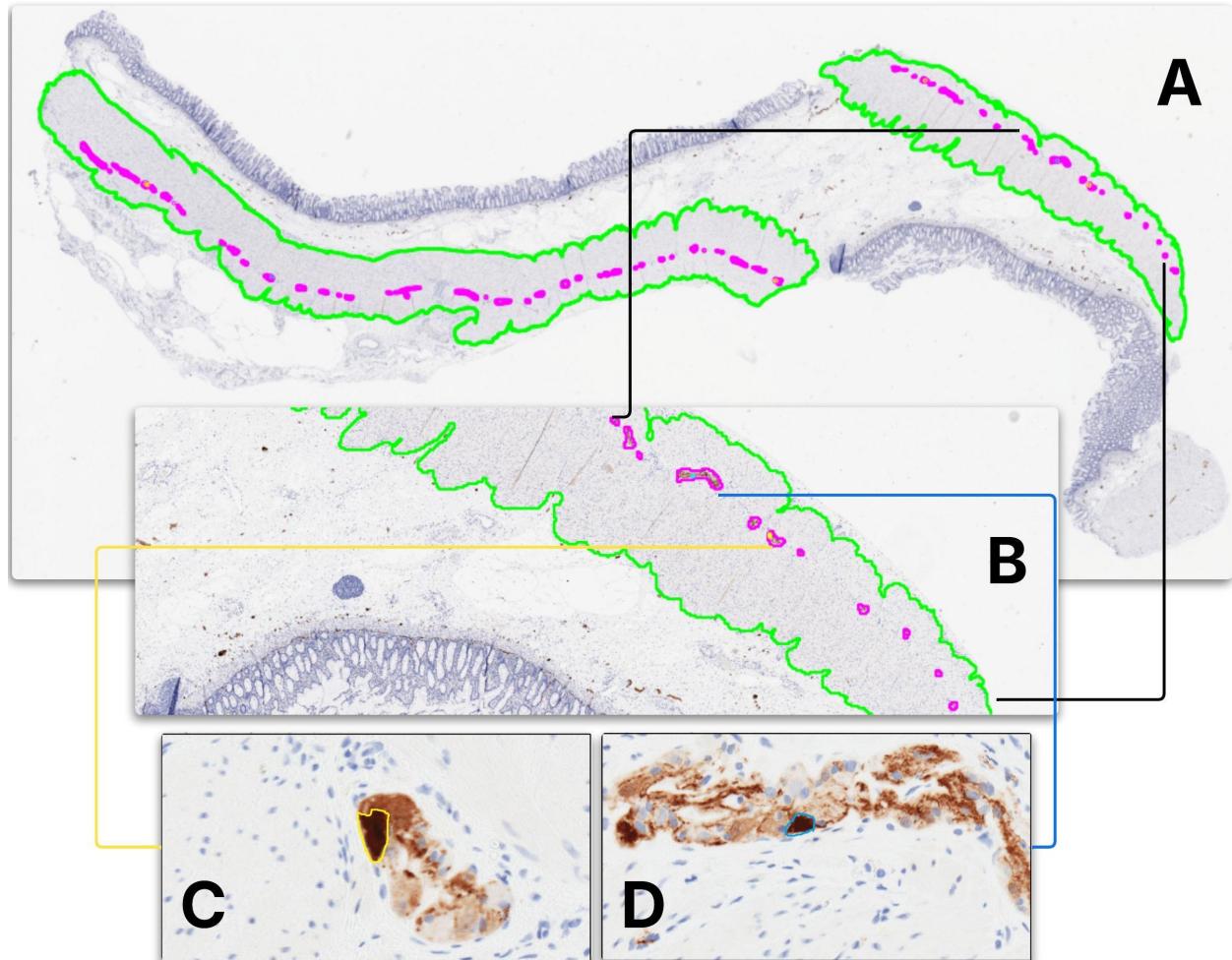
to medical image data [5], as shown in the myenteric plexus classification study in [16]. This is possible because ViTs have inherent self-attention, which enables them to model long-range dependencies and capture global contextual relationships across all parts of an image. Thereby, recognizing patterns that span large areas, such as tissue structure in WSI histopathology.

In this study, we propose a hierarchical segmentation pipeline utilizing ViTs to segment the muscularis layer, which encloses the myenteric plexus. This is followed by segmenting the myenteric plexus layer and finally identifying the presence of ganglion cells, typically found in the myenteric plexus, to further confirm the onset of HD. The study used 30 WSIs of colon sections, manually annotated by a pediatric pathologist. The primary metrics used for evaluating segmentation performance of the muscularis include Dice-Sørensen coefficient (Dice coefficient), precision, recall, and Plexus Inclusion Rate (PIR), ultimately assessing the similarity between the ground truth dataset and the prediction produced. For evaluating the myenteric plexus, evaluation metrics include precision, recall, and Ganglia Inclusion Rate (GIR). Finally, considering the presence of ganglion cells, it was also necessary to assess recall and precision.

## 2. Methodology

### 2.1. Dataset

The dataset used in this research study originates from the Children's Hospital of Eastern Ontario (CHEO), comprising 30 WSIs corresponding to colon sections from 26 different patients, all of whom were diagnosed with HD. All high-resolution images were extracted from prepared tissue slides after being scanned using the Aperio ScanScope CS (Aperio Technologies) at a 20 $\times$  magnification level (0.50  $\mu\text{m}/\text{pixel}$ ). All 30 WSIs were saved in an SVS format. For each layer, there exists ground truth annotations that a neonatal pathologist manually segmented (Figure 1). For



**Figure 2:** Overview of tissue regions on a whole slide section. **A:** Low magnification view (2x resolution) with the muscularis Propria tissue boundary outlined in green and the myenteric plexus path marked in magenta. **B:** Intermediate magnification showing myenteric plexus regions highlighted in magenta within the annotated muscularis propria segment. **C:** Close view of a ganglion cell (within the myenteric plexus regions) with high certainty outlined in yellow. **D:** Close view of a ganglion cell (within the myenteric plexus regions) with low certainty in cyan.

instance, the muscularis propria was manually delineated to represent the ground-truth segmentation, the myenteric plexus region was roughly manually segmented (a visually noticeable amount of tissue around the plexus regions was also included), followed by ganglion cells, which were also roughly segmented. A confidence level accompanies each ganglion cell annotation. A high-confidence level indicates strong certainty that the annotated object is a ganglion cell ([Figure 2C](#)). In contrast, a low-confidence level reflects that the object is believed to be a ganglion cell but with some uncertainty ([Figure 2D](#)), due to the difficulty in recognizing and selecting physiological features of a ganglion cell [\[15\]](#). An overview of these tissue regions at different magnification levels is shown in [Figure 2](#).

## 2.2. Data Preprocessing

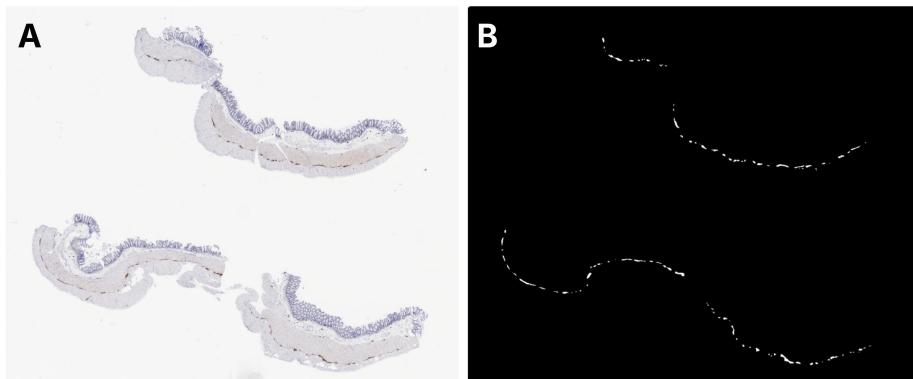
### 2.2.1. Muscularis Propria

Preprocessing for the muscularis propria region involved three steps. First, WSIs were down-sampled by a factor of 10

from the original 20x (0.50  $\mu\text{m}/\text{pixel}$ ) magnification, to the final resolution of 5  $\mu\text{m}/\text{pixel}$ . The lower magnification was chosen because the muscularis propria occupies a large, continuous region of the colon, so a more zoomed-out view of the tissue was better able to capture its global structure. Second, colour normalization was applied using the Macenko normalization method [\[18\]](#) to correct for staining variation. Finally, the 30 WSIs were tiled into 224x224-pixel patches (1000 tiles per WSI, for a total of 30,000 tiles), enabling the model to learn the characteristic texture, morphology, and appearance of the muscularis layer across the slide. An example of the muscularis ground truth annotations used in this stage is presented in [Figure 1B](#).

### 2.2.2. Myenteric Plexus

Preprocessing for the myenteric plexus region followed the same general procedure but was adjusted based on the plexus's anatomy. The WSIs were down-sampled by a factor of 4 from 20x (0.50  $\mu\text{m}/\text{pixel}$ ) to a final magnification of 5x



**Figure 3:** Reference annotations for tissue compartments in a whole slide section. **A:** Original WSI. **B:** Ground truth mask marking the myenteric plexus regions.

( $2.50 \mu\text{m}/\text{pixel}$ ). This level was chosen as it offers a compromise between resolution and spatial context, preserving finer details while still providing a sufficiently large field of view. Macenko normalization [18] was applied to correct for staining variation. The WSIs were tiled into  $224 \times 224$ -pixel patches, and the 30,000 tiles were selected if they contained at least one pixel within the muscularis propria region. This step was included because the myenteric plexus is strictly contained within the muscularis layer, so there is no need to tile other tissue regions. A reference example of plexus ground truth annotations is shown in Figure 1C, which shows how the plexus ground truth is exclusively within the muscularis segmentation, and in 3B.

### 2.2.3. Ganglion Cells

Preprocessing for ganglion cell segmentation followed the same general workflow but was modified to be suitable for cellular-level analysis. The WSIs were down-sampled by a factor of 2 from  $20\times$  ( $0.50 \mu\text{m}/\text{pixel}$ ) to a final magnification of  $10\times$  ( $1.0 \mu\text{m}/\text{pixel}$ ). The higher resolution provides a closer view of the tissue, preserving the fine morphological features required to identify ganglion cells. Macenko normalization [18] was used to correct for staining variation. The images were tiled into  $224 \times 224$  pixel patches, and only tiles containing at least one pixel from the myenteric plexus region were selected. Since ganglion cells are exclusively found within the plexus, this ensured that the model was trained only on relevant regions.

## 2.3. Training & Testing

### 2.3.1. Model Architecture

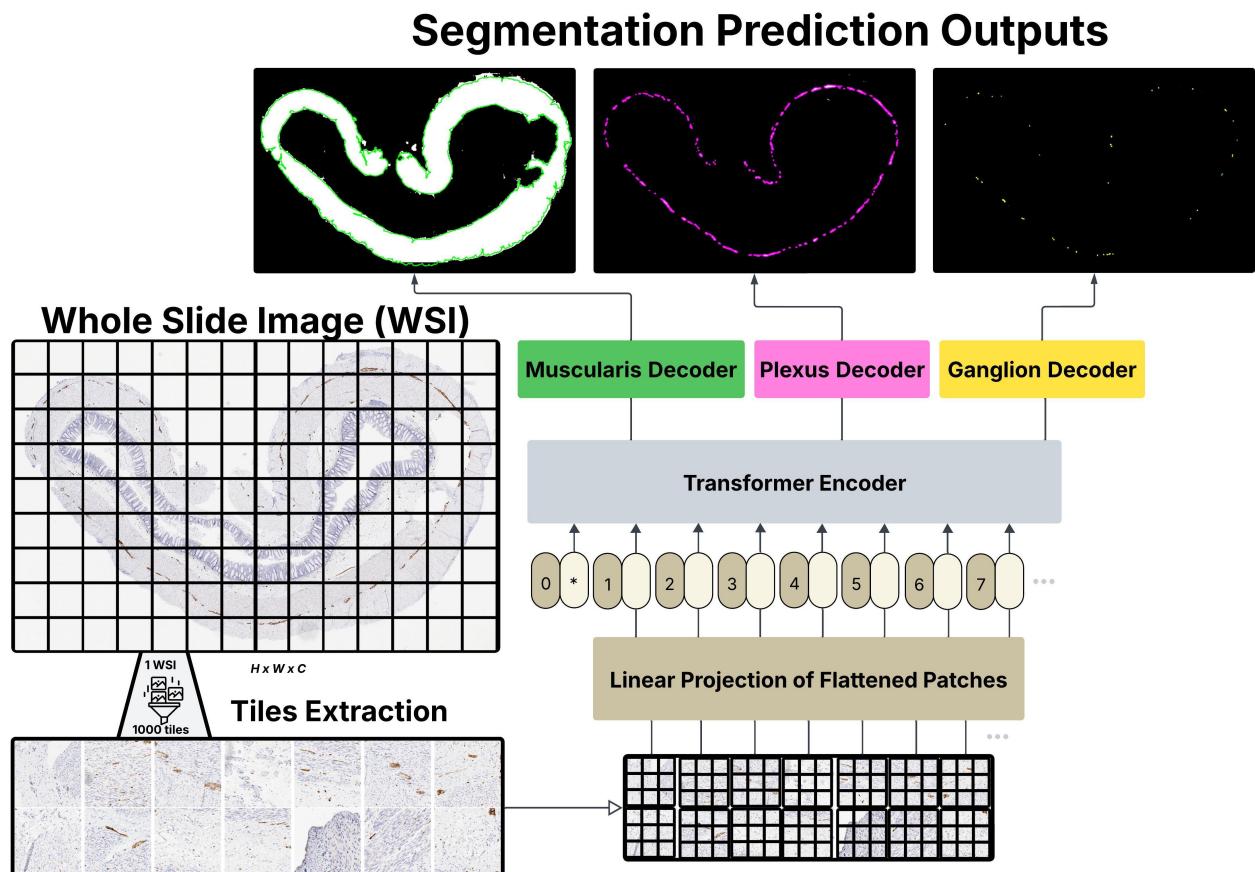
The segmentation framework was based on a ViT-B/16 backbone [17]. Specifically, each  $224 \times 224$  tile was divided into fixed-sized patches that were linearly projected to patch embeddings and then fed to the transformer encoder. The encoder generated a global contextual representation shared among the three downstream segmentation tasks. To enable the hierarchical workflow, three task-specific decoders were employed to independently generate prediction maps for the muscularis, plexus, and ganglion. For a given task, the decoder received the encoded representations and produced

a dense pixel output for the target layer. An overview of the complete three-stage pipeline is provided in Figure 4.

### 2.3.2. Muscularis Propria

To ensure robust segmentation of the muscularis propria and reduce the risk of overfitting, especially given the small dataset size, a 5-fold cross-validation strategy was employed. During training, the dataset was split into five equal folds, each comprising 6 WSIs (6000 tiles) for validation. This process was applied iteratively, so that each WSI was used both as a training set in some folds and as a test set in others. Thereby, making better use of the dataset. To increase dataset diversity, both sets of tiles were augmented using four methods: scaling, vertical flipping, horizontal flipping, and random rotation. The scaling operation simulated zoom effects by randomly enlarging or shrinking the visible region before resizing the tile back to  $224 \times 224$  pixels. This augmentation enabled the model to learn from tiles representing different effective magnifications. Rotation entailed rotating the tiles by  $90^\circ$ ; horizontal and vertical flipping consisted of flipping the images by  $180^\circ$ , with a  $\frac{1}{2}$  chance of being flipped. Due to the re-stitched prediction masks being binary, with white representing regions of interest and black representing areas not of interest, any rotation operation on the original ground truth mask tiles could result in inverting the areas of interest to be considered as regions not of interest by adding extra black space on the tile borders, thus introducing noise to the true performance. To overcome this, a function was deployed to trim the border regions of the augmented mask, thereby preserving the alignment between the training tile and the ground truth binary mask and reducing potential inconsistent tissue placement. For model training, the ViT-B/16 utilized a base learning rate of  $5e-4$  as the first hyperparameter. An AdamW optimizer was employed, along with a weight decay of  $1e-4$ . Furthermore, a cosine learning rate was used, gradually adjusting the learning rate over time, which converges faster and potentially yields better results. Training used five epochs for warm-up, followed by a complete run of 50 epochs with a batch size of 64.

Since the muscularis is the largest region within the colon, the tiles that were accepted correspond to random



**Figure 4:** Overview of the model pipeline for WSI processing and multi-target segmentation. The WSI is divided into fixed-size tiles that are linearly projected into patch embeddings before entering the transformer encoder. The encoded representations are passed to three task-specific decoders responsible for muscularis segmentation, plexus segmentation, and ganglion detection. The outputs of these decoders generate the corresponding prediction masks shown at the top.

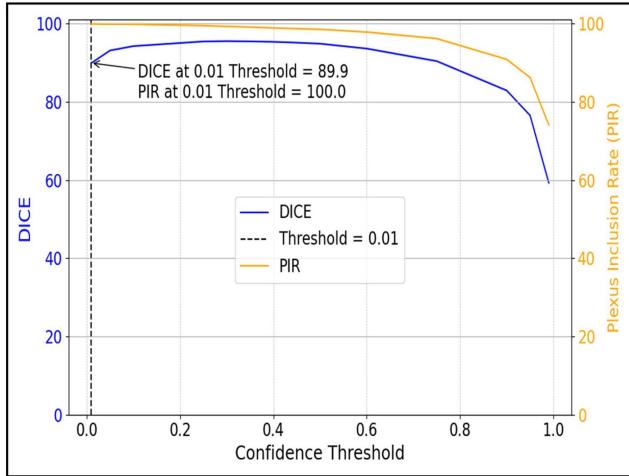
tiles; therefore, the tiles chosen for the model’s training purposes weren’t necessarily selected based on whether they contained plexus regions. After inputting each tile into the ViT model, the output consists of a pixel-wise raw logit, which is then converted to a probability between 0 and 1 after passing through a SoftMax activation function, representing the probability of the muscularis propria layer. The center region of the prediction tile is then retained and later used for re-stitching to form the full WSI prediction map. To evaluate the model’s predictive performance for each layer, a manual sweep operation was performed to ensure proper alignment and no missing regions between the ground truth segmentation binary map and the prediction segmentation binary mask. In terms of determining the most optimal threshold to accept a tile corresponding to a muscularis propria pixel, a threshold sweep of 0.01-0.99 was applied, and the threshold which produced the highest PIR and Dice coefficient was utilized. The effect of varying the prediction threshold on both Dice and PIR is summarized in Figure 5.

### 2.3.3. Myenteric Plexus

Similar to the training procedure for segmenting the muscularis propria, a 5-fold cross-validation strategy was also applied throughout training to minimize the risk of the

pre-trained ViT model overfitting to the dataset. This was achieved by utilizing five different folds, each consisting of 6 WSIs. Therefore, each WSI was used for both training and testing, thereby allowing the model to make full use of each WSI. Identical data augmentation techniques were applied when preparing the tiles for training purposes. This also included scaling, a 90° random rotation, and horizontal and vertical flipping of approximately 180°, each with a 50% chance of occurrence. To mitigate any potential effects that rotation operation could have had on performance, the border edges of the binary mask were trimmed to preserve alignment between the original tiles and ground truth tiles, specifically in regions of interest. The ViT-based model also utilized an AdamW optimizer. Unlike the segmentation of the muscularis propria, segmenting the myenteric plexus utilized a learning rate of 2e-4 and a weight decay of 1e-3. Similarly, a cosine learning rate scheduler was also used with five warmup epochs followed by a full training session of 50 epochs.

The ViT-based model output comprised pixel-wise logits per tile, which were converted to probabilities (0-1) via a SoftMax activation function, corresponding to the likelihood that a pixel represented the myenteric plexus layer. The tiles are then stitched after preserving their central regions,



**Figure 5:** Impact of confidence threshold on Dice coefficient and PIR.

forming the WSI binary prediction map. A manual sweep is conducted, comparing the model’s prediction to the ground-truth mask to evaluate segmentation performance. As with the muscularis propria layer, a threshold sweep was conducted to determine the optimal pixel value to accept as part of the myenteric plexus layer.

#### 2.3.4. Ganglion Cells

Similar to the previous segmentation steps, a 5-fold cross-validation scheme was employed to maximize data use while preventing overfitting. Five folds were defined, each with six WSIs. This procedure was repeated for all five folds, such that every WSI was used as a test sample exactly once and was never present in the training set of the same fold.

To increase data diversity, the  $224 \times 224$ -pixel tiles were augmented in four different ways. The first augmentation technique was scaling, followed by vertical and horizontal flipping by  $180^\circ$ , and rotation, which rotated an image by  $180^\circ$  with a  $\frac{1}{2}$  chance. Since the augmentations were also applied to binary masks, border trimming was performed to prevent flipping operations from inverting regions of interest and creating inconsistencies. Hyperparameters included a learning rate of  $5e-4$ , a weight decay of  $1e-4$ , and the AdamW optimizer. Training consisted of 5 warm-up epochs followed by 50 full training epochs, with a batch size of 64.

The ViT-based model output consisted of raw pixel logits, which were converted to probabilities between 0 and 1 via a SoftMax activation. The center regions of the tiles were preserved and re-stitched to reconstruct the full prediction mask. Tiles were then accepted as ganglion cell predictions if their probabilities exceeded the predetermined threshold. To evaluate predictive performance, a manual sweep was conducted to compare the stitched prediction maps with the ground truth segmentation maps.

## 2.4. Data Postprocessing

### 2.4.1. Myenteric Plexus

Binary mask tiles produced by the ViT model were first processed by retaining their center regions and re-stitching them to reconstruct the full WSI prediction map. Connected component labelling was applied to group adjacent positive pixels into individual plexus regions, and any predicted component with an area smaller than 50 pixels was removed to eliminate small spurious detections. This object-wise filtering step ensured that only anatomically plausible plexus regions were retained in the final prediction mask.

### 2.4.2. Ganglion Cells

Postprocessing for ganglion cell prediction followed the same initial procedure of retaining tile centers and re-stitching them to reconstruct the full WSI-level mask. Predicted ganglion regions were restricted to the previously segmented myenteric plexus by applying a binary mask, ensuring that detections were only evaluated within anatomically valid regions. Connected component labeling was then used to group positive pixels into individual ganglion candidates. To suppress small spurious detections, any predicted component with an area smaller than 10 pixels was removed. This object-level filtering step yielded a cleaner, anatomically consistent set of ganglion predictions without altering the morphology of the retained objects.

## 2.5. Performance Metrics

To evaluate segmentation performance across all three layers, several common metrics were used. Precision and recall quantify the model’s ability to correctly identify positive pixels, defined as

$$\text{Precision} = \frac{TP}{TP + FP} \quad (1)$$

$$\text{Recall} = \frac{TP}{TP + FN} \quad (2)$$

where a True Positive (TP) corresponds to a correctly identified positive pixel, a False Positive (FP) corresponds to a pixel incorrectly labeled as positive, and a False Negative (FN) represents a pixel that should have been labeled positive but was not.

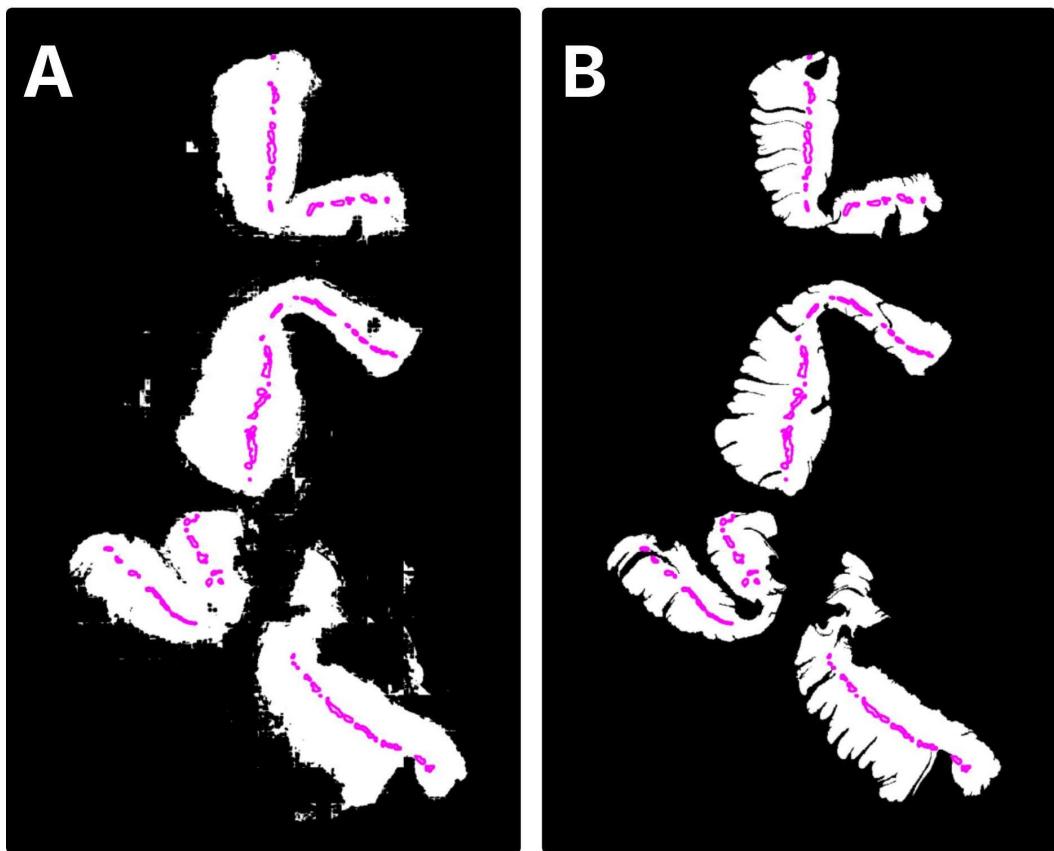
Layer-specific metrics were also used where appropriate. The Dice coefficient measures spatial overlap:

$$\text{Dice} = \frac{2|X_{ViT} \cap Y_{GT}|}{|X_{ViT}| + |Y_{GT}|} \quad (3)$$

where  $X_{ViT}$  is the predicted mask and  $Y_{GT}$  is the ground truth mask.

For evaluating the inclusion of relevant anatomical regions, two additional metrics were used. The Plexus Inclusion Rate (PIR) for muscularis segmentation is defined as

$$\text{PIR} = \frac{N_{ViT}}{N_{GT}}, \quad (4)$$



**Figure 6:** Muscularis propria layer segmentation on a whole slide section with plexus annotations overlaid. **A:** Stitched model prediction for the muscularis layer shown in white, with ground truth plexus regions overlaid in magenta. **B:** Ground truth muscularis mask with the same plexus regions overlaid in magenta. The magenta plexus annotations appear consistently in both panels, reflecting that these regions are fully captured by the model, yielding a perfect PIR.

where  $N_{GT}$  is the total number of ground truth plexus regions and  $N_{ViT}$  is the number of those plexus regions that intersect the predicted muscularis mask at least once. Similarly, the Ganglia Inclusion Rate (GIR) for plexus segmentation is defined as

$$\text{GIR} = \frac{N_{ViT}}{N_{GT}}, \quad (5)$$

where  $N_{GT}$  is the total number of ground truth ganglion regions and  $N_{ViT}$  is the number of those ganglion regions that are fully contained within, or intersect, the predicted plexus mask.

### 2.5.1. Muscularis Propria

For muscularis segmentation, performance was assessed using the Dice coefficient, PIR, precision, and recall. The Dice coefficient measured the overlap between the predicted and ground-truth muscularis regions. At the same time, PIR quantified how effectively the segmented muscularis retained the plexus regions, which is necessary for downstream layers. Precision and recall quantified pixel-level accuracy. A threshold sweep identified 0.01 as the optimal threshold, producing the highest Dice coefficient and PIR (Figure 5).

### 2.5.2. Myenteric Plexus

For myenteric plexus segmentation, precision, recall, and GIR were used. Precision measured how many predicted plexus pixels were correct, while recall measured how completely the plexus regions were captured. GIR quantified how many ganglion-containing plexus regions were retained, which directly affects the success of ganglion detection. The optimal threshold identified for plexus segmentation was 0.25.

### 2.5.3. Ganglion Cells

For ganglion cell detection, precision and recall were used to evaluate object-level accuracy. Precision represented the proportion of predicted ganglion cells that were correctly identified, while recall measured how completely the model captured the ganglion cells present in the ground truth. A threshold sweep determined that 0.175 provided the best balance between precision and recall for ganglion.

## 2.6. Baseline Model

### 2.6.1. Muscularis Propria

The baseline model's methodology to successfully segment the muscularis propria consisted of utilizing Calretinin-stained WSI of the colon cross-sections and deploying a colour-based k-means clustering algorithm to segment the

**Table 1**Performance comparison (%) of segmentation models on the **Muscularis Propria** layer.

Model	Dice	Precision	Recall	PIR
K-means [14]	70.7	70.6	78.9	77.4
CNN [13, 14]	89.2	81.9	96.2	96.0
<b>ViT-B/16 (proposed) [15]</b>	<b>89.9</b>	<b>82.4</b>	<b>99.7</b>	<b>100</b>

colon section [14]. Additionally, a CNN-based algorithm that was pretrained on ImageNet-1k was also utilized to segment the muscularis propria layer [13, 14], using tiles with dimensions of 256×256 pixels. Postprocessing techniques were also applied, which included morphological smoothing [19]. The segmentation performance was evaluated utilizing the Dice coefficient, PIR, in addition to precision and recall.

### 2.6.2. Myenteric Plexus

The segmentation of the myenteric plexus involved receiving the binary masks, corresponding to the model's prediction output, and utilizing these images for the segmentation process. Unlike the muscularis propria segmentation, the myenteric plexus was segmented first, using a colour-based thresholding to differentiate between brown-stained areas [14]. This was then followed by applying a pixel threshold to accept and reject pixels that could correspond to being a part of the myenteric plexus layer. Furthermore, morphological filtering [19] was also applied to filter out pixels that potentially correspond to true or false myenteric plexus regions for the final stage of segmentation. To evaluate the segmentation accuracy of the myenteric plexus layer, the GIR was also utilized, alongside precision and recall.

### 2.6.3. Ganglion Cells

The final step included classifying true ganglion cells within each plexus region, also utilizing a colour intensity-based threshold to determine candidates for ganglion cells [14], fitting categories with different certainties, high-certainty is represented by being compact, with a visible nuclei. Whereas, low-certainty candidates are borderline ganglion, which do appear to be ganglion; however, they do lack typical morphology and ganglion physical traits. Furthermore, physical characteristics such as area, circularity, colour, and gradient contrast were used as descriptors for each candidate, and then used to train an LDA model to classify true ganglion cells from false candidates.

## 3. Results

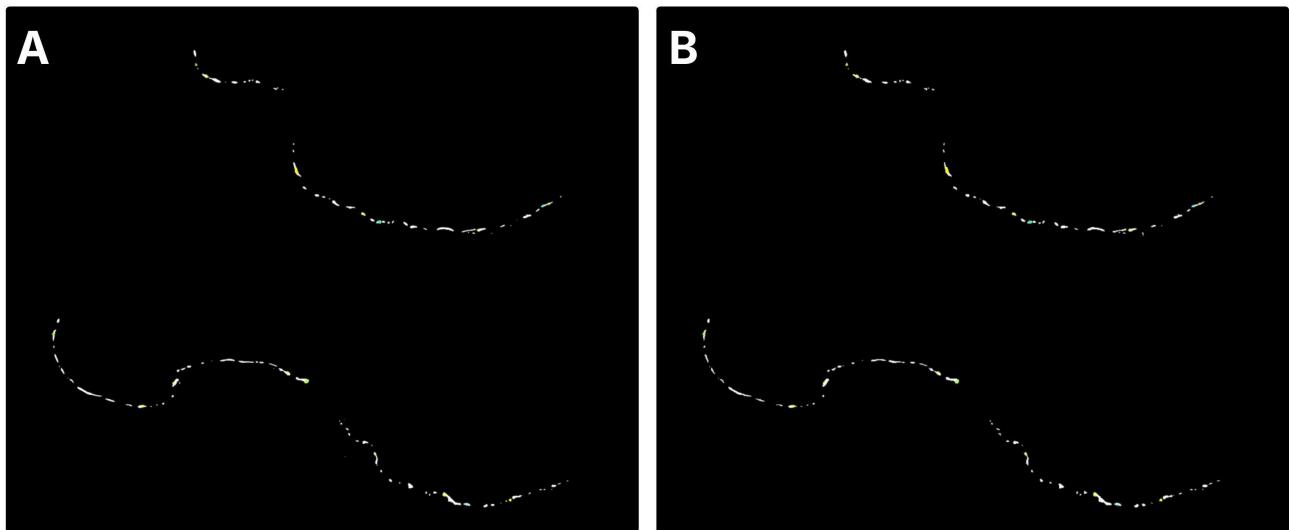
This research study's results section aims to evaluate the segmentation performance of the ViT-B/16 model across three layers of the colon: muscularis propria, myenteric plexus, and ganglion cells. The metrics presented pertain to each segmented layer.

**Table 2**Performance comparison (%) of segmentation performance on **Myenteric Plexus** layer.

Model	Precision	Recall	GIR
Colour-based thresholding [11, 14]	73.8	85.9	99.2
<b>ViT-B/16 (proposed)</b>	<b>84.2</b>	<b>94.8</b>	<b>99.7</b>

The segmentation of the muscularis propria region is of significant importance, as it shapes the anatomical boundaries for the later identification of ganglion cells located within the myenteric plexus. Therefore, isolating the muscularis layer allows segmentation of the plexus region, where ganglion cells are typically found. Furthermore, this narrows the region of interest from a much broader WSI to only tissue regions clinically more relevant to diagnosing HD, as evidenced by a high Dice coefficient. Moreover, confidence thresholds, which determine the model's ability to include a tile as part of the muscularis propria, were experimented with. A confidence threshold of 0.4 (Figure 5) yields model results that are similar to those of the ground truth. However, a lower confidence threshold of 0.01 achieves a superior PIR with a lower Dice coefficient, giving further certainty that all tiles include a plexus region, which serves the ultimate goal of detecting ganglion cells better, as shown in Figure 5.

The results, as shown in Tables 1, demonstrate a high Dice coefficient of 89.9%, indicating a more accurate segmentation performance of the muscularis layer compared to the ground truth manual segmentation. This not only shows a high overall accuracy but also superior performance compared to the baseline CNN and K-means models. Additionally, the model achieved a precision score of 82.4%, indicating that the muscularis region corresponds to muscularis tissue, mislabeling only a few non-muscularis tissues. The model also achieved a recall score of 99.7%, demonstrating its superior ability to identify all positive tissue pixels. The model also achieved a PIR of 100%, which was deemed highly important, as all ganglion cells that needed to be identified later resided within the myenteric plexus. That said, all model-predicted outputs contained plexus regions, which play an instrumental role in the subsequent segmentation stage. A qualitative comparison of predicted muscularis segmentation with plexus overlays is shown in Figure 6.



**Figure 7:** Myenteric plexus region segmentation on a whole slide section with ganglion cell annotations overlaid. **A:** Stitched model prediction for the plexus regions displayed in white. **B:** Corresponding ground truth mask for the plexus regions. Yellow and cyan objects represent high- and low-certainty ganglion cells, respectively. All 28 ganglion cells fall within the predicted plexus regions, demonstrating complete ganglion inclusion in the plexus segmentation.

The following segmentation stage included segmenting the myenteric plexus section of the colon. This stage is considered the most crucial segmentation phase before ganglion identification, as ganglion cells are located within the myenteric plexus layer. The ViT-B/16 proposed model achieves high segmentation accuracy in the plexus layer, which prepares for optimal ganglia identification in the downstream task. A confidence threshold of 0.25 was used to select the predicted pixels in the binary output prediction map. Although the chosen confidence threshold was flexible, it yielded a higher recall at the expense of precision, prioritizing a higher GIR score. This threshold yielded the optimal GIR, as well as precision and recall. A representative plexus prediction example with ganglion cell annotations is shown in Figure 7.

The results, as shown in Tables 2, demonstrate a high recall score of 94.8%, indicating the model's ability to detect almost all positive samples, corresponding to regions that contain myenteric plexus. Furthermore, a high precision score of 84.2% demonstrates the model's capability of detecting regions that not only accurately pertain to the myenteric plexus but also limit false positives by being confined to regions enclosed by the muscularis. Finally, the ViT model achieved a high GIR score of 99.7% when segmenting the myenteric plexus. Given the myenteric plexus' thin and long anatomical structure, and its thin, elongated, and discontinuous nature, the model's segmentation performance, with high GIR (Figure 7), precision, and recall, allows for higher chances of success in later ganglion cell identification. A detailed breakdown of per-slide plexus segmentation outcomes is provided in Table 3.

Once the myenteric plexus was segmented, the next stage involved segmenting and identifying ganglion cells to further enhance clinical relevance to HD. This segmentation

stage involved transitioning from region-level segmentation to cell-level (object) prediction. The threshold utilized for segmenting and detecting the ganglion cells is 0.175. An example visualization of predicted high- and low-certainty ganglion cells across a whole slide is shown in Figure 8.

The ViT-B/16 proposed model achieved a high recall average score of 78.6% (Table 4) by utilizing the combined low-certainty high-certainty ganglion cells, demonstrating strong capabilities in detecting the majority of ganglion cells within the myenteric plexus region. Furthermore, using a more conservative threshold for higher certainty yields a precision of 62.1%, indicating the model's ability to identify clear, morphologically distinct ganglion cells. Using the same high-certainty threshold, the model achieved a high recall of 89.1%, thereby identifying the overwhelming majority of ganglion cells with a typical ganglion-cell morphology. Per-slide high-certainty ganglion identification results are summarized in Table 5. Combined high and low certainty object-level results for each WSI are listed in Table 6. Overall, the DL model used for segmenting and identifying ganglion cells outperformed the baseline model in ganglion cell detection.

## 4. Discussion

### 4.1. Interpretation of Results

This research study proposed a multi-stage pipeline to segment three colon layers, with the overall goal of identifying the presence or absence of ganglion cells. The segmentation process used a ViT-based model to first segment the muscularis propria, then the myenteric plexus, and finally the ganglion cells within the myenteric plexus. The proposed approach achieved high Dice coefficient and PIR scores when segmenting the muscularis, a high GIR

**Table 3**

Per-WSI segmentation metrics including true positives (TP), false positives (FP), false negatives (FN), precision, recall, and GIR at a prediction threshold of 0.25. All metrics correspond to region-wise evaluation within the segmented Myenteric Plexus.

No.	WSI Name	TP	FP	FN	Precision	Recall	GIR
1	S14-580	158	55	17	0.742	0.903	0.978
2	S00-1910	79	5	1	0.940	0.988	1.0
3	S02-410	101	8	2	0.927	0.981	1.0
4	S02-484	141	24	0	0.855	1.000	1.0
5	S03-2391	96	7	9	0.932	0.914	1.0
6	S01-18	182	32	9	0.850	0.953	1.0
7	S03-3178 D2	63	8	1	0.887	0.984	1.0
8	S03-3178 D3	47	16	0	0.746	1.000	1.0
9	S03-3178 D4	59	6	6	0.908	0.908	1.0
10	S04-52	82	16	2	0.837	0.976	1.0
11	S04-910	156	33	8	0.825	0.951	1.0
12	S07-1808	124	30	1	0.805	0.992	1.0
13	S08-2215	161	26	10	0.861	0.942	1.0
14	S09-2723	43	21	7	0.672	0.860	1.0
15	S04-1840	118	13	4	0.901	0.967	0.976
16	S07-1465	162	37	17	0.814	0.905	1.0
17	S14-1715	62	12	2	0.838	0.969	1.0
18	S09-2909	71	20	0	0.780	1.000	1.0
19	S14-3414	108	8	5	0.931	0.956	1.0
20	S14-2038	164	33	2	0.832	0.988	1.0
21	S15-1442	127	23	0	0.847	1.000	1.0
22	S15-1518	111	9	4	0.925	0.965	1.0
23	S16-567	118	34	4	0.776	0.967	1.0
24	S16-1197 B1	51	16	0	0.761	1.000	1.0
25	S11-1760	223	14	17	0.941	0.929	1.0
26	S16-1467	84	10	3	0.894	0.966	0.950
27	S16-1197 B3	70	8	8	0.897	0.897	1.0
28	S16-1197 B2	100	35	1	0.741	0.990	1.0
29	S97-2054	90	6	53	0.938	0.629	1.0
30	S16-1415	36	18	2	0.667	0.947	1.0

score when segmenting the myenteric plexus, and, finally, high precision and recall scores for identifying ganglion cells within the plexus region. Demonstrating the model's capability to detect fine colon boundaries and detect ganglia with confidence. Performance evaluations surpassed those of conventional CNN methods and other shallow machine learning methods, indicating superior performance for segmenting colon tissue. These findings further demonstrate the potential of DL algorithms to aid HD diagnosis.

The muscularis propria is a smooth muscle layer embedded within the walls of organs such as the GI tract and bladder, composed of an inner circular layer and an outer longitudinal layer, giving it a more uniform structure [26]. The ViT has an inherent self-attention mechanism, enabling it to attend to other image patches and capture long-range dependencies. For instance, it could relate to the appearance of the muscularis propria and its location relative to different regions of the image that correspond to the lamina propria

**Table 4**Performance comparison (%) of segmentation performance on **Ganglion Cell** identification.

Model	Precision (high-certainty)	Recall (high-certainty)	Precision (combined)	Recall (combined)
LDA model [14]	60.9	82.1	64.8	<b>80.2</b>
<b>ViT-B/16 (proposed)</b>	<b>62.1</b>	<b>89.1</b>	<b>67.0</b>	78.6

or the submucosa. The segmentation performance of the ViT model yielded a Dice coefficient of 89.9%, indicating strong agreement between the predicted and ground truth maps. These results demonstrate that global-context information, combined with local feature representations, enables ViTs to achieve high segmentation accuracy, as evidenced by the Dice coefficients obtained. Additionally, segmentation performance on this layer yields a PIR of 100%, ensuring that the plexus segmentation is applied within the anatomical region of interest for later accurate segmentation of ganglion cells.

The myenteric plexus region is embedded within the inner layer of the muscularis propria and consists of nerve fibres responsible for several autonomous digestive functions. Unlike the muscularis propria, the myenteric plexus is a smaller, less structured region and thus appears more fragmented across WSIs [26, 27], making segmentation more difficult. The myenteric plexus is of great importance because it is the anatomical site where ganglion cells are present [27], which are needed for future segmentation and identification. The segmentation performance yielded recall and precision scores of 94.8% and 84.2%, respectively. Indicating the model's robust capability of segmenting most of the pixels that correspond to the plexus region. Additionally, among the tiles segmented, fewer false plexus segmentations (false positives) are present. Furthermore, the model produced a GIR of 99.7%, indicating that it not only retained accurate segmentations but also identified relevant areas containing ganglion cells. Because the myenteric plexus is embedded within the muscularis propria, producing textural variability [26], this allows ViTs to leverage their inherent self-attention mechanisms, achieving superior results. For instance, the attention layers embedded within the ViT enable the model to distinguish between the thick, continuous, and homogeneous fibres of the muscularis layer and the thin, irregular fibres of the plexus, even when the muscularis propria and the myenteric plexus share the exact colour tones [28]. The ViT computes attention scores between patch pairs, enabling it to identify which patches are similar or different based on the embeddings it produces. This allows the model to distinguish between slight variations in texture. Furthermore, given the fragmented myenteric plexus throughout the WSI, it is discontinuous across several patches, providing ViTs with an opportunity to leverage their inherent understanding of global contextual relationships. Unlike conventional CNNs, which primarily rely on local receptive fields, this may prevent recognizing that spatially separated patches belong to the same structure, leading

to fragmented predictions, especially when there are mild structural differences between patches. However, since ViTs operate globally, the model can label patches with slightly different structures as belonging to the plexus region, even when their embeddings are similar. That being said, high segmentation accuracy and high GIR of the myenteric plexus are critical to optimizing ganglion cell segmentation and identification.

The final stage of segmentation is considered the most biologically significant, aiming to identify the presence or absence of ganglion cells, which is directly attributed to the clinical diagnosis of HD. Unlike the two previous stages of segmentation, ganglion cell segmentation transitions from region-level segmentation to object-level segmentation of cells. Due to morphological differences between different ganglion cells, such as their shape, size and sparse distribution within the myenteric plexus [29], their segmentation is increasingly complex compared to previous layers. Segmenting the ganglion cells used two different manual annotations: one for ganglion cells classified as "high-certainty", thus possessing typical ganglion morphological characteristics. Whilst "low-certainty" is utilized for ganglion cells that lack typical morphological representations. At a high confidence level, the model achieves 62.1% precision and 89.1% recall. The average of the high- and low-certainty results yields recall and precision scores of 78.6% and 67.0%, respectively. Based on these results, the model is more conservative, producing fewer but highly reliable detections. On the other hand, a less strict threshold ensures that most ganglion cells are predicted, at the cost of limited reliability. Unlike previous layers for segmentation, where the ViT model leveraged its innate self-attention mechanism to distinguish local regions, for this segmentation task, the ViT leverages its innate mechanism to understand morphological cues corresponding to ganglion cell characteristics by integrating fine-grained texture into feature embeddings. For instance, the self-attention mechanism for this task interprets patterns of nuclei, dendrites, and axons, thus focusing on local micro-relationships rather than long-range dependencies. Furthermore, the ViT architecture implements certainty thresholds for ganglion cells, where the embedded attention heads agree on whether the prediction is most likely a ganglion cell; in low-certainty mode, the attention heads produce more dispersed attention on the segmented prediction. This two-tier framework shows how these results, in terms of certainty, might be helpful for diagnostic confirmation or further review, similar to applications in a clinical setting. Ganglion cell identification using DL demonstrates the

**Table 5**

Object-wise **ganglion (high-certainty only)** identification results for each WSI at a prediction threshold of 0.175. Metrics (TP, FP, FN, precision, recall) are computed per ganglion region.

No.	WSI Name	TP	FP	FN	Precision	Recall
1	S14-580	83	34	9	0.709	0.902
2	S00-1910	4	2	0	0.667	1.0
3	S02-410	23	17	0	0.575	1.0
4	S02-484	32	46	0	0.410	1.0
5	S03-2391	13	9	4	0.591	0.765
6	S03-3178 D2	5	2	0	0.714	1.0
7	S03-3178 D3	3	9	0	0.250	1.0
8	S03-3178 D4	1	7	0	0.125	1.0
9	S04-52	8	11	8	0.421	0.50
10	S04-910	30	7	0	0.811	1.0
11	S07-1808	82	33	0	0.713	1.0
12	S08-2215	50	18	1	0.735	0.980
13	S09-2723	22	24	2	0.478	0.917
14	S04-1840	41	4	0	0.911	1.0
15	S07-1465	32	74	0	0.302	1.0
16	S14-1715	4	1	6	0.800	0.40
17	S09-2909	43	41	1	0.512	0.977
18	S14-3414	50	17	3	0.746	0.943
19	S14-2038	70	30	2	0.700	0.972
20	S15-1442	82	35	1	0.701	0.988
21	S15-1518	63	20	16	0.759	0.797
22	S16-567	117	23	29	0.836	0.801
23	S16-1197 B1	51	14	0	0.785	1.0
24	S11-1760	24	19	2	0.558	0.923
25	S16-1467	19	4	1	0.826	0.950
26	S16-1197 B3	35	8	3	0.814	0.921
27	S16-1197 B2	32	5	0	0.865	1.0
28	S97-2054	0	1	4	0.0	0.0
29	S16-1415	22	2	0	0.917	1.0
30	S01-18	10	15	0	0.400	1.0

potential to automate the most critical step in diagnosing HD, reducing diagnostic variability and further assisting pathologists.

Previous research has used similar approaches to segment and detect the potential presence of ganglion cells [13, 14]. For instance, a traditional CNN model [13, 14] was utilized to segment the muscularis propria layer alongside a shallow learning model, a K-means clustering algorithm [14]. This resulted in the CNN model producing a Dice

coefficient of 89.2%, a precision of 81.9%, a recall of 96.2% and a PIR of 96%. Compared to our 89.9%, 82.4%, 99.7% and 100%, respectively. Performance regarding the myenteric plexus tissue in [11, 14], utilized a colour-based threshold framework and produced a recall score of 85.9%, precision of 73.8%, GIR of 99.2%. Compared to our 94.8%, 84.2%, and 99.7%, respectively. Regarding the ganglion segmentation, for the LDA model [14], the combined certainty ganglia (high and low certainty) produced a recall of 80.2% and

**Table 6**

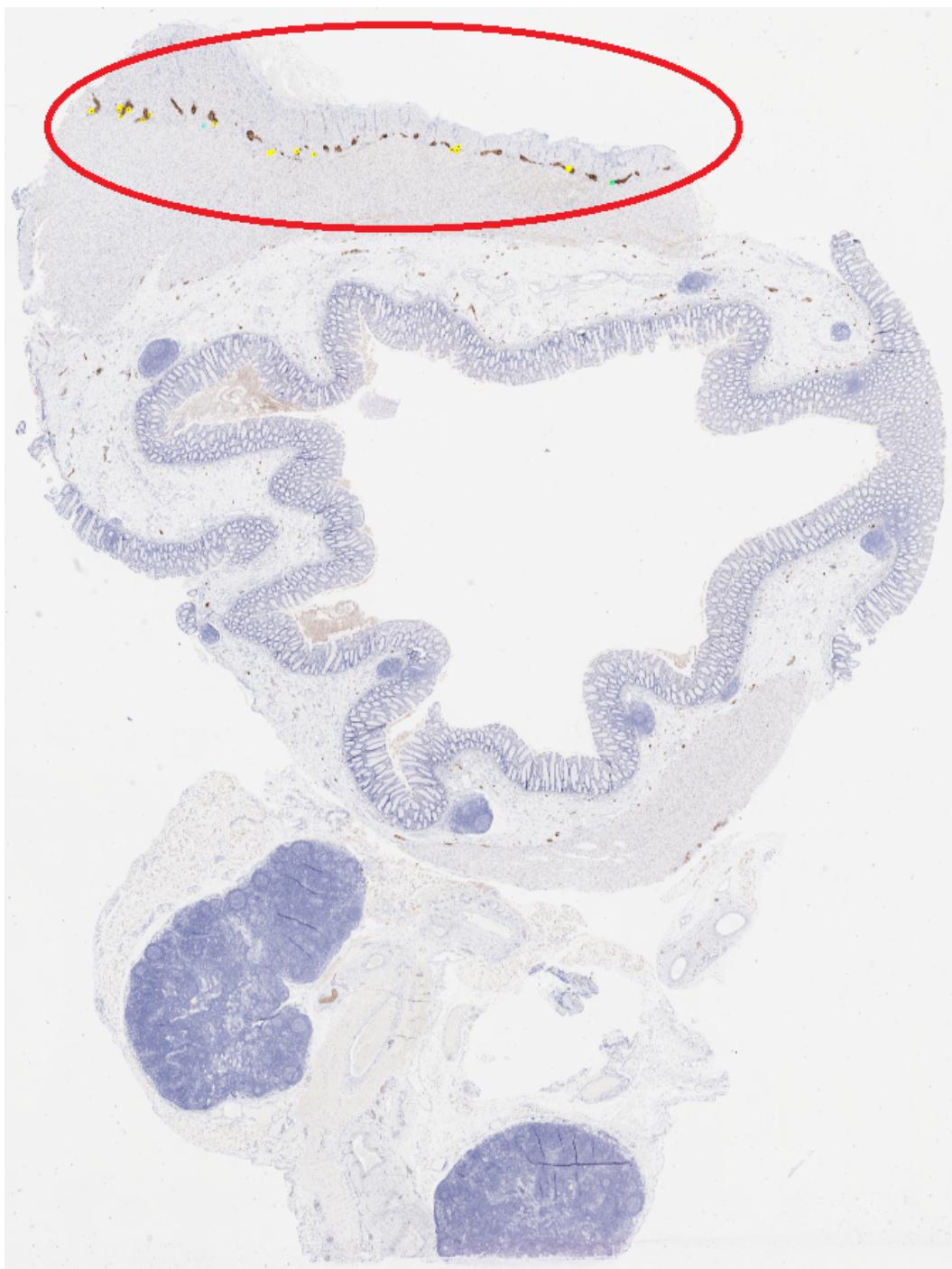
Region-level **ganglion identification performance for combined (high + low certainty) cells** at a prediction threshold of 0.175. Metrics represent object-wise TP, FP, and FN counts, along with precision and recall.

No.	WSI Name	TP	FP	FN	Precision	Recall
1	S14-580	86	31	14	0.735	0.860
2	S00-1910	5	1	1	0.833	0.833
3	S02-410	26	15	2	0.634	0.929
4	S02-484	40	37	2	0.519	0.952
5	S03-2391	13	9	6	0.591	0.684
6	S03-3178 D2	5	2	2	0.714	0.714
7	S03-3178 D3	4	8	0	0.333	1.0
8	S03-3178 D4	4	4	2	0.50	0.667
9	S04-52	12	7	13	0.632	0.480
10	S04-910	34	3	3	0.919	0.919
11	S07-1808	90	25	8	0.783	0.918
12	S08-2215	56	13	3	0.812	0.949
13	S09-2723	25	22	7	0.532	0.781
14	S04-1840	44	1	13	0.978	0.772
15	S07-1465	37	69	1	0.349	0.974
16	S14-1715	4	1	10	0.80	0.286
17	S09-2909	46	40	12	0.535	0.793
18	S14-3414	51	16	3	0.761	0.944
19	S14-2038	74	27	22	0.733	0.771
20	S15-1442	86	32	6	0.729	0.935
21	S15-1518	72	11	25	0.867	0.742
22	S16-567	120	20	37	0.857	0.764
23	S16-1197 B1	55	10	14	0.846	0.797
24	S11-1760	27	16	15	0.628	0.643
25	S16-1467	19	4	3	0.826	0.864
26	S16-1197 B3	39	4	8	0.907	0.830
27	S16-1197 B2	32	5	5	0.865	0.865
28	S97-2054	0	1	6	0.0	0.0
29	S16-1415	22	2	2	0.917	0.917
30	S01-18	14	11	0	0.560	1.0

precision of 64.8%, compared to our 78.6% and 67.0%, respectively. For the high-certainty ganglion segmentation results reported in [14], precision was 60.9% and recall was 82.1%, compared to our 62.1% and 89.1%. Results from this study highlight the ViT's ability to leverage its inherent self-attention mechanism to learn global dependencies across an image. This allows the model to accurately isolate the muscularis propria layer, consistently linking relationships within the fragmented myenteric plexus and understanding

local microcellular structure at the microscale. In contrast to conventional CNNs or other shallow learning models, which typically rely on localized features within each map, ViT's architecture enables global reasoning across different structures and textures within the same image or region, thereby achieving higher performance than traditional CNNs for segmentation tasks.

This research study proposed a multi-stage hierarchical segmentation pipeline. First, segment the muscularis propria



**Figure 8:** WSI with predicted ganglion cells overlaid. Cyan objects represent low certainty predictions, and yellow objects represent high certainty predictions. All visible ganglion cells in this section are correctly identified by the model.

layer, then the myenteric plexus layer, and finally segment and identify the presence of ganglion cells within the plexus layer, mirroring the workflow of a pathologist when detecting ganglion cells in the plexus region. Further, the consistent results highlight the potential of DL being utilized in a clinical setting as a confirmatory tool for diagnosing HD. The advantages of this study could be reduced time to

interpret WSIs and an automated solution. Reducing inter-observer variability also provides a uniform way to segment and confirm the presence of HD.

#### 4.2. Study Limitations

Although this study demonstrated strong segmentation performance and is deemed promising for future steps, it does include certain limitations that could limit its generalizability to future unseen data. For instance, the data used

for this research study were obtained solely from CHEO, as mentioned in Section 2.1, thereby limiting variability in structural and staining differences across WSIs that would have been acquired from different datasets and potentially increasing model generalizability. Furthermore, this study uses a small dataset of 30 WSIs, which increases the potential for overfitting during training, especially since large DL models such as CNNs and ViTs typically yield consistent improvements on larger datasets. Furthermore, since a single pathologist labeled the ground truth masks, this could introduce annotation bias; that said, a variety of expert manual annotations could improve model generalization in future work. Addressing these limitations in future studies through larger multi-center datasets with a variety of staining techniques and structures, a range of expert annotations, and increased dataset size could further enhance model robustness and clinical relevance.

### 4.3. Future Work

Although the limitations mentioned above limit the model's generalization ability, the ViT-based model used for segmentation and ganglion cell detection nevertheless produced highly accurate segmentations compared to conventional LDA and other shallow learning models. That said, future research work could benefit from optimizing the utility of ViTs' performance to increase its clinical applicability by incorporating data corresponding to similar forms of colon disease, such as megacystis microcolon intestinal hypoperistalsis syndrome and internal anal sphincter achalasia, giving the model opportunities to learn from different colonic structures and anatomical locations. Furthermore, future models of ViTs could also utilize other forms of self-supervised learning, such as distillation with no labels (DINO I & II) [20, 21] or other forms of Masked Autoencoding [22], such as a curriculum masking schedule [23], both allowing the model to learn essential features corresponding to the structure, texture and location of the colon from solving complex pretext tasks. Collectively, these directions will advance the integration of computer vision into digital pathology, enabling reliable, objective utility for HD diagnosis.

## 5. Conclusion

This study presents a hierarchical segmentation framework that employs a ViT-B/16 model to segment the muscularis propria and myenteric plexus, enabling subsequent detection and identification of ganglion cells, which satisfies a key diagnostic criterion for HD. The workflow mirrors the approach used by pathologists, sequentially isolating the muscularis layer, delineating the plexus regions, and identifying ganglion cells within them. The framework achieved a Dice coefficient of 89.9% and a PIR of 100% for muscularis segmentation, a GIR of 99.7% for plexus segmentation, and a precision and recall of 62.1% and 89.1% for high-certainty ganglion cell identification. These findings illustrate ViTs' ability to use self-attention to capture relevant features and

global contextual relationships, enabling the model to distinguish fine textural details at both regional and cellular scales. The results further suggest that such a framework could reduce assessment time and diagnostic variability. Although future work would benefit from larger multi-center datasets and multiple expert annotations, this study demonstrates the promise of DL based methods for supporting histological assessment in HD.

### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## References

- [1] Y. Chen, X. Yuan, Y. Li, S. Wu, X. Miao, J. Gong, and Y. Huang, "The prevalence and clinical presentation of Hirschsprung's disease in preterm infants: a systematic review and meta-analysis," *Pediatric Surgery International*, vol. 38, no. 4, pp. 523–532, 2022, doi: 10.1007/s00383-021-05054-2.
- [2] J. Kessmann, "Hirschsprung's disease: diagnosis and management," *American Family Physician*, vol. 74, no. 8, pp. 1319–1322, 2006.
- [3] G. Campanella, M. G. Hanna, L. Geneslaw et al., "Clinical-grade computational pathology using weakly supervised deep learning on whole slide images," *Nature Medicine*, vol. 25, pp. 1301–1309, 2019, doi: 10.1038/s41591-019-0508-1.
- [4] J. Yan, Y. Zeng, J. Lin, Z. Pei, J. Fan, C. Fang, and Y. Cai, "Enhanced object detection in pediatric bronchoscopy images using YOLO-based algorithms with CBAM attention mechanism," *Heliyon*, vol. 10, no. 12, p. e32678, 2024, doi: 10.1016/j.heliyon.2024.e32678.
- [5] F. Li, L. Sun, K. Y. Lam, S. Zhang, Z. Sun, B. Peng, H. Xu, and L. Zhang, "Segmentation of human aorta using 3D nnU-net-oriented deep learning," *Review of Scientific Instruments*, vol. 93, no. 11, p. 114103, 2022, doi: 10.1063/5.0084433.
- [6] M. Duci, A. Magoni, L. Santoro, A. P. Dei Tos, P. Gamba, F. Uccheddu, and F. Fascetti-Leon, "Enhancing diagnosis of Hirschsprung's disease using deep learning from histological sections of post pull through specimens: preliminary results," *Pediatric Surgery International*, vol. 40, no. 1, p. 12, 2023, doi: 10.1007/s00383-023-05590-z.
- [7] S. Lotfollahzadeh, M. Taherian, and S. Anand, "Hirschsprung disease," in *StatPearls, StatPearls Publishing*, 2023.
- [8] A. Mukherjee et al., "The placental distal villous hypoplasia pattern: interobserver agreement and automated fractal dimension as an objective metric," *Pediatric and Developmental Pathology*, vol. 19, no. 1, pp. 31–36, 2016.
- [9] Y. Megahed, R. Ducharme, A. Erman, M. Walker, S. Hawken, and A. D. C. Chan, "USF-MAE: Ultrasound Self-Supervised Foundation Model with Masked Autoencoding," *arXiv preprint arXiv:2510.22990*, 2025. doi:10.48550/arXiv.2510.22990.
- [10] A. J. Demetris et al., "Intraobserver and interobserver variation in the histopathological assessment of liver allograft rejection," *Hepatology*, vol. 14, no. 5, pp. 751–755, 1991.
- [11] J. Kurian et al., "Image Processing and Analysis of Histopathological Images Relating to Hirschsprung's Disease," *CMBES Proceedings*, vol. 41, 2018.
- [12] M. T. K. Law, A. D. C. Chan, and D. El Demellawy, "Color image processing in Hirschsprung's disease diagnosis," in *2016 IEEE EMBS International Student Conference (ISC)*, 2016, pp. 1–4.
- [13] C. McKeen, F. Zabihollahy, J. Kurian, A. D. Chan, D. El Demellawy, and E. Ukwatta, "Machine learning-based approach for fully automated segmentation of muscularis propria from histopathology images of intestinal specimens," in *Medical Imaging 2019: Digital Pathology*, vol. 10956, pp. 146–151, Mar. 2019.

- [14] J. A. Kurian, "Automated Identification of Myenteric Ganglia in Histopathology Images for the Study of Hirschsprung's Disease", M.A.Sc. thesis, Carleton University, 2021.
- [15] Y. Megahed, . A. . Fuller, . S. . Abou-Alwan, . D. El Demellawy, and A. Chan, "Segmentation of Muscularis Propria in Colon Histopathology Images Using Vision Transformers for Hirschsprung's Disease," *CMBES Proc.*, vol. 47, no. 1, May 2025.
- [16] Y. Megahed, A. Madi, D. El Demellawy, and A. D. C. Chan, "Knowledge-Driven Vision-Language Model for Plexus Detection in Hirschsprung's Disease," *arXiv preprint arXiv:2510.21083*, Oct. 2025. doi: 10.48550/arXiv.2510.21083.
- [17] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16×16 words: Transformers for image recognition at scale," in *Proc. Int. Conf. Learn. Representations (ICLR)*, 2021. Available: <https://openreview.net/forum?id=YicbFdNTTy>.
- [18] M. Macenko et al., "A method for normalizing histology slides for quantitative analysis," in *2009 IEEE International Symposium on Biomedical Imaging: From Nano to Macro*, 2009, pp. 1107–1110.
- [19] J. Serra, *Image Analysis and Mathematical Morphology*. U.K.: Academic Press, 1982, ISBN: 978-0126372410.
- [20] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, "Emerging properties in self supervised Vision Transformers," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2021.
- [21] M. Oquab et al., "DINOv2: Learning robust visual features without supervision," *arXiv preprint arXiv:2304.07193*, 2023.
- [22] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, "Masked autoencoders are scalable vision learners," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, New Orleans, LA, USA, pp. 16000–16009, 2022.
- [23] A. Jarca, F.-A. Croitoru, and R. T. Ionescu, "CBM: Curriculum by Masking," *arXiv preprint arXiv:2407.05193*, 2024.
- [24] A. G. Alharthi and S. M. Alzahrani, "Multi-Slice Generation sMRI and fMRI for Autism Spectrum Disorder Diagnosis Using 3D-CNN and Vision Transformers," *Brain Sciences*, vol. 13, no. 11, p. 1578, 2023, doi: 10.3390/brainsci13111578.
- [25] X. Gao and Y. Xu, "Vision transformer and complex network analysis for autism spectrum disorder classification in T1 structural MRI," *Japanese Journal of Radiology*, vol. 43, no. 11, pp. 1788–1802, 2025, doi: 10.1007/s11604-025-01832-3.
- [26] R. K. Goyal and I. Hirano, "The enteric nervous system," *The New England journal of medicine*, vol. 334, no. 17, pp. 1106–1115, 1996, doi: 10.1056/NEJM199604253341707.
- [27] N. Bernardini, C. Segnani, C. Ippolito, R. De Giorgio, R. Colucci, M. S. Faussone-Pellegrini, M. Chiarugi, D. Campani, M. Castagna, L. Mattiù, C. Blandizzi, and A. Dolfi, "Immunohistochemical analysis of myenteric ganglia and interstitial cells of Cajal in ulcerative colitis," *Journal of cellular and molecular medicine*, vol. 16, no. 2, pp. 318-327, 2012, doi: 10.1111/j.1582-4934.2011.01298.x.
- [28] A. E. Bharucha and S. J. H. Brookes, "Neurophysiologic mechanisms of human large intestinal motility," in *Physiology of the Gastrointestinal Tract*, 5th ed., L. R. Johnson, F. K. Ghishan, J. D. Kaunitz, J. L. Merchant, H. M. Said, and J. D. Wood, Eds. Academic Press, 2012, pp. 977-1022, ISBN 9780123820266.
- [29] H. Iwase, S. Sadahiro, S. Mukoyama, H. Makuuchi, and M. Yasuda, "Morphology of myenteric plexuses in the human large intestine: comparison between large intestines with and without colonic diverticula," *Journal of clinical gastroenterology*, vol. 39, no. 8, pp. 674-678, 2005, doi: 10.1097/01.mcg.0000173856.84814.37.
- [30] Y. Megahed, I. Lee, R. Ducharme, A. Erman, O. X. Miguel, K. Dick, A. D. C. Chan, S. Hawken, M. Walker, and F. Moretti, "Deep learning analysis of prenatal ultrasound for identification of ventriculomegaly," *arXiv preprint arXiv:2511.07827*, 2025.