

GastroVision: A Multi-class Endoscopy Image Dataset for Computer Aided Gastrointestinal Disease Detection

Debesh Jha^{*1}, Vanshali Sharma^{*2}, Neethi Dasu³, Nikhil Kumar Tomar¹, Steven Hicks⁴, M.K. Bhuyan², Pradip K. Das², Michael A. Riegler⁴, Pål Halvorsen⁴, Ulas Bagci^{†1}, and Thomas de Lange^{†5}

¹ Department of Radiology, Northwestern University, Chicago, USA
{debesh.jha, nikhil.tomar, ulas.bagci}@northwestern.edu

² Indian Institute of Technology Guwahati, Assam, India
{vanshalisharma, mkb, pkdas}@iitg.ac.in

³ Department of Gastroenterology, Jefferson Health NJ, Cherry Hill, USA
nrdasu4@gmail.com

⁴ SimulaMet, Oslo, Norway
{steven, michael, paalh}@simula.no

⁵ Department of Medicine and Emergencies - Mölndal Sahlgrenska University Hospital, Region Västra Götaland, Sweden & Department of Molecular and Clinical Medicin, Sahlgrenska Academy, University of Gothenburg, Sweden
thomas.de.lange@gu.se

Abstract. Integrating real-time artificial intelligence (AI) systems in clinical practices faces challenges such as scalability and acceptance. These challenges include data availability, biased outcomes, data quality, lack of transparency, and underperformance on unseen datasets from different distributions. The scarcity of large-scale, precisely labeled, and diverse datasets are the major challenge for clinical integration. This scarcity is also due to the legal restrictions and extensive manual efforts required for accurate annotations from clinicians. To address these challenges, we present *GastroVision*, a multi-center open-access gastrointestinal (GI) endoscopy dataset that includes different anatomical landmarks, pathological abnormalities, polyp removal cases and normal findings (a total of 27 classes) from the GI tract. The dataset comprises 8,000 images acquired from Bærum Hospital in Norway and Karolinska University Hospital in Sweden and was annotated and verified by experienced GI endoscopists. Furthermore, we validate the significance of our dataset with extensive benchmarking based on the popular deep learning based baseline models. We believe our dataset can facilitate the development of AI-based algorithms for GI disease detection and classification. Our dataset is available at <https://osf.io/84e7f/>.

Keywords: Medical image · GastroVision · Gastrointestinal diseases.

* These authors contributed equally to this work.

† Shared senior authorship.

1 Introduction

Gastrointestinal (GI) cancers account for 26% of cancer incidences and 35% of cancer-related deaths worldwide. In 2018, there were approximately 4.8 million new cases of GI cancer and 3.4 million deaths [7]. The five major types of GI cancers are colo-rectal (1.93 million cases; third most common cancer), pancreas (466,003 deaths; lowest survival rate), liver (905,677 cases), stomach (1.09 million cases), and esophagus (604,100 cases) [13]. These cancer cases are predicted to increase by 58%, and related deaths could show a 73% rise by 2040 [7]. Early detection of such cancers and their precursors can play an important role in improving the outcome and make the treatment less invasive. Some of the common examinations performed for GI cancer detection include endoscopy or esophagogastroduodenoscopy (EGD), capsule endoscopy, colonoscopy, imaging studies (MRI, X-ray, ultrasound, CT scan, or PET scan) or endoscopic ultrasound (EUS). Endoscopy is widely accepted as the gold standard for detecting abnormalities of the bowel lumen and mucosa, upper endoscopy for esophagus, stomach, and the duodenum and colonoscopy for the large bowel and rectum GI tract for abnormalities, respectively.

The endoscopies are performed by nurses or doctor endoscopists. The assessment of the endoscopy examinations is operator dependent, and the assessment and therapeutic decision vary between endoscopists. Consequently, the quality and accuracy of detection and diagnosis of lesions are attributed to the level of the operator skills and efforts of the endoscopists. Despite various measures taken to provide guidance for operators, significant lesion miss rates are still reported. For example, there is evidence of colon polyp miss rates of up to 27% due to polyps and operator characteristics [2,20]. Considering the shortcomings of the manual review process, various automated systems are adopted to provide AI-based real-time support to clinicians to reduce lesion miss rates and misinterpretation of lesions to ultimately increase detection rates. A microsimulation study reports a 4.8% incremental gain in the reduction of colorectal cancer incidence when colonoscopy screening was combined with AI tools [6]. Such findings motivated research work in healthcare to adopt AI as a potential tool for GI cancer detection. Gastric cancer, inflammatory bowel disease (IBD), and esophageal neoplasia are some of the GI tract findings already being investigated using AI techniques [1]. Despite AI being adopted in some hospitals for clinical applications, the integration of AI into the extensive clinical setting is still limited. Integrating AI techniques with regular clinical practices is multifactorial and poses serious concerns regarding its implementation in large-scale healthcare systems. One of the significant factors is the algorithmic bias, which worsens when the system learns from the annotations handled by a single, non-blinded endoscopist who may have personal thresholds to label the findings.

Moreover, most existing AI models depend on data acquired from a single center, which makes them less valid when faced with a varied patient population. This leads to spectrum bias under which AI systems encounter performance drops due to the significant shift in the original clinical context and the test sample population. In such cases, unexpected outcomes and diagnostic accuracy

could be obtained using automated tools. Such bias issues could reach the clinical systems at any point of the process, including data collection, study design, data entry and pre-processing, algorithm design, and implementation. The very beginning of the process, i.e., data collection, is of utmost importance for reproducibility and to perform validations on images from a diverse population, different centers, and imaging modalities. To develop scalable healthcare systems, it is vital to consider these challenges and perform real-time validations. However, the scarcity of comprehensive data covering a range of real-time imaging scenarios arising during endoscopy or colonoscopy makes it difficult to develop a robust AI-based model. Although much progress has been made on automated cancer detection and classification [16,19], it is still challenging to adapt such models into real-time clinical settings as they are tested on small-sized datasets with limited classes.

Some classes in the dataset could be scarce because some conditions or diseases occur less often. Consequently, such findings are not frequently captured and remain unexplored despite requiring medical attention. AI-based detection of these findings, even with a small sample count, can significantly benefit from techniques like one-shot or few-shot learning. These techniques allow the AI models to learn patterns and features indicative of the condition, thus, enabling accurate diagnosis with minimal training data. Apart from the above-mentioned limitations, many existing datasets are available on request, and prior consenting is required, which delays the process and does not guarantee accessibility. Therefore, in this paper, we publish *GastroVision*, an open-access multi-class endoscopy image dataset for automated GI disease detection that does not require prior consenting and can be downloaded easily with a single click. The data covers a wide range of classes that can allow initial exploration of many anatomical landmarks and pathological findings.

The main contributions of this work are summarized below:

- We present an open-access multi-class GI endoscopy dataset, namely, *Gastrovision*, containing 8,000 images with 27 classes from two hospitals in Norway and Sweden. The dataset exhibits a diverse range of classes, including anatomical landmarks, pathological findings, polyp removal cases and normal or regular findings. It covers a wide range of clinical scenarios encountered in endoscopic procedures.
- We evaluated a series of deep learning baseline models on standard evaluation metrics using our proposed dataset. With this baseline, we invite the research community to improve our results and develop novel GI endoscopy solutions on our comprehensive set of GI finding classes. Additionally, we encourage computer vision and machine learning researchers to validate their methods on our open-access data for a fair comparison. This can aid in developing state-of-the-art solutions and computer-aided systems for GI disease detection and other general machine learning classification tasks.

Table 1. List of the existing datasets within GI endoscopy.

Dataset	Data type	Size	Accessibility
Kvasir-SEG [17]	Polyps	1,000 images [†] [*]	Public
HyperKvasir [11]	GI findings	110,079 images & 374 videos	Public
Kvasir-Capsule [24]	GI findings [◊]	4,741,504 images	Public
Kvasir [22]	GI findings	8,000 images	Public
CVC-ColonDB [10]	Polyps	380 images [†] [‡]	As per request [*]
ETIS-Larib Polyp DB [23]	Polyps	196 images [†]	Public
EDD2020 [4,3]	GI lesions	386 images [†] [*]	Public
CVC-ClinicDB [9]	Polyps	612 images [†]	Public
CVC-VideoClinicDB [8]	Polyps	11,954 images [†]	As per request
ASU-Mayo polyp database [25]	Polyps	18,781 images [†]	As per request [*]
KID [18]	Angiectasia, bleeding, inflammations [◊]	> 2500 images, 47 videos	Public [*]
PolypGen [5]	Polyps	1,537 images [†] [*] & 2,225 video sequence, 4,275 negative frame	Public
SUN Database [21]	Polyps	158,690 video frames [*]	As per request
GastroVision (ours)	GI findings	8,000 image frames	Public

[†]Segmentation ground truth ^{*}Not available now [‡]Contour [◊]Video capsule endoscopy ^{*} Bounding box information

2 Related Work

Table 1 shows the list of the existing dataset along with data type, their size, and accessibility. It can be observed that most of the existing datasets in the literature are from colonoscopy procedures and consist of polyps still frames or videos. These are mostly used for segmentation tasks. Most of the existing datasets are small in size and do not capture some critical anatomical landmarks or pathological findings. In the earlier GI detection works, the CVC-ClinicDB [9] and CVC-ColonDB [10] were widely used. **CVC-ClinicDB** is developed from 23 colonoscopy video studies acquired with white light. These videos provide 31 video sequences, each containing one polyp, which finally generates 612 images of size 576×768 . **CVC-ColonDB** consists of 300 different images obtained from 15 random cases. Similarly, **ETIS-Larib Polyp DB** [23] is a colonoscopy dataset consisting of 196 polyp frames and their corresponding segmentation masks. Recently, **Kvasir-SEG** [17] dataset has been introduced that comprises of 1,000 colonoscopy images with segmentation ground truth and bounding box coordinate details. This dataset offers a diverse range of polyp frames, including multiple diminutive polyps, small-sized polyps, regular polyps, sessile or flat polyps collected from varied cohort populations. The dataset is open-access and is one of the most commonly used datasets for automatic polyp segmentation.

The **ASU-Mayo Clinic Colonoscopy Video (c) database** [25] is a copy-righted dataset and is considered the first largest collection of short and long video sequences. Its training set is composed of 10 positive shots with polyps inside and 10 negative shots with no polyps. The associated test set is provided with 18 different unannotated videos. **CVC-VideoClinicDB** [8] is extracted from more than 40 long and short video sequences. Its training set comprises 18 different sequences with an approximate segmentation ground truth and Paris

classification for each polyp. **SUN Colonoscopy Video Database** comprises 49,136 polyp frames and 109,554 non-polyp frames. Unlike the datasets described above, this dataset includes pathological classification labels, polyp size, and shape information. It also includes bounding box coordinate details. The **Polyp-Gen** [5] dataset is an open-access dataset that comprises 1,537 polyp images, 2,225 positive video sequences, and 4,275 negative frames. The dataset is collected from six different centers in Europe and Africa. Altogether the dataset provides 3,762 positive frames and 4,275 negative frames. These still images and video frames are collected from varied populations, endoscopic systems, and surveillance expert in Norway, France, United Kingdom, Egypt, and Italy and is one of the comprehensive open-access datasets for polyp detection and segmentation.

Apart from the lower GI-related datasets, there are a few datasets that provide combined samples of upper and lower GI findings. For example, **HyperKvasir** [11] is a multi-class GI endoscopy dataset that covers 23 classes of anatomical landmarks. It contains 110,079 images out of which 10,662 are labeled and 99,417 are unlabeled images. The **EDD2020** dataset [4,3] is a collection of five classes and 386 still images with detection and segmentation ground truth. The classes are divided into 160 non-dysplastic Barrett’s, 88 suspicious precancerous lesions, 74 high-grade dysplasia, 53 cancer, and 127 polyps with overall 503 ground truth annotations. The **Kvasir-Capsule** [24] is a video capsule endoscopy dataset comprising 4,741,504 image frames extracted from 117 videos. From the total frames, 4,694,266 are unlabeled, and 47,238 frames are annotated with a bounding box for each of the 14 classes. Similarly, **KID** [18] is a capsule endoscopy dataset with 47 videos and over 2,500 images. The images are annotated for normal, vascular, inflammatory, lymphangiectasias, and polypoid lesions.

The literature review shows that most GI-related datasets focus on a single specific finding, such as colon polyps. Some of the datasets are small in size and have ignored non-lesion frames, which are essential for developing algorithms to be integrated into clinical settings. Additionally, many of these datasets are available on request and require approval from the data providers resulting in further delays. A few datasets like Kvasir, HyperKvasir, Kvasir-Capsule and KID provide multiple GI findings. However, Kvasir-Capsule and KID are video capsule endoscopy datasets. The Kvasir dataset has only eight classes, whereas Hyperkvasir has 23 classes. In contrast, our *GastroVision* dataset has 27 classes and covers more labeled classes of anatomical landmarks, pathological findings, and normal findings. Additionally, we establish baseline results on this dataset for GI disease detection and classification, offering valuable research resources for advancing GI endoscopy studies.

3 GastroVision

Here, we provide detailed information about the dataset, acquisition protocol, ethical and privacy aspects of data and suggested metrics.

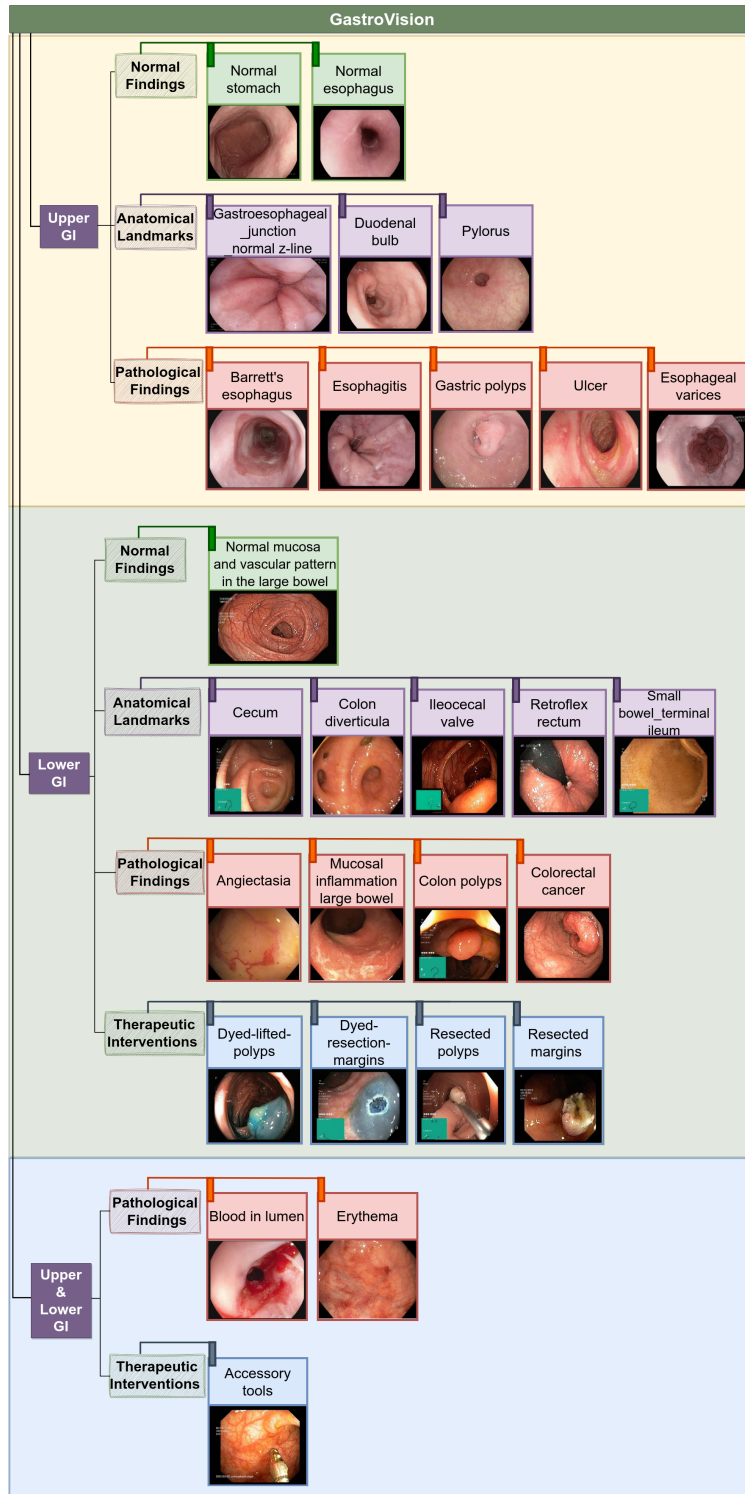


Fig. 1. Example images from the gastrointestinal tract showing distinct findings from the upper and lower GI tract.

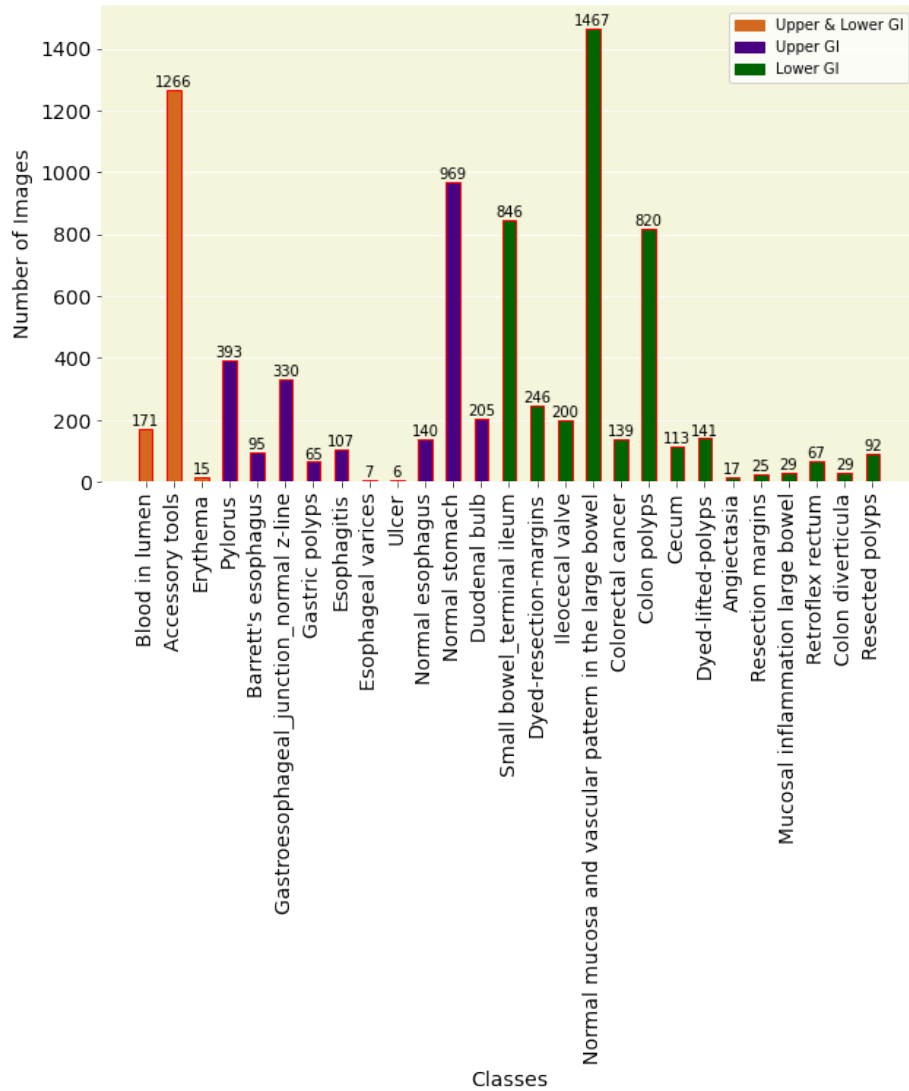


Fig. 2. The figure shows the number of images per class. Some classes have few samples because of the rarity of the findings and the technical challenges associated with obtaining such samples in endoscopic settings.

3.1 Dataset details

GastroVision is an open-access dataset that incorporates 8,000 images pertaining to 27 different labeled classes (Fig. 1). Most images are obtained through White Light Imaging (WLI), while a few samples are acquired using Narrow Band Imaging (NBI). These classes are categorized into two broad categories:

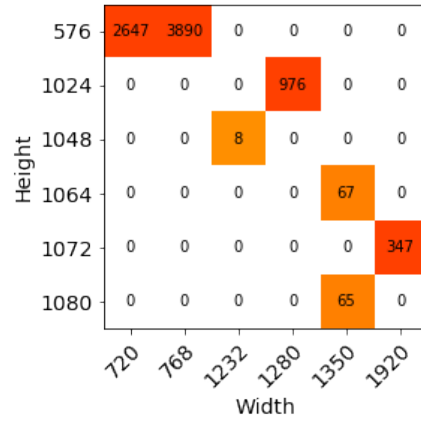


Fig. 3. Resolutions of the 8,000 images of GastroVision.

Upper GI tract and *Lower GI tract*. The number of images under each class is presented in Fig. 2. These classes indicate findings acquired from GI tract. It can be observed that the sample count is not balanced across classes, which is generally experienced in the medical image acquisition process as some findings occur less often. Releasing these classes in the dataset will allow the researchers to leverage the fast-emerging AI techniques to develop methods for detecting such rare but clinically significant anomalies. All the images are stored in JPG format, and the overall size is around 1.8 GB. The resolution details of the images can be found in Fig. 3. *GastroVision* is provided as a structured directory, with each class having a specific folder. For example, the ‘*Accessory tools*’ folder contains all images featuring diagnostic and therapeutic tools.

Upper GI Tract: Upper GI endoscopy examines the esophagus, stomach, and duodenum. The various classes covered in this GI tract are discussed below as three subcategories: *normal findings*, *anatomical landmarks*, and *pathological findings*. A detailed categorization is shown in Fig. 1. The **normal stomach** serves as a critical site for initial digestion, while the **duodenal bulb**, the first part of the small intestine, is critical for nutrient absorption. Anatomical landmarks are used as reference points to indicate a specific location and assist in navigating during endoscopy procedures. The **gastroesophageal junction** is an anatomical area where **esophagus** joins the **stomach** also alining to the **normal z-line**, a transitional point where the esophagus’s squamous epithelium and the stomach’s columnar mucosa lining join. **Pylorus** is a sphincter connecting the stomach and the duodenum, the first part of the small intestine.

Apart from these anatomical landmarks, any pathological conditions may be encountered during endoscopy. **Esophagitis**, the most common abnormality, is characterized by an inflammation of the esophagus. This disease is graded based on its severity according to the Los Angeles classification. For example, grade

B refers to the condition when the mucosal break is limited to the mucosal fold and is more than 5 mm long. In grade D, mucosal break affects 75% of the esophageal circumference. Long standing esophagitis may cause **Barett's esophagus**, a condition in which the cells of the esophagus's lining start to change, and tissues appear red. This is a precancerous condition. Other frequent lesions observed are **polyps**, abnormal tissue growth or ulcers. **Gastric polyps** are abnormal growths in the stomach lining. **Ulcers** are the open sores in the stomach or duodenum that can lead to discomfort and bleeding. **Esophageal varices** result from portal hypertension, causing swollen veins in the esophagus. **Erythema** refers to redness, often indicating inflammation and **blood in the lumen** denotes bleeding. **Accessory tools** aid in investigating and diagnosing upper and lower GI tract conditions for targeted treatment.

Lower GI Tract: The lower GI tract is examined by colonoscopy to investigate any abnormalities in the colon, the rectum, and the terminal ileum (the last part of the small bowel). Here, we covered one more subcategory, *therapeutic interventions*, in addition to *normal findings*, *anatomical landmarks*, and *pathological findings*. A detailed class-wise division is shown in Fig. 1.

The **normal mucosa and vascular pattern in the large bowel** is essential for absorbing water and electrolytes. The different anatomical landmarks associated with lower GI include **cecum** (first part of the large intestine), visualizing the appendiceal orifice, **ileocecal valve** (sphincter muscle between ileum and colon), and the **small bowel**. During the colonoscopy, these anatomical landmarks act as reference points to prove complete examination. Retroflexion in the rectum is performed to visualize a blind zone, using the bending section of the colonoscope to visualize the distal area of the colon, called **retroflex-rectum**. The **terminal ileum**, the last part of the small intestine, aids in nutrient absorption. **Colon diverticula**, small pouch-like protrusions, can form along the colon's weakened wall, often in the sigmoid colon [12].

During the colonoscopy, the endoscopist navigates through these landmarks and looks for abnormalities such as **polyps**, **angiectasia**, and inflammation like **ulcerative colitis**. **Angiectasia** is a common lesion representing abnormal blood vessels and is responsible for obscure recurrent lower GI bleeding. These can easily be distinguished from the **normal vessels** shown in Fig. 1. **Colorectal cancer** occurs in the colon or rectum. One of the early signs of this colorectal cancer can be detected through **colon polyps**. **Mucosal inflammation in the large bowel** may be caused by different factors, such as infections or chronic inflammatory conditions.

Apart from the aforementioned pathological conditions, several therapeutic interventions are adopted to treat the detected anomalies effectively. It frequently involves the removal of the lesion/polyp. The surrounding of such **resected polyps**, also called the **resection margins** or resection sites, are then considered for biopsies. To enhance lesion demarcation, a solution containing indigo carmine is injected, making resection easier. The appearance of blue color underneath the **dyed-lifted-polyp** provides accurate polyp margins. After resecting

such polyps, the underlying region, known as **dyed-resection-margin**, appears blue. These margins are important to examine for any remaining tissue of the resected polyp.

3.2 Dataset acquisition, collection and construction

Data acquisition and collection: The dataset images are acquired from two centers (Department of Gastroenterology, Bærum Hospital, Vestre Viken Hospital Trust (VV), Norway and Karolinska University Hospital, Stockholm, Sweden) using standard endoscopy equipment from Olympus (Olympus Europe, Germany) and Pentax (Pentax Medical Europe, Germany). A team of expert gastroenterologists, one junior doctor, and two computational scientists were involved in the labelling of the images and the related review process. It is worth noting that for dataset collection, we labeled some of the unlabeled images from the HyperKvasir dataset and included them in our dataset. Additionally, we labeled the images acquired from the Karolinska University Hospital to their respective classes for developing a diverse and multi-center ‘‘GastroVision’’ dataset.

Ethical and privacy aspects of the data: The dataset is constructed while preserving the patients’ anonymity and privacy. All videos and images from Bærum hospitals were fully anonymized, following the GDPR requirements for full anonymization. Hence, it is exempted from patient consent. The files were renamed using randomly generated filenames. The Norwegian Privacy Data Protection Authority approved this export of anonymized images for research purposes. As the dataset development procedure involved no interference with the medical treatment or care of the patient, it has also been granted an exemption for approval by Regional Committee for Medical and Health Research Ethics - South East Norway. Similarly, the data collection process at Karolinska University Hospital, Sweden, is completely anonymized as per the GDPR requirements.

3.3 Suggested metrics

Standard multi-class classification metrics, such as Matthews Correlation Coefficient (MCC), micro and macro averages of recall/sensitivity, precision, and F1-score, can be used to validate the performance using our dataset. MCC provides a balanced measure even in cases with largely varying class sizes. A macro-average will compute the metric independently for each class and then take the average, whereas a micro-average will aggregate the contributions of all classes to compute the metric. Recall presents the ratio of correctly predicted positive observations to all the original observations in the actual class. Precision is the ratio of correctly predicted positive observations to all the positive predicted observations. F1-score integrates both recall and precision and calculates a weighted average/harmonic mean of these two metrics.

Table 2. Results for all classification experiments on the Gastrovision dataset.

Method	Macro Average			Micro Average			MCC
	Prec.	Recall	F1	Prec.	Recall	F1	
ResNet-50 [14]	0.4373	0.4379	0.4330	0.6816	0.6816	0.6816	0.6416
Pre-trained ResNet-152 [14]	0.5258	0.4287	0.4496	0.6879	0.6879	0.6879	0.6478
Pre-trained EfficientNet-B0 [26]	0.5285	0.4326	0.4519	0.6759	0.6759	0.6759	0.6351
Pre-trained DenseNet-169 [15]	0.6075	0.4603	0.4883	0.7055	0.7055	0.7055	0.6685
Pre-trained ResNet-50 [14]	0.6398	0.6073	0.6176	0.8146	0.8146	0.8146	0.7921
Pre-trained DenseNet-121 [15]	0.7388	0.6231	0.6504	0.8203	0.8203	0.8203	0.7987

4 Experiments and Results

In this section, we describe the implementation details, technical validation and the limitation of the dataset.

4.1 Implementation Details

All deep learning diagnostic models are trained on NVIDIA TITAN Xp GPU using PyTorch 1.12.1 framework. A stratified sampling is performed to preserve the similar distribution of each class during 60:20:20 training, validation, and testing split formation. The images are resized to 224×224 pixels, and simple data augmentations, including random rotation and random horizontal flip, are applied. All models are configured with similar hyperparameters, and a learning rate of $1e^{-4}$ is initially set with 150 epochs. An Adam optimizer is used with the *ReduceLROnPlateau* scheduler. More description about the implementation details and dataset can be found on our GitHub page [†].

4.2 Technical Validation

To evaluate the presented data for technical quality and classification tasks, we performed a series of experiments using some state-of-the-art deep learning models. The purpose of this preliminary validation is to provide baseline results that can be referred to for comparison by future researchers. We carried out multi-class classification using CNN-based models, namely, ResNet-50 [14], ResNet-152 [14], EfficientNet-B0 [26], DenseNet-121 [15], and DenseNet-169 [15], considering their competent performance in GI-related image-based tasks in the literature [27]. Note that we have only included classes with more than 25 samples in the experiments, which resulted in 22 classes in total. However, we also release the other classes with fewer samples to welcome new interesting findings in areas similar to one-shot learning.

The different experiments performed include (a) *ResNet-50*: The model is randomly initialized, and an end-to-end training is done, (b) *Pre-trained ResNet-50* and (c) *Pre-trained DenseNet-121*: The models are initialized with pre-trained weights, and then all layers are fine-tuned, (d) *Pre-trained ResNet-152*,

[†] <https://github.com/DebeshJha/GastroVision>

Class	Pre-trained DenseNet-121			Support
	Precision	Recall	F1-score	
Accessory tools	0.93	0.96	0.95	253
Barrett’s esophagus	0.55	0.32	0.4	19
Blood in lumen	0.86	0.91	0.89	34
Cecum	0.33	0.17	0.23	23
Colon diverticula	1	0.33	0.5	6
Colon polyps	0.78	0.87	0.82	163
Colorectal cancer	0.63	0.41	0.5	29
Duodenal bulb	0.72	0.76	0.74	41
Dyed-lifted-polyps	0.86	0.86	0.86	28
Dyed-resection-margins	0.94	0.92	0.93	49
Esophagitis	0.5	0.23	0.31	22
Gastric polyps	0.6	0.23	0.33	13
Gastroesophageal_junction_normal z-line	0.65	0.85	0.74	66
Ileocecal valve	0.74	0.7	0.72	40
Mucosal inflammation large bowel	1	0.33	0.5	6
Normal esophagus	0.72	0.82	0.77	28
Normal mucosa and vasular pattern in the large bowel	0.81	0.87	0.84	293
Normal stomach	0.9	0.86	0.88	194
Pylorus	0.8	0.92	0.86	78
Resected polyps	0.33	0.11	0.17	18
Retroflex rectum	0.75	0.43	0.55	14
Small bowel_terminal ileum	0.86	0.85	0.85	169

Table 3. Class-wise performance associated with the best outcome obtained using pre-trained DenseNet-121.

(e) *Pre-trained EfficientNet-B0* and (f) *Pre-trained DenseNet-169*: The models are initialized with pre-trained weights, and only the updated last layer is fine-tuned. All the above pre-trained models use ImageNet weights. The associated results are shown in Table 2. It can be observed that the best outcome is obtained using the pre-trained DenseNet-121. A class-wise analysis using the same model is provided in Table 3 and Fig. 4. It shows that while most classes achieved satisfactory prediction outcomes, a few proved to be very challenging for the classification model. For a more detailed analysis, we plotted a two-dimensional t-SNE embedding for *GastroVision* (Fig. 5). The classes like *Normal stomach*, *Dyed-resection-margins*, which present a clear distinction in the t-SNE embedding, are less often misclassified. The above points could be the reasons for the F1-score of 0.88 and 0.93 in the case of *Dyed-resection-margins* and *Normal stomach* classes, respectively. On the other hand, there are some overlapping classes such as *Cecum* and *Normal mucosa and vascular pattern in the large bowel* or *Colorectal cancer* and *Colon polyps* which do not present clear demarcation with each other and hence, are likely to be misclassified.

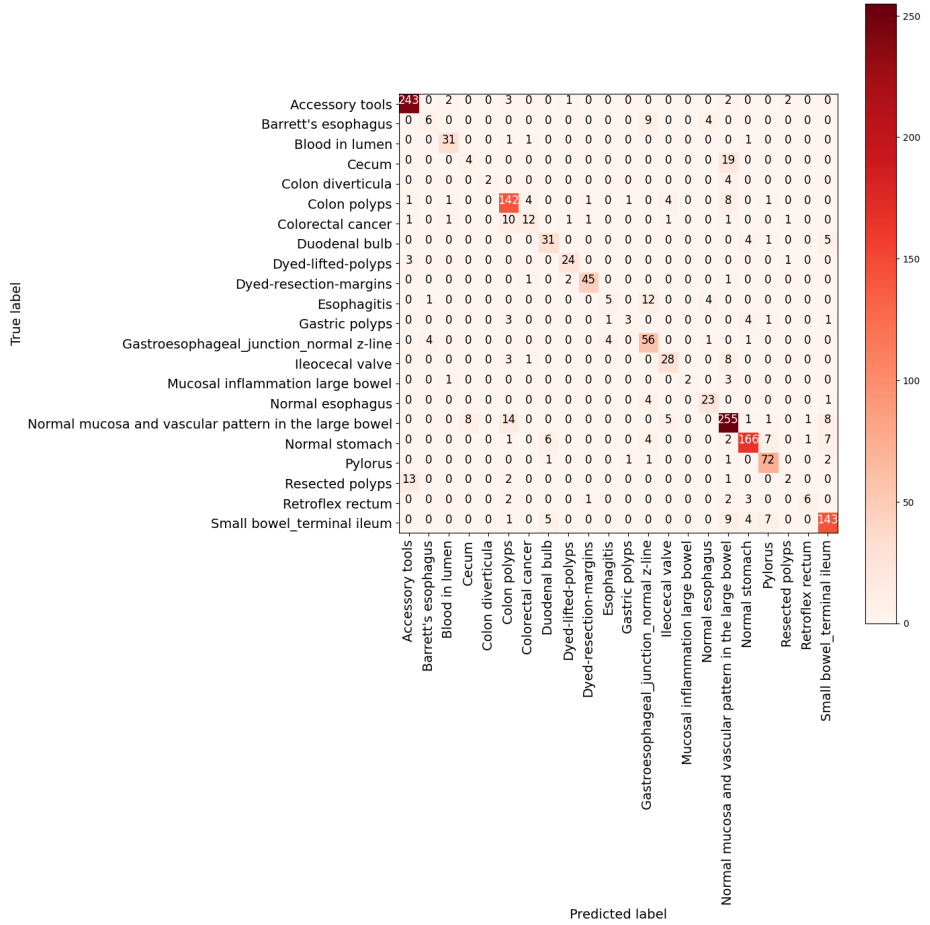


Fig. 4. Confusion matrix for the best outcome obtained using pre-trained DenseNet-121.

Considering the overall results and many overlapping classes (without distinct clustering), it can be inferred that classifying GI-related anatomical landmarks and pathological findings is very challenging. Many abnormalities are hard to differentiate, and the rarely occurring findings have higher chances of getting misclassified. This presents the challenge of developing a robust AI system that could address multiple aspects important for GI image classification, e.g., many findings are subtle and difficult to be identified, and some findings are not easily acquired during the endoscopy procedure, which results in less number of data samples. Such underrepresented classes need to be explored with some specific algorithms specially designed to leverage the availability of a few hard-to-find samples. Thus, the potential of the baseline results and associated issues and challenges motivate the need to publish this dataset for further investigations.

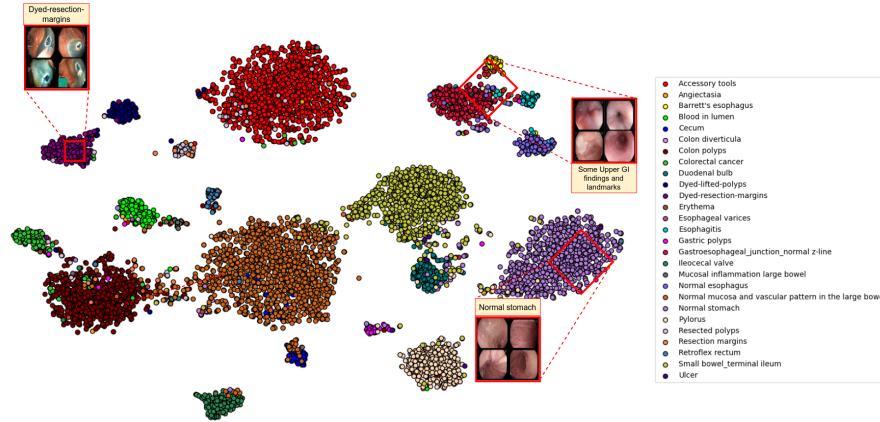


Fig. 5. Two-dimensional t-SNE embedding for GastroVision. The pre-trained DenseNet-121 model, which is further trained on our training set, is used to extract features. Some sample images are shown with either a specific or a broader (due to multiple overlapping classes) categorization.

4.3 Limitation of the dataset

Our dataset, *GastroVision*, is a unique and diverse dataset with the potential to explore a wide range of anatomical and pathological findings using automated diagnosis. Although this labeled image data can enable the researchers to develop methods to detect GI-related abnormalities and other landmarks, it lacks segmented annotations in the current version, which could further enhance the treatment experience and surgical procedures. It is important to note that some classes (for example, colon diverticula, erythema, cecum, esophagitis, esophageal varices, ulcer and pylorus) have only a few images. Despite this limitation, our dataset is well suited for one-shot and few-shot learning approaches to explore some GI-related conditions that have still not received attention in medical image analysis. In the future, we plan to extend the dataset by including more classes and a larger number of samples, along with ground truth for some of the classes that could be used for segmentation purposes as well as images with higher resolution from the most recent endoscopy systems.

5 Conclusion

In this paper, we presented a new multi-class endoscopy dataset, *GastroVision*, for GI anomalies and disease detection. We have made the dataset available for the research community along with the implementation details of our method.

The labeled image data can allow researchers to formulate methodologies for classifying different GI findings, such as important pathological lesions, endoscopic polyp removal cases, and anatomical landmarks found in the GI tract. We evaluated the dataset using some baseline models and standard multi-class classification metrics. The results motivate the need to investigate better specific techniques for GI-related data. Having a diverse set of categories labeled by expert endoscopists from two different centers, *GastroVision* is unique and valuable for computer-aided GI anomaly and disease detection, patient examinations, and medical training.

Acknowledgements

D. Jha is supported by the NIH funding: R01-CA246704 and R01-CA240639.V. Sharma is supported by the INSPIRE fellowship (IF190362), DST, Govt. of India.

References

1. Abadir, A.P., Ali, M.F., Karnes, W., Samarasena, J.B.: Artificial intelligence in gastrointestinal endoscopy. *Clinical endoscopy* **53**(2), 132–141 (2020)
2. Ahn, S.B., Han, D.S., Bae, J.H., Byun, T.J., Kim, J.P., Eun, C.S.: The miss rate for colorectal adenoma determined by quality-adjusted, back-to-back colonoscopies. *Gut and liver* **6**(1), 64 (2012)
3. Ali, S., Dmitrieva, M., Ghatwary, N., Bano, S., Polat, G., Temizel, A., Krenzer, A., Hekalo, A., Guo, Y.B., Matuszewski, B., Gridach, M., Voiculescu, I., Yoganand, V., Chavan, A., Raj, A., Nguyen, N.T., Tran, D.Q., Huynh, L.D., Boutry, N., Rezvy, S., Chen, H., Choi, Y.H., Subramanian, A., Balasubramanian, V., Gao, X.W., Hu, H., Liao, Y., Stoyanov, D., Daul, C., Realdon, S., Cannizzaro, R., Lamarque, D., Tran-Nguyen, T., Bailey, A., Braden, B., East, J.E., Rittscher, J.: Deep learning for detection and segmentation of artefact and disease instances in gastrointestinal endoscopy. *Medical Image Analysis* **70**, 102002 (2021)
4. Ali, S., Ghatwary, N., Braden, B., Lamarque, D., Bailey, A., Realdon, S., Cannizzaro, R., Rittscher, J., Daul, C., East, J.: Endoscopy disease detection challenge 2020. arXiv preprint arXiv:2003.03376 (2020)
5. Ali, S., Jha, D., Ghatwary, N., Realdon, S., Cannizzaro, R., Salem, O.E., Lamarque, D., Daul, C., Riegler, M.A., Anonsen, K.V., et al.: A multi-centre polyp detection and segmentation dataset for generalisability assessment. *Scientific Data* **10**(1), 75 (2023)
6. Areia, M., Mori, Y., Correale, L., Repici, A., Bretthauer, M., Sharma, P., Taveira, F., Spadaccini, M., Antonelli, G., Ebigbo, A., et al.: Cost-effectiveness of artificial intelligence for screening colonoscopy: a modelling study. *The Lancet Digital Health* **4**(6), e436–e444 (2022)
7. Arnold, M., Abnet, C.C., Neale, R.E., Vignat, J., Giovannucci, E.L., McGlynn, K.A., Bray, F.: Global burden of 5 major types of gastrointestinal cancer. *Gastroenterology* **159**(1), 335–349 (2020)
8. Bernal, J., Aymeric, H.: MICCAI endoscopic vision challenge polyp detection and segmentation (2017)

9. Bernal, J., Sánchez, F.J., Fernández-Esparrach, G., Gil, D., Rodríguez, C., Vilarino, F.: Wm-dova maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from physicians. *Computerized Medical Imaging and Graphics* **43**, 99–111 (2015)
10. Bernal, J., Sánchez, J., Vilarino, F.: Towards automatic polyp detection with a polyp appearance model. *Pattern Recognition* **45**(9), 3166–3182 (2012)
11. Borgli, H., Thambawita, V., Smedsrud, P.H., Hicks, S., Jha, D., Eskeland, S.L., Randel, K.R., Pogorelov, K., Lux, M., Nguyen, D.T.D., et al.: Hyperkvasir, a comprehensive multi-class image and video dataset for gastrointestinal endoscopy. *Scientific Data* **7**(1), 1–14 (2020)
12. Crafa, P., Diaz-Cano, S.J.: Changes in colonic structure and mucosal inflammation. In: *Colonic Diverticular Disease*, pp. 41–61 (2022)
13. Globocan: Cancer today (2020), <https://gco.iarc.fr/today/fact-sheets-cancers>
14. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*. pp. 770–778 (2016)
15. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 4700–4708 (2017)
16. Jha, D., Ali, S., Tomar, N.K., Johansen, H.D., Johansen, D., Rittscher, J., Riegler, M.A., Halvorsen, P.: Real-time polyp detection, localization and segmentation in colonoscopy using deep learning. *IEEE Access* **9**, 40496–40510 (2021)
17. Jha, D., Smedsrud, P.H., Riegler, M.A., Halvorsen, P., Lange, T.d., Johansen, D., Johansen, H.D.: Kvasir-SEG: a segmented polyp dataset. In: *Proceedings of the International Conference on Multimedia Modeling (MMM)*. pp. 451–462 (2020)
18. Koulaouzidis, A., Iakovidis, D.K., Yung, D.E., Rondonotti, E., Kopylov, U., Plevris, J.N., Toth, E., Eliakim, A., Johansson, G.W., Marlicz, W., et al.: Kid project: an internet-based digital video atlas of capsule endoscopy for research purposes. *Endoscopy international open* **5**(6), E477 (2017)
19. Li, K., Fathan, M.I., Patel, K., Zhang, T., Zhong, C., Bansal, A., Rastogi, A., Wang, J.S., Wang, G.: Colonoscopy polyp detection and classification: Dataset creation and comparative evaluations. *arXiv preprint arXiv:2104.10824* (2021)
20. Mahmud, N., Cohen, J., Tsourides, K., Berzin, T.M.: Computer vision and augmented reality in gastrointestinal endoscopy. *Gastroenterology report* **3**(3), 179–184 (2015)
21. Misawa, M., Kudo, S.e., Mori, Y., Hotta, K., Ohtsuka, K., Matsuda, T., Saito, S., Kudo, T., Baba, T., Ishida, F., et al.: Development of a computer-aided detection system for colonoscopy and a publicly accessible large colonoscopy video database (with video). *Gastrointestinal endoscopy* **93**(4), 960–967 (2021)
22. Pogorelov, K., Randel, K.R., Griwodz, C., Eskeland, S.L., de Lange, T., Johansen, D., Spampinato, C., Dang-Nguyen, D.T., Lux, M., Schmidt, P.T., et al.: Kvasir: A multi-class image dataset for computer aided gastrointestinal disease detection. In: *Proceedings of the 8th ACM on Multimedia Systems Conference*. pp. 164–169 (2017)
23. Silva, J., Histace, A., Romain, O., Dray, X., Granado, B.: Toward embedded detection of polyps in wce images for early diagnosis of colorectal cancer. *International journal of computer assisted radiology and surgery* **9**(2), 283–293 (2014)
24. Smedsrud, P.H., Thambawita, V., Hicks, S.A., Gjestang, H., Nedrejord, O.O., Næss, E., Borgli, H., Jha, D., Berstad, T.J.D., Eskeland, S.L., et al.: Kvasir-capsule, a video capsule endoscopy dataset. *Scientific Data* **8**(1), 1–10 (2021)

25. Tajbakhsh, N., Gurudu, S.R., Liang, J.: Automated polyp detection in colonoscopy videos using shape and context information. *IEEE transactions on medical imaging* **35**(2), 630–644 (2015)
26. Tan, M., Le, Q.: Efficientnet: Rethinking model scaling for convolutional neural networks. In: *Proceedings of the International conference on machine learning*. pp. 6105–6114 (2019)
27. Thambawita, V., Jha, D., Riegler, M., Halvorsen, P., Hammer, H.L., Johansen, H.D., Johansen, D.: The medico-task 2018: Disease detection in the gastrointestinal tract using global features and deep learning. In: *Proceedings of the MediaEval 2018 Workshop* (2018)