

# Patch-level instance-group discrimination with pretext-invariant learning for colitis scoring

Ziang Xu<sup>1,2</sup>, Sharib Ali<sup>1,3</sup>, Soumya Gupta<sup>1,2</sup>, Simon Leedham<sup>4,5</sup>, James E East<sup>4,5</sup>, Jens Rittscher<sup>1,2,4</sup>

<sup>1</sup> Institute of Biomedical Engineering, University of Oxford, Oxford, UK

<sup>2</sup> Big Data Institute, University of Oxford, Li Ka Shing Centre for Health Information and Discovery, Oxford, UK

<sup>3</sup> School of Computing, University of Leeds, Leeds, UK

<sup>4</sup> NIHR Oxford Biomedical Research Centre, Oxford, UK

<sup>5</sup> Translational Gastroenterology Unit, Experimental Medicine Div., John Radcliffe Hospital, University of Oxford, Oxford, UK

**Abstract.** Inflammatory bowel disease (IBD), in particular ulcerative colitis (UC), is graded by endoscopists and this assessment is the basis for risk stratification and therapy monitoring. Presently, endoscopic characterisation is largely operator dependant leading to sometimes undesirable clinical outcomes for patients with IBD. We focus on the Mayo Endoscopic Scoring (MES) system which is widely used but requires the reliable identification of subtle changes in mucosal inflammation. Most existing deep learning classification methods cannot detect these fine-grained changes which make UC grading such a challenging task. In this work, we introduce a novel patch-level instance-group discrimination with pretext-invariant representation learning (PLD-PIRL) for self-supervised learning (SSL). Our experiments demonstrate both improved accuracy and robustness compared to the baseline supervised network and several state-of-the-art SSL methods. Compared to the baseline (ResNet50) supervised classification our proposed PLD-PIRL obtained an improvement of 4.75% on hold-out test data and 6.64% on unseen center test data for top-1 accuracy.

**Keywords:** Colonoscopy · Inflammation · Self-supervised learning · Classification · Group discrimination · Colitis

## 1 Introduction

Ulcerative colitis (UC) is a chronic intestinal inflammatory disease in which lesions such as inflammation and ulcers are mainly located in the colon and rectum. UC is more common in early adulthood, the disease lasts for a long time, and the possibility of further cancerous transformation is high [15]. It is therefore important to diagnose ulcerative colitis early. Colonoscopy is a gold standard clinical procedure widely used for early screening of disease. Among the various colonoscopic evaluation methods proposed, the Mayo Endoscopic Score (MES) is considered to be the most widely used evaluation indicators to measure the

UC activity [16,6]. MES divides UC into three categories, namely mild (MES-1), moderate (MES-2) and severe (MES-3). MES-2 and MES-3 indicate that an immediate follow up is required. However, the grading of UC in colonoscopy is dependent on the level of experience. Differences in assessment amongst endoscopists have been observed that can affect patient management. Automated systems based on artificial intelligence can help identify subtle abnormalities that represent UC, improve diagnostic quality and minimise subjectivity.

Deep learning models based on Convolutional Neural Networks (CNN) [9] have already been used to build UC MES scoring systems [11,2]. But rather than formulating the problem as a 3-way classification task that separates the three MES categories (mild, moderate, severe), existing methods resort to learning binary classifiers to deal with the high degree of intra-class similarity. Consequently, a number of different models needs to be trained.

We propose to amplify the classification accuracy of a CNN network for a 3-way classification using an invariant pretext representation learning technique in a self-supervised setting that exploits patch-based image transformations and additionally use these patches for instance-based group discrimination by grouping same class together using  $k$ -means clustering, referred to as “PLD-PIRL”. The idea is to increase the intera-class separation and minimise the intra-class separation. The proposed technique uses a CNN model together with unsupervised  $k$ -means clustering for achieving this objective. The subtlety in the mucosal appearances are learnt by transforming images into a jigsaw puzzle and computing contrastive losses between feature embedding (*aka* representations). In addition, we also explore the introduction of an attention mechanism in our classification network to further boost classification accuracy. We would like to emphasise that the UC scoring is a complex classification task as image samples are very similar and often confusing to experts (especially between grades 1 and 2). Thus, developing an automated system for this task has tremendous benefit in clinical support system.

The related work on UC scoring based on deep learning is presented in Section 2. In Section 3, we provide details of the proposed method. Section 4 consists of implementation details, dataset preparation and results. Finally, a conclusion is presented in Section 5.

## 2 Related work

### 2.1 CNN-based UC grading

Most research work on UC grading is based on MES scores. Mokter et al. [11] propose a method to classify UC severity in colonoscopy videos by detecting vascular (vein) patterns using three CNN networks and use a training dataset comprising of over 67k frames. The first CNN is used to discriminate between a high and low density of blood vessels. Subsequently they use two CNNs separately for the subsequent UC classification each in binary two class configuration. Such a stacked framework can minimise false positives but does not enhance the

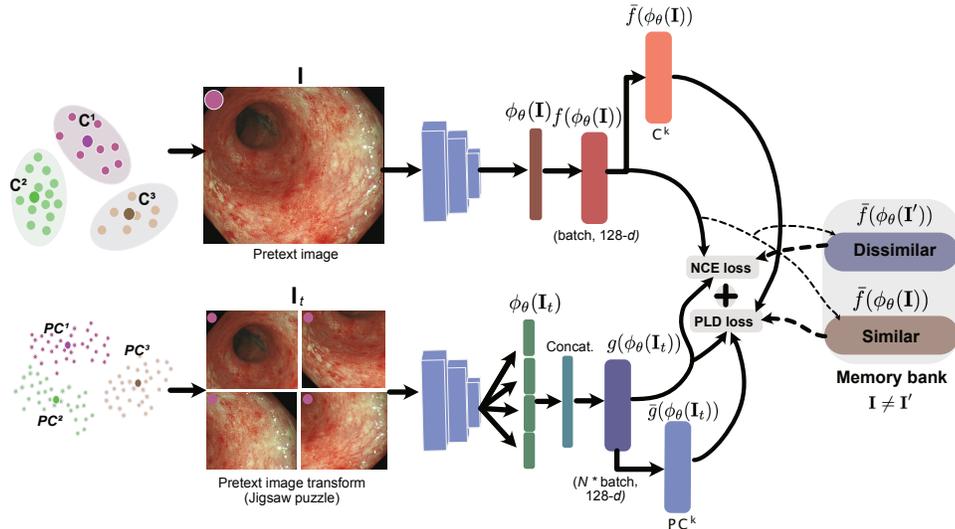
model’s ability to understand variability of different MES scores. Similarly, Stidham *et al.* [14] use the Inception V3 model to train and evaluate MES scores in still endoscopic frames. They used 16k UC images and obtained an accuracy of 67.6%, 64.3% and 67.9% for the three MES classes. UCS-CNN [1] includes several preprocessing steps such as interpretation of vascular patterns, patch extraction techniques and CNN-based classifier for classification. A total of 92000 frames were used for training obtaining accuracy of 53.9% 62.4% and 78.9% for mild, moderate and severe classes. Similarly, Ozawa *et al.* [13] use a CNN for binary classification only on still frames comprising of 26k training images, which first between normal (MES 0) and mucosal healing state (MES 1) while next between moderate (MES 2) and severe (MES 3). Gutierrez *et al.* [2] used CNN model to predict a binary version of the MES scoring of UC.

One common limitation of existing CNN-based UC scoring literature is the use of existing multiple CNN models in an ensemble configuration, simplifying MES to a binary problem and use of very large in-house datasets for training. In contrast, we aim to develop a single CNN model-based approach for a 3-way MES scoring that is clinically relevant. In addition, we use a publicly available UC dataset [3] to guarantee reproducibility of our approach. Furthermore, this dataset consists of only 851 image samples and therefore poses a small data problem.

## 2.2 Self-supervised approach for classification

Self-supervised learning (SSL) uses pretext tasks to mine self-supervised information from large-scale unsupervised data, thereby learning valuable image representations for downstream tasks. By doing so, the limitation of network performance on predefined annotations are greatly reduced. In SSL, the pretext task typically applies a transformation to the input image and predicts properties of the transformation from the transformed image. Chen *et al.* [4] proposed the SimCLR model, which performs data enhancement on the input image to simulate the input from different perspectives of the image. Contrastive loss is then used to maximize the similarity of the same object under different data augmentations and minimised the similarity between similar objects. Later, the MoCo model proposed by He *et al.* [7] also used contrastive loss to compare the similarity between a query and the keys of a queue to learn feature representation. The authors used a dynamic memory, rather than static memory bank, to store feature vectors used in training. In contrast to these methods that encourages the construction of covariant image representations to the transformations, pretext-invariant representation learning (PIRL) [10] pushes the representations to be invariant under image transformations. PIRL computes high similarity to the image representations that are similar to the representation of the transformed versions and low similarity to representations for the different images. Jigsaw puzzle [12] was used as pretext task for PIRL representation learning.

Inspired by PIRL, we propose a novel approach that exploits the invariant representation learning together with patch-level instance-group discrimination.



**Fig. 1:** Pretext invariant patch-level instance group discrimination for ulcerative colitis (UC) scoring. Two identical classification networks are used to compute image-level and patch-level embedding. Three Mayo Endoscopic Scoring (MES) for UC from 1 up to 3 (mild: 1, moderate: 2 and severe: 3) are presented as three separate clusters for both images and patches. The memory bank contains the moving average of representations for all images in the dataset. Here,  $\mathbf{I}$  represent an image sample while  $\mathbf{I}_t$  is a transformed puzzle of that image and  $\mathbf{I}'$  represent negative sample.

Here, the idea is to increase the inter-class separation and minimise the intra-class separation. An unsupervised  $k$ -means clustering is used to define feature clusters for  $k$ -class categories. We demonstrate the effectiveness of this approach on ulcerative colitis (UC) dataset. UC scoring remains a very challenging classification task, while being very important task for clinical decision making and minimising current subjectivity.

### 3 Method

We propose to increase inter-class separation and minimise the intra-class distance by jointly minimising two loss functions that are based on contrastive loss. In contrast to widely used image similarity comparisons we use patch-level and image-level configurations. Additionally, we propose a novel instance group level discriminative loss. A memory bank is used to store moving average embedding of negative samples for efficient memory management. The block diagram of our proposed MES-scoring for ulcerative colitis classification framework is shown in Figure 1.

### 3.1 Pretext invariant representations

Let the ulcerative colitis dataset consists of  $N$  image samples denoted as  $\mathcal{D}_{uc} = \{\mathbf{I}_1, \mathbf{I}_2, \dots, \mathbf{I}_N\}$  for which a transformation  $\mathcal{T}$  is applied to create and reshuffle  $m$  number of image patches for each image in  $\mathcal{D}_{uc}$ ,  $\mathcal{P}_{uc} = \{\mathbf{I}_{1t}^1, \dots, \mathbf{I}_{1t}^m, \dots, \mathbf{I}_{Nt}^1, \dots, \mathbf{I}_{Nt}^m\}$  with  $\mathcal{T} \in t$ . We train a convolutional neural network with free parameters  $\theta$  that embody representation  $\phi_\theta(\mathbf{I})$  for a given sample  $\mathbf{I}$  and  $\phi_\theta(\mathbf{I}_t)$  for patch  $\mathbf{I}_t$ . For image patches, representations of each patch constituting the image  $\mathbf{I}$  is concatenated. A unique projection heads,  $f(\cdot)$  and  $g(\cdot)$ , are applied to re-scale the representations to a 128-dimensional feature vector in each case (see Figure 1). A memory bank is used to store positive and negative sample embedding of a mini-batch  $B$  (in our case,  $B = 32$ ). Negative refer to embedding for  $I' \neq I$  that is required to compute our contrastive loss function  $\mathcal{L}(\cdot, \cdot)$ , measuring the similarity between two representations. The list of negative samples, say  $\mathcal{D}_n$ , grows with the training epochs and are stored in a memory bank  $\mathcal{M}$ . To compute a noise contrastive estimator (NCE), each positive samples has  $|\mathcal{D}_n|$  negative samples and minimizes the loss:

$$\mathcal{L}_{NCE}(\mathbf{I}, \mathbf{I}_t) = -\log[h(f(\phi_\theta(\mathbf{I})), g(\phi_\theta(\mathbf{I}_t)))] - \sum_{\mathbf{I}' \in \mathcal{D}_n} \log[1 - h(f(\phi_\theta(\mathbf{I}')), g(\phi_\theta(\mathbf{I}_t)))] \quad (1)$$

For our experiments we have used both ResNet50 [8] and a combination of convolutional block attention (CBAM, [17]) with ResNet50 model (ResNet50<sup>+cbam</sup>) for computing the representation  $f(\cdot)$  and  $g(\cdot)$ . In Eq. (1),  $h(\cdot, \cdot)$  is the *cosine* similarity between the representations with a temperature parameter  $\tau$ , and for  $h(f(\cdot), g(\cdot))$ :

$$h(f, g) = \frac{\exp \langle f, g \rangle / \tau}{\exp \langle f, g \rangle / \tau + |\mathcal{D}_n| / N} \quad (2)$$

The presented loss encourages the representation of image  $\mathbf{I}$  to be similar to its corresponding transformed patches  $\mathbf{I}_t$  while increasing the distance between the dissimilar image samples  $\mathbf{I}'$ . This enables network to learn invariant representations.

### 3.2 Patch-level instance group discrimination loss

Let  $\bar{f}(\cdot)$  and  $\bar{g}(\cdot)$  be the mean embedding for classes  $k$  with cluster centers  $C^k$  and  $PC^k$  respectively for the image  $\mathbf{I}$  and patch samples  $\mathbf{I}_t$ .  $k$ -means clustering is used to group the embedding into  $k$  ( $= 3$ ) class instances. The idea is to then compute the similarity of each patch embedding  $g(\cdot)$  with the mean image embedding  $\bar{f}(\cdot)$  for all  $k$  classes and vice-versa using Eq. (2). A cross-entropy (CE) loss is then computed that represent our proposed  $\mathcal{L}_{PLD}(\cdot, \cdot)$  loss function given as:

$$\begin{aligned} \mathcal{L}_{PLD}(\mathbf{I}, \mathbf{I}_t) = & -0.5 \sum_{\forall k} C^k \log(h(\bar{f}(\phi_\theta(\mathbf{I})), g(\phi_\theta(\mathbf{I}_t)))) \\ & -0.5 \sum_{\forall k} PC^k \log(h(\bar{g}(\phi_\theta(\mathbf{I}_t)), f(\phi_\theta(\mathbf{I})))) \end{aligned} \quad (3)$$

The proposed patch-level group discrimination loss  $\mathcal{L}_{PLD}(\mathbf{I}, \mathbf{I}_t)$  takes  $k$  class instances into account. As a result not only the similarity between the group (mean) embedding and a single sample embedding for the same class is maximised but also it guarantees inter-class separation. The final loss function with empirically set  $\lambda = 0.5$  is minimised in our proposed PLD-PIRL network and is given by:

$$\mathcal{L}_{final}(\mathbf{I}, \mathbf{I}_t) = \mathcal{L}_{NCE}(\mathbf{I}, \mathbf{I}_t) + \lambda \mathcal{L}_{PLD}(\mathbf{I}, \mathbf{I}_t) \quad (4)$$

## 4 Experiments and results

### 4.1 Implementation Details.

For training of pretext tasks in self-supervised learning, we use a learning rate (lr) of  $1e^{-3}$  and a SGD optimizer. 3000 epochs with a batch size of 32 is used to train pretext tasks presented in all experiments. For PLD loss, we set  $k = 3$  for number of clusters and test the effect of different temperature  $\tau$  and  $\lambda$  on the model performance.

For the next downstream classification task, we use finetune the model with the learning rate of  $1e^{-4}$ , the SGD optimizer, batch size of 32, and the learning rate decay strategy with a learning rate decay of 0.9 times per 30 epochs. Our experimental results showed that most of the models converged around 150 epochs. For the baseline supervised models training converged to higher epochs of nearly 200 epochs. The stopping criteria was based on minimal loss improvement of 0.000001 over 20 consecutive epochs. The proposed method is implemented on a server deployed with an NVIDIA Quadro RTX 6000 graphics card using the PyTorch framework. All input images were resized to  $224 \times 224$  pixels.

### 4.2 Datasets and evaluation.

We have used both publicly available and in-house dataset. HyperKvasir [3] public dataset was used for model training, validation and as hold-out test samples (referred as Test-I). The available dataset includes MES scores (1,2 and 3) and three additional scoring levels categorising into scores 0-1, 1-2 and 2-3, totaling to 6 UC categories and 851 images. After re-examination by expert colonoscopist, the final data was divided into three different grades: mild, moderate and severe. In the experiment, 80% of the data is used for training, 10% is used for validation and 10% is used for testing. Furthermore, to evaluate the efficacy of the proposed PLD-PIRL method on unseen center data (referred as Test-II), we used one in-house dataset as the test set. This dataset contains 151 images from 70 patient videos. We manually selected frames containing UC from the videos, which were then labeled as mild, moderate, and severe by an expert colonoscopist. All datasets will be made public upon acceptance of the paper.

We have used standard top- $k$  accuracy (percentage of samples predicted correctly), F1-score ( $= \frac{tp}{tp+fp}$ , tp: true positive, fp: false positive), specificity ( $= \frac{tn}{tp+fn}$ ) and sensitivity ( $= \frac{tn}{tn+fp}$ ) and for our 3-way classification task of MES-scoring for UC.

### 4.3 Comparison with SOTA methods

Result of baseline fully supervised classification and self-supervised learning model (SSL) for UC classification on two test datasets (Test-I and Test-II) are presented in Table 1. ResNet50 and ResNet50<sup>+cbam</sup> are established as the baseline model for supervised learning and the same are also used for other state-of-the-art (SOTA) SSL comparisons. In Table 1 for Test-I dataset, it can be observed that the proposed PLD-PIRL approach using ResNet50<sup>+cbam</sup> model achieves the best results with 69.04%, 68.98%, 84.71% and 67.35%, respectively, for top 1 accuracy, F1 score, specificity and sensitivity. Compared to the supervised learning based baseline models i.e., ResNet50 and ResNet50<sup>+cbam</sup>, the top 1 accuracy is improved by 4.75% and 3.57%, respectively, for these models using our proposed PLD-PIRL. We also compared the proposed PLD-PIRL approach with other SOTA self-supervised learning methods including popular

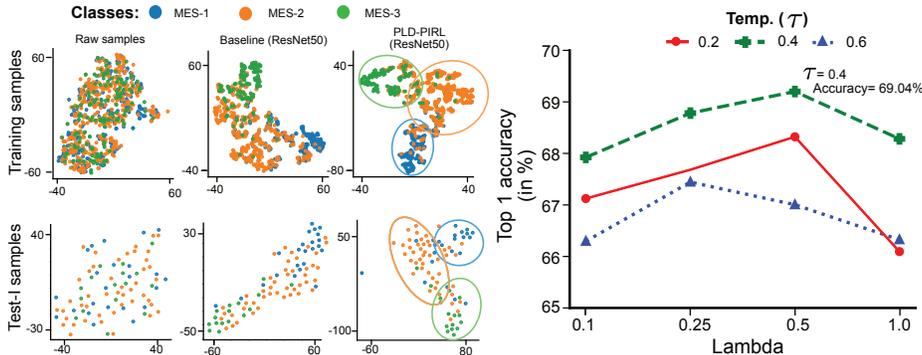
**Table 1:** Experimental results on hold-out test set (Test-I) and unseen center test set (Test-II) for a three way classification of MES scores (1, 2 and 3). Both top 1 and top 2 accuracies are provided together with F1-score, specificity and sensitivity values (in percentage). Two best results are shown in bold. Both classical supervised approach (baseline) and different self-supervised methods are compared with our proposed PLD-PIRL classification method.

Method	Model	Top 1	Top 2	F1	Sen.	Spec.
<b>Test-I</b>						
<b>Baseline</b>	ResNet50	64.29%	91.63%	62.77%	60.56%	81.02%
	ResNet50 <sup>+cbam</sup>	65.47%	93.59%	64.38%	63.58%	81.95%
SimCLR [4]	ResNet50	61.91%	94.93%	59.87%	58.49%	79.60%
SimCLR [4]	ResNet50 <sup>+cbam</sup>	63.09%	92.27%	60.90%	59.31%	79.84%
SimCLR + DCL [5]	ResNet50	64.28%	<b>94.98%</b>	62.34%	62.78%	80.35%
SimCLR + DCL [5]	ResNet50 <sup>+cbam</sup>	64.79%	94.06%	62.01%	62.39%	79.56%
MOCO + CLD [7]	ResNet50	66.96%	93.12%	65.79%	66.38%	84.26%
MOCO + CLD [7]	ResNet50 <sup>+cbam</sup>	67.32%	92.52%	66.38%	65.91%	<b>84.53%</b>
PIRL [10]	ResNet50	65.93%	92.89%	64.87%	64.26%	81.03%
PIRL [10]	ResNet50 <sup>+cbam</sup>	66.67%	93.21%	65.91%	66.47%	82.59%
PLD-PIRL (ours)	ResNet50	<b>67.85%</b>	93.98%	<b>67.48%</b>	<b>66.93%</b>	83.41%
PLD-PIRL(ours)	ResNet50 <sup>+cbam</sup>	<b>69.04%</b>	<b>96.31%</b>	<b>68.98%</b>	<b>67.35%</b>	<b>84.71%</b>
<b>Test-II</b>						
<b>Baseline</b>	ResNet50	57.61%	88.36%	57.03%	56.69%	71.10%
	ResNet50 <sup>+cbam</sup>	60.92%	90.49%	59.38%	58.81%	73.79%
SimCLR [4]	ResNet50	56.95%	85.88%	55.91%	54.69%	71.26%
SimCLR [4]	ResNet50 <sup>+cbam</sup>	57.31%	85.21%	56.29%	56.50%	71.98%
SimCLR + DCL [5]	ResNet50	58.94%	87.92%	57.36%	57.29%	73.29%
SimCLR + DCL [5]	ResNet50 <sup>+cbam</sup>	59.60%	90.34%	58.42%	59.58%	74.19%
MOCO + CLD [7]	ResNet50	60.61%	90.71%	60.52%	59.88%	75.69%
MOCO + CLD [7]	ResNet50 <sup>+cbam</sup>	60.93%	92.12%	60.61%	59.29%	77.33%
PIRL [10]	ResNet50	61.59%	92.61%	60.55%	60.53%	75.98%
PIRL [10]	ResNet50 <sup>+cbam</sup>	62.25%	<b>93.92%</b>	61.96%	60.92%	78.03%
PLD-PIRL (ours)	ResNet50	<b>62.90%</b>	92.93%	<b>62.81%</b>	<b>61.79%</b>	<b>80.23%</b>
PLD-PIRL(ours)	ResNet50 <sup>+cbam</sup>	<b>64.24%</b>	<b>95.32%</b>	<b>64.38%</b>	<b>62.99%</b>	<b>80.09%</b>

Sen. - sensitivity; Spec. - specificity

SimCLR [4], SimCLR+DCL [18], MOCO+CLD [7] and PIRL [10] methods. Our proposed network (ResNet50) clearly outperforms all these methods with at least nearly 1.1% (MOCO+CLD) up to 6% (SimCLR) on top-1 accuracy.

Similarly, for out-of-sample unseen center Test-II dataset (see Table 1), the proposed model outperforms the baseline fully supervised models by a large margin accounting to nearly 5% for ResNet50 and 4% for ResNet50<sup>+cbam</sup>. A similar trend is observed for all SOTA SSL methods ranging from 3.63% for MOCO+CLD upto 5.3% for SimCLR with ResNet50. A clear boost of 1.99% can be seen for the best PLD-PIRL (ResNet50<sup>+cbam</sup>) model compared to the PIRL (ResNet50<sup>+cbam</sup>). Our experiments on all the existing approaches with ResNet and CBAM (ResNet<sup>+cbam</sup>) backbone showed nearly 1% improvement over ResNet50 on test-I with other SOTA (e.g.,



**Fig. 2:** (Left) Classified clusters for three MES classes obtained from fully supervised baseline and proposed PLD-PIRL on both training (top) and test-I (down) samples. Raw sample distributions are also shown. A  $t$ -distributed stochastic neighbor embedding is used for the point plots of image samples embedding. (Right) Experiments for finding best values for hyper-parameters temperature  $\tau$  and  $\lambda$  weights in the loss function.

top 1 accuracies for SimCLR: 63.09%, SimCLR+DCL: 64.79% and for MOCO+CLD: 67.32%) and around 0.5% on test-II (SimCLR: 57.31%, SimCLR+DCL: 59.60% and for MOCO+CLD: 60.93%). Figure 2 (left) representing  $t$ -SNE plots demonstrate an improved separation of sample points in both training and test sets compared to fully supervised baseline approach. Confusion matrix for both Test-I and Test-II are provided in supplementary Figure 1 that shows that proposed PLD-PIRL were able to classify more samples compared to the baseline method. Similarly, it can be observed from supplementary Figure 2 that the wrongly classified ones are only between adjacent classes which at times categories as both classes by the clinical endoscopists.

#### 4.4 Ablation study

Our experiments indicate that the settings of  $\lambda$  and temperature  $\tau$  parameters in the proposed PLD-PIRL approach will affect the model performance. Therefore, we conducted an ablation study experiment to further study the performance of PLD-PIRL under different parameter settings. We set  $\tau = \{0.2, 0.4, 0.6\}$  and  $\lambda = \{0.1, 0.25, 0.5, 1.0\}$ . As can be observed from the plot in Figure 2 (right) that for  $\tau = 0.4$  PLD-PIRL maintained high accuracy at different  $\lambda$  values, and is better than other parameter settings. The best value is obtained at  $\lambda = 0.5$  with the top1 accuracy of 69.04%.

## 5 Conclusion

Our novel self-supervised learning method using pretext-invariant representation learning with patch-level instance-group discrimination (PLD-PIRL) applied to the UC classification task overcomes the limitations of previous approaches that rely on binary classification tasks. We have validated our method on a public dataset and an unseen dataset. Our experiments show that compared with other SOTA classification

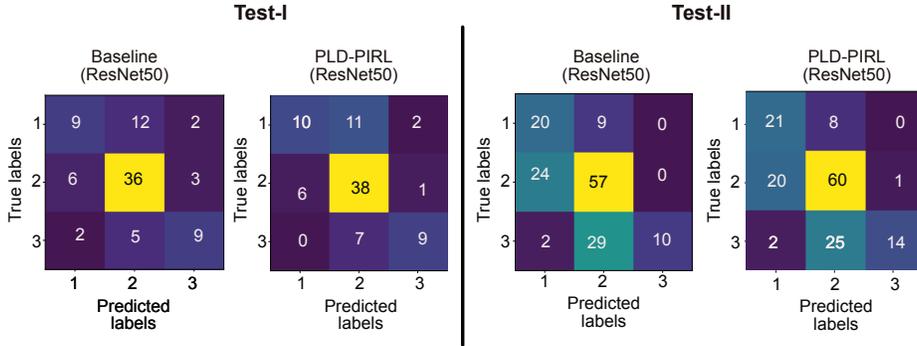
methods that include fully supervised baseline models, our proposed method obtained large improvements in all metrics. The test results on the unseen dataset provides an evidence that our proposed PLD-PIRL method can learn to capture the subtle appearance of mucosal changes in colonic inflammation and the learnt feature representations together with instance-group discrimination allows improved accuracy and robustness for clinically use Mayo Endoscopic Scoring of UC.

## References

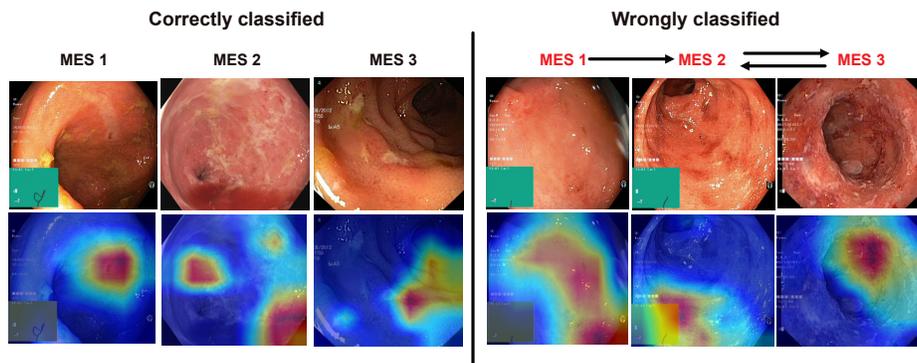
1. Alammari, A., Islam, A.R., Oh, J., Tavanapong, W., Wong, J., De Groen, P.C.: Classification of ulcerative colitis severity in colonoscopy videos using CNN. In: Proceedings of the 9th international conference on information management and engineering. pp. 139–144 (2017)
2. Becker, B.G., Arcadu, F., Thalhammer, A., Serna, C.G., Feehan, O., Drawnel, F., Oh, Y.S., Prunotto, M.: Training and deploying a deep learning model for endoscopic severity grading in ulcerative colitis using multicenter clinical trial data. *Therapeutic advances in gastrointestinal endoscopy* **14** (2021)
3. 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)
4. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: International conference on machine learning. pp. 1597–1607. PMLR (2020)
5. Chuang, C.Y., Robinson, J., Lin, Y.C., Torralba, A., Jegelka, S.: Debaised contrastive learning. *Advances in neural information processing systems* **33**, 8765–8775 (2020)
6. D’haens, G., Sandborn, W.J., Feagan, B.G., Geboes, K., Hanauer, S.B., Irvine, E.J., Lémann, M., Marteau, P., Rutgeerts, P., Schölmerich, J., et al.: A review of activity indices and efficacy end points for clinical trials of medical therapy in adults with ulcerative colitis. *Gastroenterology* **132**(2), 763–786 (2007)
7. He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 9729–9738 (2020)
8. 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. pp. 770–778 (2016)
9. LeCun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proceedings of the IEEE* **86**(11), 2278–2324 (1998)
10. Misra, I., Maaten, L.v.d.: Self-supervised learning of pretext-invariant representations. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6707–6717 (2020)
11. Mokter, M.F., Oh, J., Tavanapong, W., Wong, J., Groen, P.C.d.: Classification of ulcerative colitis severity in colonoscopy videos using vascular pattern detection. In: International Workshop on Machine Learning in Medical Imaging. pp. 552–562. Springer (2020)
12. Noroozi, M., Favaro, P.: Unsupervised learning of visual representations by solving jigsaw puzzles. In: European conference on computer vision. pp. 69–84. Springer (2016)

13. Ozawa, T., Ishihara, S., Fujishiro, M., Saito, H., Kumagai, Y., Shichijo, S., Aoyama, K., Tada, T.: Novel computer-assisted diagnosis system for endoscopic disease activity in patients with ulcerative colitis. *Gastrointestinal endoscopy* **89**(2), 416–421 (2019)
14. Stidham, R.W., Liu, W., Bishu, S., Rice, M.D., Higgins, P.D., Zhu, J., Nallamotheu, B.K., Waljee, A.K.: Performance of a deep learning model vs human reviewers in grading endoscopic disease severity of patients with ulcerative colitis. *JAMA network open* **2**(5), e193963–e193963 (2019)
15. Torres, J., Halfvarson, J., Rodríguez-Lago, I., Hedin, C.R., Jess, T., Dubinsky, M., Croitoru, K., Colombel, J.F.: Results of the seventh scientific workshop of ecco: precision medicine in ibd—prediction and prevention of inflammatory bowel disease. *Journal of Crohn’s and Colitis* **15**(9), 1443–1454 (2021)
16. Vashist, N.M., Samaan, M., Mosli, M.H., Parker, C.E., MacDonald, J.K., Nelson, S.A., Zou, G., Feagan, B.G., Khanna, R., Jairath, V.: Endoscopic scoring indices for evaluation of disease activity in ulcerative colitis. *Cochrane Database of Systematic Reviews* (1) (2018)
17. Woo, S., Park, J., Lee, J.Y., Kweon, I.S.: Cbam: Convolutional block attention module. In: *Proceedings of the European conference on computer vision (ECCV)*. pp. 3–19 (2018)
18. Yeh, C.H., Hong, C.Y., Hsu, Y.C., Liu, T.L., Chen, Y., LeCun, Y.: Decoupled contrastive learning. *arXiv preprint arXiv:2110.06848* (2021)

## Supplementary material



**Supplementary Figure 1:** Confusion matrix for our proposed method (PLD-PIRL) and the baseline fully supervised method both using ResNet50 model on test-I dataset (same distribution, left from the solid line) and test-II dataset (out-of-sample distribution, right from the solid line).



**Suupplementary Figure 2:** Qualitative results showing the model attentions for correctly (left from the solid line) and incorrectly (right from the solid line) classified Mayo Endoscopic Scores. Arrows in wrongly classified ones points to the predicted score for the given sample.