
Unsupervised Domain Adaptation for Label-Free Polyp Segmentation in Capsule Endoscopy

Ang Chen Quanzhi Fu Tianshuo Yan Butian Xiong Nuohan Chen Yulin Li
Chinese University of Hong Kong, Shenzhen

Abstract

1 Capsule endoscopy (CE) is crucial for the diagnosis of gastrointestinal diseases.
2 Due to the limited datasets available for CE and the prohibitive cost of labeling
3 them, there are still relatively few studies that leverage deep learning to analyze
4 CE data. In this project, we propose to use unsupervised domain adaptation
5 (UDA) to achieve label-free polyp segmentation for CE by transferring domain
6 knowledge from an annotated, large-scale traditional endoscopy dataset. We
7 experimented with two types of DA methods: (i) Fourier Domain Adaptation,
8 and (ii) Entropy Minimization, both of which achieved better performance than
9 the unadapted baseline. Our best model, AdvEnt, increased the segmentation
10 Dice score on a customized CE dataset by 6.2%. (Code: [https://github.com/](https://github.com/Depersonalizc/vce-patho-discovery)
11 Depersonalizc/vce-patho-discovery)

12

1 Introduction

13 Capsule endoscopy was designed to solve problems of traditional endoscopy, such as gastroscopy
14 and colonoscopy. These endoscopies has a camera connected with a tube to take photos in human
15 body. Although these information are crucial for diagnosis, these methods are invasive and have
16 limited detective range. Capsule endoscopy solved these problems by load a camera into a capsule.
17 The camera could take photos for the whole GI tract of the patients therefore provide more valuable
18 information for gastroenterologists.

19 The capsule endoscopy also face some challenge. First, examining the recorded video is time
20 consuming. It may take gastroenterologists multiple hours to examine a capsule endoscopy video.
21 Meanwhile, critical information are easy to be missed, as it may only appears in the video for a few
22 seconds. Therefore, researchers are keep develop models to help gastroenterologists increase their
23 efficiency and precision on analyzing capsule endoscopic video. This is also what we try to do.

24 Polyps is also a hot topic in computer aided diagnosis. Polyps are projecting growth of mucous
25 membrane from surfaces in the body. It is an important indicator to cancerwhile always overlooked
26 in endoscopic analysis. Data shows that about 14%-30% polyps fail to be detected during diagnosis.
27 There are matured model for polyps detection for traditional endoscopy. While research focus on
28 capsule endoscopy are still limited.

29 The main reason for the small number of studies in this direction is the lack of data. A capsule
30 endoscopy provides a video about 2 hours. And its movement and camera perspective are pretty
31 random, which result in high cost in data annotation. In a VCE dataset called Kvasir-Capsule which
32 contains more than hundreds of thousand images, has only 55 images with polyps labeled. For CE
33 data, only unlabeled images are available. However, the data from traditional endoscopy is sufficient.
34 There are a amount of dataset for polyp segmentation in traditional colonoscopy.

35 In this project, we use a technique called domain adaptation and developed deep learning models
36 for polyp segmentation in Capsule endoscopy. Assume the traditional endoscopy images come from
37 a source domain and capsule endoscopy images come from a target domain, our models have no

38 requirements on annotated capsule endoscopy data, which provides a data-efficient solution for polyp
39 detection in capsule endoscopy video.

40 The remaining content of the paper is organized as follows. In section 2, we introduce some related
41 works in semantic segmentation and domain adaptation. In section 3, we demonstrate in details the
42 domain adaptation methods we use. In section 4, we report the experimental results and provide an
43 analysis. In section 5, we conclude our project and propose some future directions of the study.

44 2 Related Works

45 2.1 Semantic Segmentation

46 Semantic segmentation has benefited from the continuous evolution of DNN architectures [15][10]. A
47 segmentation architecture ResUNet++ was proposed in 2019 and achieved state of the art performance
48 in polyp segmentation problem [7]. Semantic Segmentation models usually trained on dataset with
49 pixel-level annotations. Different with classification task where human only need to provide a label to
50 each image, manual annotation is not scalable [16]. Thus, techniques like transfer learning or domain
51 adaptation were used to solve the problem of lack of data.

52 2.2 Domain Adaptation

53 Unsupervised Domain Adaptation (UDA) is a well researched topic for multiple tasks like classifi-
54 cation and detection. While it was not intensely explored in semantic segmentation until recently. It
55 aims to reduce the gap between two distributions[11].

56 Adversarial Learning for domain adaptation is one of the most well explored topics. Adversarial
57 learning uses a discriminator trained to maximize the confusion between source and target repre-
58 sentations, thus reduce the domain discrepancy [3]. It involves two networks. The first network
59 predict the segmentation maps from a input from source or target domain. Another network is a
60 discriminator. The discriminator takes the feature maps from the segmentation network and tries
61 to preidict domain of the input. The optimization gaol for segmentation network is to deceive the
62 discriminator. Optimize such a problem can making the features from source and target domain have
63 similar distribution [14].

64 Some methods use a additional generative network to generate target images conditioned on the
65 source. Like Cycle-Consistent Adversarial Domain Adaptation proposed by Hoffman et al [5] and
66 Cycle-GAN [17]. Although these networks shows excellent performance, they usually require more
67 memory and hard to be trained with limited computational resources.

68 Self-training is also an approach. The idea is to use the prediction from an ensembled model or a
69 previous state of model as pseudo-labels for the unlabeled data to train the current model [14]. It can
70 be used with many semi-supervised methods [9] and can dealing with class imbalance and embedded
71 with spatial prior [18].

72 3 Methods

73 In this project, we leverage two promising unsupervised domain adaptation (UDA) methods to tackle
74 the problem of label-free polyp segmentation in capsule endoscopy: (i) a Fourier transformation-based
75 DA technique termed Fourier Domain Adaptation (FDA), (ii) entropy-based DA methods: Direct
76 Entropy Minimization (MinEnt) and Adversarial Entropy Training (AdvEnt).

77 In the setting of unsupervised domain-adapted semantic segmentation, we are given a source (labeled)
78 dataset $D^s = \{(x_i^s, y_i^s) \sim P(x^s, y^s)\}_{i=1}^N$ where $x^s \in \mathbb{R}^{H \times W \times 3}$ is a color image, and $y^s \in$
79 $\mathbb{R}^{H \times W \times K}$ is the one-hot semantic map associated with x^s . Meanwhile, $D^t = \{x_i^t\}_{i=1}^{N_t}$ denotes the
80 target dataset, where the ground-truth semantic labels are absent. Typically, the segmentation network
81 trained on D^s will experience a significant drop in performance when tested on D^t and domain
82 adaptation algorithms are proposed to reduce the performance drop in D^t .

83 **3.1 Fourier Domain Transformation [15]**

84 Fourier domain adaptation [15] is based on a simple assumption. Data from different domains have
 85 only differences in some low-level features: like color, sharpness, contrast, etc. But the high-level
 86 semantic information are same. For example, in Figure 1 the left image is taken from a famous video
 87 game GTA5. And the right image is a real photo taken in a US city. Although they are largely not
 88 same: you can quickly realize that the left image is a screenshot from a video game. But the semantic
 89 information are almost same. They are all city landscape, both contain road, vehicles, building,
 90 etc. Human being can easily ignore the low-level differences and finding similarities in these two
 91 images. However, those low-level differences causes confusion to deep learning model and make its'
 92 performance drop dramatically.



Figure 1: left: City scene from a video game GTA5; right: City scene from CityScape: a real city image dataset. They have similar high-level semantic components while different low-level features like contrast, color, etc. (Source: [15])

93 To shrink the gap generated by different low-level features. Fourier Domain Adaptation utilizing a
 94 principle in image processing, which is low-level features are encoded in low-frequency region in
 95 frequency domain. With this principle, FDA bridges the domain gap by replacing the low-frequency
 96 region of Source by Target. A illustration of the process can be seen in Figure 2. Again, we have a
 97 source image from GTA5 and a target image of real cityscape. To transfer the source to the target
 98 domain. We first do fourier transformation on both image to get their frequency map. Then we replace
 99 source's low frequency region, and apply inverse fourier transformation to obtain the transformed
 100 image. As you can see, we finally got a image with source content but have very similar low-level
 101 information to target image.

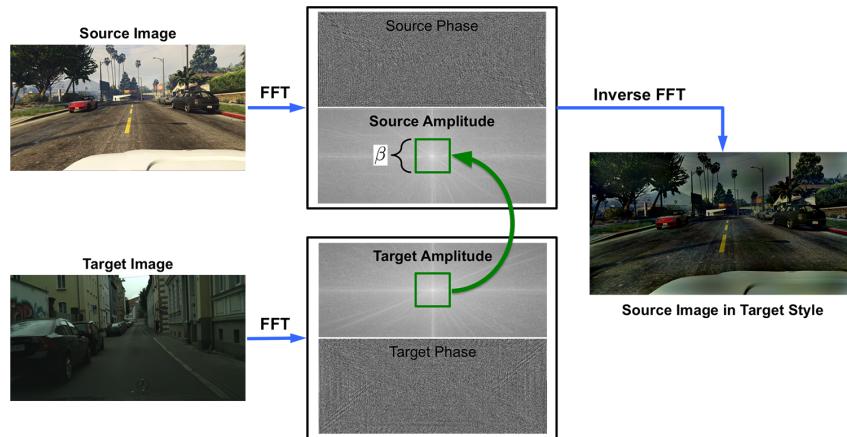


Figure 2: The process of Fourier Domain Adaptation. We apply Fourier transformation to images in both region and replace source image's low-frequency region with the target image's low-frequency region. Then we apply inverse Fourier transformation to obtain the transformed source image. (Source: [15])

102 The rigorous exchange process can be expressed as follows. Let $\mathcal{F}^A, \mathcal{F}^P : \mathbb{R}^{H \times W \times 3} \rightarrow \mathbb{R}^{H \times W \times 3}$
 103 be the amplitude and phase components of the Fourier transform \mathcal{F} of an RGB image. We applied

104 Fourier transform to each channel:

$$\mathcal{F}(x)^{(m,n)} = \sum_{h,w} x^{(h,w)} e^{-j2\pi(\frac{h}{H}m + \frac{w}{W}n)} \quad (1)$$

105 where j denotes the imaginary unit.

106 Given two random sampled images $x^s \sim D^s$, $x^t \sim D^t$, Fourier Domain Adaptation can be formalized
107 as:

$$x^{s \rightarrow t} = \mathcal{F}^{-1}([M_\beta \circ \mathcal{F}^A(x^t) + (1 - M_\beta) \circ \mathcal{F}^A(x^s), \mathcal{F}^P(x^s)]) \quad (2)$$

108 where $x^{s \rightarrow t}$ denotes the transformed image, \circ denotes element-wise multiplication. $M_\beta \in \mathbb{R}^{H \times W \times 3}$
109 denotes the exchange mask, where elements in index range $[-\beta H : \beta H], [-\beta W : \beta W]$ are 1 and 0
110 otherwise. \mathcal{F}^{-1} denotes inverse Fourier transform.

111 In Figure 3 the result of applying Fourier domain adaptation to the Kvasir-SEG source dataset and
Video-Kvasir-Capsule target dataset. The training process of FDA is very similar to the traditional

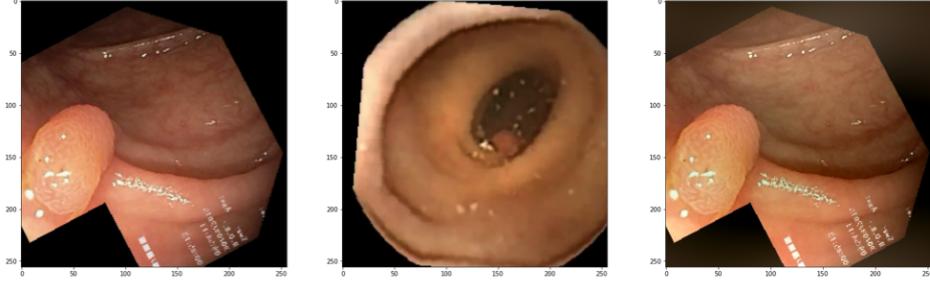


Figure 3: left: source image from Kvasir-SEG dataset; middle: target image from VCE dataset; right:
transformed source image with Kvasir-SEG content and VCE style.

112 semantic segmentation algorithm. Cross entropy is selected as loss function.
113

$$\mathcal{L}(\Phi_\theta; D^{s \rightarrow t}) = - \sum_i \langle y_i^s, \log(\Phi_\theta(x_i^{s \rightarrow t})) \rangle \quad (3)$$

114 where Φ_θ denotes a semantic segmentation network with parameters θ , and $\langle \cdot, \cdot \rangle$ denotes the inner
115 product. The training process can be summarized as follows, in each iteration,

- 116 1. Random sampling batch of source and target images X^s, X^t and the segmentation mask of
117 source image Y^s with same batch size from two datasets respectively.
- 118 2. Apply equation 2 to X^s and X^t to do the domain transform and got transformed images
119 $X^{s \rightarrow t}$.
- 120 3. Use $X^{s \rightarrow t}$ and Y^s to compute cross entropy loss using formula 3
- 121 4. Compute gradients for the network parameters w based on the loss and update the network
122 parameters w .
- 123 5. Repeat step 1-4 until convergence

124 3.2 Entropy-Based Domain Adaptation

125 Entropy-based domain adaptation aims to bridge the domain gap by enforcing the prediction entropy
126 distribution of the model to be similar between source and target domains. Given a semantic
127 segmentation model Φ_θ with parameters θ , and a sample $x \in \mathbb{R}^{H \times W \times 3}$, we denote by $\hat{y} := \Phi_\theta(x) \in$
128 $\mathbb{R}^{H \times W \times K}$ the soft segmentation map where each $(\hat{y}^{(h,w)}) \in \mathbb{R}^K$ is the outputted class probabilities
129 vector of the pixel $x^{(h,w)}$. (For our purpose of polyp segmentation, $K = 2$.) In information theory, a
130 useful quantity is the self-information $I := -\log \hat{y}$, which is closed connected to the notion of the
131 entropy map, \mathcal{H} , by

$$\mathcal{H}^{(h,w)} = \langle \hat{y}^{(h,w)}, I^{(h,w)} \rangle, \quad \forall h \forall w \quad (4)$$

132 where $\langle \cdot, \cdot \rangle$ denotes the inner product. [14] observes that a model trained solely on the source domain
 133 produces over-confident predictions in the source domain and under-confident predictions in the target
 134 domain. This corresponds to low-entropy source predictions and high-entropy target predictions
 (Figure 4).

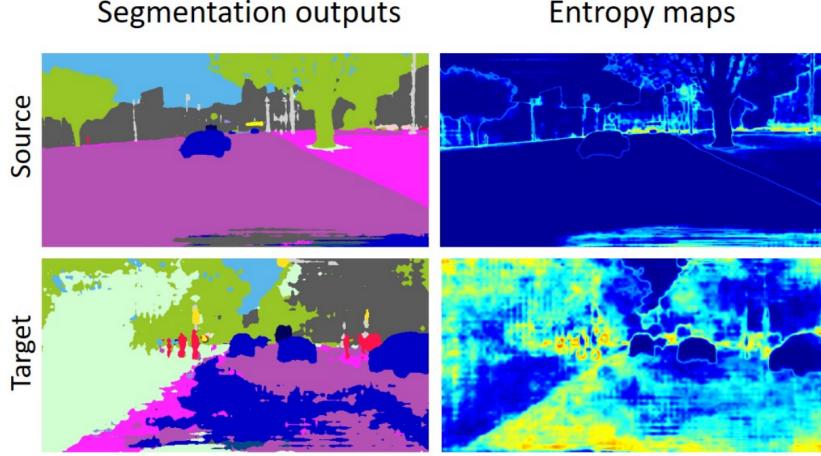


Figure 4: Predicted segmentation maps and entropy maps from the model trained solely in the source domain. (Source: [14])

135
 136 Inspired by this observation, we apply two entropy minimization techniques to unsupervisedly adapt
 137 the model trained on the source domain to the target: (i) Direct Entropy Minimization (MinEnt)
 138 and (ii) Adversarial Entropy Training (AdvEnt). We shall examine each in more details in the following
 139 sections.

140 3.2.1 Direct Entropy Minimization

141 Since the source-trained model produces under-confident predictions on the target domain, MinEnt
 142 puts an entropy regularization term to encourage the model to output confident predictions similar to
 143 that in the source domain. Specifically, we introduce the entropy loss

$$\mathcal{L}_{\text{Ent}}(x^t) = \frac{1}{HW} \sum_{h,w} \mathcal{H}^{(h,w)}(x^t) \quad (5)$$

144 to each unlabeled samples $x^t \in D^t$. On each labeled samples $(x^s, y^s) \in D^s$, we apply the conventional
 145 cross-entropy segmentation loss

$$\mathcal{L}_{\text{Seg}}(x^s, y^s) = \frac{1}{HW} \sum_{h,w} \langle -\log(y^s)^{(h,w)}, (\hat{y}^s)^{(h,w)} \rangle \quad (6)$$

146 where $\log(\cdot)$ is applied element-wise. During MinEnt training, we jointly optimize the segmentation
 147 loss in the source domain and the entropy loss in the target domain, by solving

$$\min_{\theta} \mathbb{E}_{(x^s, y^s) \sim D^s} [\lambda_{\text{Seg}} \mathcal{L}_{\text{Seg}}(x^s, y^s)] + \mathbb{E}_{x^t \sim D^t} [\lambda_{\text{Ent}} \mathcal{L}_{\text{Ent}}(x^t)] \quad (7)$$

148 where $\lambda_{\text{Ent}}, \lambda_{\text{Seg}}$ are the weighting factors.

149 3.2.2 Adversarial Entropy Training

150 Observe from Eq. (4) that the entropy map H can be interpreted as a sum of the self-information
 151 I across classes, weighted by the class probabilities. Directly minimizing the entropy thus ignores
 152 the internal structure of the class information. It is desirable to have a network automatically learn
 153 the most efficient way to combine the predicted self-information. Furthermore, minimizing the
 154 aggregated entropy over all pixel locations may not be the best way to match model's predictions in

155 the source domain. Learning to mimic complex patterns within the entropy map should benefit the
156 performance.

157 To this end, we use adversarial training that simultaneously optimizes the segmentation model Φ_θ and
158 a discriminator D_ω that tries to distinguish source/target predictions from the model Φ_θ , while Φ_θ
159 attempts to fool D_ω . In this work, the discriminator is implemented as a fully-convolutional network
160 that takes as input a class-information map (that could either be from a source or target prediction)

$$\mathcal{I}^{(h,w)} = \hat{y}^{(h,w)} \circ I^{(h,w)}, \quad \forall h \forall w \quad (8)$$

161 and outputs the patch-wise probabilities $D_\omega(\mathcal{I}) \in \mathbb{R}^{H' \times W'}$ of each patch being from a source
162 prediction. Here, \circ is multiplication applied element-wise. Define the likelihood functions:

$$\begin{aligned} \mathcal{L}_s(x^s) &= \frac{1}{H'W'} \sum_{h,w} \log D_\omega(\mathcal{I}(x^s))^{(h,w)} \\ \mathcal{L}_t(x^t) &= \frac{1}{H'W'} \sum_{h,w} \log (1 - D_\omega(\mathcal{I}(x^t))^{(h,w)}) \end{aligned} \quad (9)$$

163 During training, we jointly optimize Φ_θ and D_ω by solving the min-max problem

$$\min_{\theta} \max_{\omega} \mathbb{E}_{x^t \sim D^t} [\lambda_{\text{Adv}} \mathcal{L}_t(x^t)] + \mathbb{E}_{(x^s, y^s) \sim D^s} [\lambda_{\text{Adv}} \mathcal{L}_s(x^s) + \lambda_{\text{Seg}} \mathcal{L}_{\text{Seg}}(x^s, y^s)] \quad (10)$$

164 4 Experiments

165 4.1 Experimental Details

166 4.1.1 Data and Preprocessing

167 In this project, we choose Kvasir-SEG, a public traditional endoscopy polyp segmentation dataset,
168 as our source-domain dataset. We manually collect video capsule endoscopy (VCE) dataset as our
169 target-domain dataset.

170 **Kvasir-SEG** The Kvasir-SEG dataset [6] contains 1000 polyp images and their corresponding
171 ground truth from the Kvasir Dataset v2. The resolution of the images contained in Kvasir-SEG
172 varies from 332x487 to 1920x1072 pixels. The images and its corresponding masks are stored in two
173 separate folders with the same filename. The image files are encoded using JPEG compression, and
174 online browsing is facilitated. The open-access dataset can be easily downloaded for research and
175 educational purposes.

176 **VCE Dataset** Kvasir-Capsule [13] contains images from inside the gastrointestinal (GI) tract via an
177 capsule endoscopy procedure. The images are carefully annotated by medical experts. Here, we only
178 use the polyps class, and three normal classes in distinct location of GI tract. Considering the limited
179 polyp samples (55) in the dataset, we also scrape another 120 CE polyp images from the internet and
180 annotate them manually (only for testing).

181 **Preprocessing** We resize all images and masks to 256×256 and applied random rotation, random
182 flips, and random translation during the training phase.

183 4.1.2 Implementation Details

184 We test all our domain adaptation methods using an state-of-the-art poly segmentation network called
185 ResUNet++[7]. We train the baseline model on Kvasir-SEG and test the performance on VCE Dataset
186 directly to get the unadapted baseline performance. For FDA, MinEnt, and AdvEnt, we use the
187 baseline model as initialization and then apply the DA technique to train for 10000 iterations.

188 We use the PyTorch framework to implement our model with the following parameter settings:
189 **baseline**: mini-batch size 16, learning rate $1e - 2$; **Fourier DA and MinEnt**: mini-batch size 16,
190 learning rate $1e - 2$; **AdvEnt**: mini-batch size 8, learning rate $1e - 3$; **MinEnt+AdvEnt**: mini-batch
191 size 8, learning rate $1e - 3$.

Models	Acc(%)	Prec(%)	Recall(%)	IoU(%)	Dice(%)
Baseline (w/o DA)	91.1	62.4	77.8	53.0	69.3
FDA	92.0	67.7	72.7	54.0	70.1
MinEnt	93.0	74.6	69.6	56.3	72.0
AdvEnt	92.9	72.3	73.1	57.1	72.7
MinEnt+AdvEnt	92.6	67.9	80.3	58.2 (+9.8%)	73.6 (+6.2%)

Table 1: Result Comparison for Different Domain Adaptation Models

192 4.2 Results

193 4.2.1 Quantitative Results

194 The experiments results are summarized in table 1. All unsupervised DA methods outperform
 195 the unadapted baseline, although improvement is more significant with the entropy minimization
 196 approaches. With the combination of direct entropy minimization and adversarial training, we are
 197 able to achieve a 9.8% increase in the IoU as well as a 6.2% increase in Dice Score. FDA shows
 198 a poorer performance compared with Entropy based methods. A possible reason is the internal
 199 complexity of the CE data. Different from cityscape images, factors including camera position or
 200 angle may result in a significant change in image style. In our experiments, we also noticed that even
 201 a very small β may lead to obvious artifacts in the transformed image. Therefore, FDA may not be
 202 the most promising approach when dealing with CE data.

203 4.2.2 Qualitative Results

204 Figure 5 shows the segmentation maps outputed by different methods. It is evident that models
 205 after domain adaption produces segmentation mask with sharper, more defined polyp boundaries.
 206 Compared with entropy-based results, FDA sometimes fail to output a segmentation. This may due
 207 to during the training process, some source images are transformed using target image with very
 208 different content. Entropy-based methods have a better performance. The improvement is more
 209 evident when we look at the entropy map (Figure 6) of the model prediction. The models with entropy
 210 based domain adaptation learns to concentrate the entropy only along the boundary of the polyps,
 211 which leads to sharper predictions.

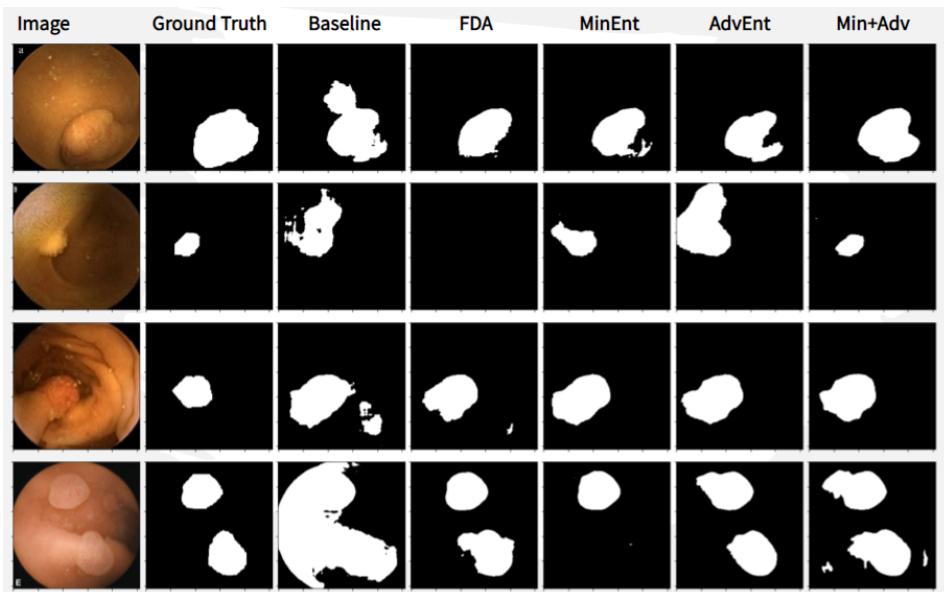


Figure 5: Predicted segmentation maps and entropy maps of the model in the source and target domains.

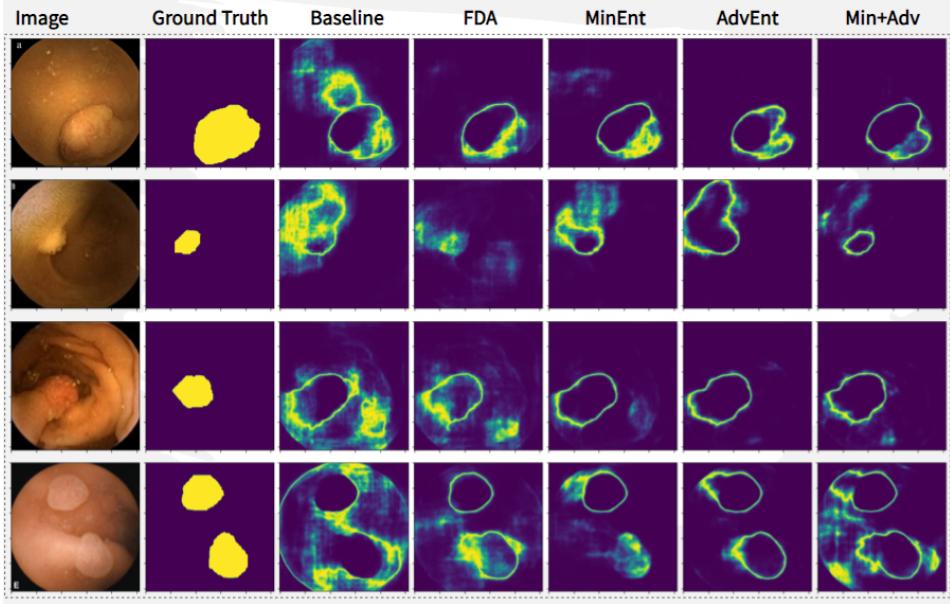


Figure 6: Predicted segmentation maps and entropy maps of the model in the source and target domains.

212 5 Conclusion & Future Work

213 In this project, we utilize three domain adaptation methods to tackle the polyp segmentation for
 214 capsule endoscopy video. All three domain adaptation shows performance increases compared with
 215 the baseline. Among three DA methods, experimental show that the FDA algorithm is not very
 216 suitable for the conversion from traditional endoscopy to capsule endoscopy. The two DA methods
 217 based on entropy minimization shown a superior performance.

218 In the future, we plan to collect a larger capsule endoscopy dataset to improve DA performance and
 219 validation accuracy. In this work, we mainly focus on inter-domain adaptation, more intra-domain
 220 adaptation are also worth to try. Incorporating more label-free learning methods, like self-supervised
 221 learning or active learning, is also a direction to further boost the model performance.

222 A Appendix

223 We did a lot of experimentation before actually choosing Domain Adaptation as our topic. These
 224 attempts include motion deblur in capsule endoscopy, direct localization of polyps, keyframe ex-
 225 traction, and more. Due to data or computing power constraints, these explorations did not end up
 226 being our subject. We attach some experimental reports from the exploration process here for future
 227 reference.

228 A.1 Feature Extraction for Important Frame Detection

229 The detection of the important frames of a video produced by Capsule Endoscopy (CE) is important
 230 for the polyps detection, because an extremely large number of frames are generated from the
 231 screening procedure of tract (just video recordings, which are able to take up to 30 minutes to 8 hours
 232 to complete the process and approximately more than 5 thousand endoscopy frames are produced) in
 233 which only a small partition of the frames contain important information – polyps [12]. However,
 234 it is a lengthy and hard process for the experts to observe each frame separately, and the abnormal
 235 frames can be easily overlooked in the procedure. Therefore, there is a requirement for automatic
 236 polyps clinical practitioners of CE-produced frames, which is usually realized by the Computer-aided
 237 diagnosing (CAD) systems based on varient machine-vision algorithms.

238 **A.1.1 Feature Extraction Strategies**

239 Although according to researchers [1], there are few researches target on automatic important frames
240 detection of polyps by feature extraction algorithms, the feature extraction strategy could gain
241 advantage on acceptable interpretability compared to other deep learning based algorithms like
242 CNN, LSTM based models. The extracted features themselves have mathematical or even biological
243 meanings indicated by the algorithms they base on, which might provide insights for the visual
244 properties of the abnormal frames. Moreover, these features later can be used for the segmentation
245 (as what our groups did), classification, and retrieval of the images. In the view of computation itself,
246 the feature extraction reduces the dimension of every matrix-form image frame into linear vector of
247 features, which significantly saves the computational source and accelerates the classification process.

248 There are various features of colon images, which could be categorized into spatial domain and
249 frequency domain. The frequency domain based algorithms have been implemented enough in the
250 other part of our research, thus here we will mainly focus on the spatial domain features. A digital
251 image is represented by a two dimension matrix in the spatial domain of image processing. The
252 spatial refers to directly manipulation and analysis of the values in the matrix (pixels). So the spacial
253 features relate to the location of the pixels, every pixel is precisely investigated [4]. Those spacial
254 features could provide the internal location structure of the pixels. Since we assume we do not have
255 prior knowledge of each image frame, we will use several commonly used spacial features of images
256 in the proposed research:

- 257 • Local Binary Patterns (LBP): representing images' texture, with the basic form calculated by
258 comparing neighboring entries with the central entry and assigns binary code to it. Then
259 these binary codes are transformed into decimal numbers. Forming a histogram, the texture
260 of image is represented locally by occurrence of these codes. An significant advantage is the
261 rotation invariance.
- 262 • Shape-Based Features: fractal dimension, smooth spiral curve, Koch snowflake, Sierpinski
263 triangle can be used for shape feature of an image. Higher order local auto-correlation
264 (HLAC) features are also used for some geometrical features. Several geometric features
265 were used to classify the colon polyps [8].
- 266 • Color Histograms: hinting the likelihood of a pixel intensity and the distribution of colors.
- 267 • Statistical Texture Features: statistical measures are widely used to represent image texture.
268 Energy, contrast, correlation and homogeneity could be calculated from the image frames as
269 the extracted features. The meaning of the features can be seen from the statistical meaning.

270 In the following part, we will explain how to utilize these extracted features for polyps detection.

271 **A.1.2 Methods and Experiments**

272 Dataset: Kvasir.

273 Kvasir is a dataset containing images from inside the gastrointestinal (GI) tract via an endoscopy
274 procedure. The images are carefully annotated by medical experts from Vestre Viken Health Trust in
275 Norway and the Cancer Registry of Norway (CRN). It contains several classes showing anatomical
276 landmarks, pathological findings or endoscopic procedures in the GI tract. The dataset consists of
277 images of different resolutions from 720x576 to 1920x1072 pixels, organized in a way that they are
278 sorted in separate folders named according to the content. Here, we only use the polyps class, and
279 three normal classes in distinct location of GI tract. The distinct locations are suitable as the negative
280 label classes, because this matches the real life application scenario. We use the images instead of
281 a video because there is no time domain information used in our method, and the data is annotated
282 more precisely.

283 Data Processing:

284 The images go through pre-processing by normalization and contrast enhancement in order to enhance
285 the performance. Then we extract the features from each image in the dataset.

286 Since the prior knowledge is not known in application scenario, all the features mentioned in above
287 part are selected for targeted extracted features, such as Local Binary Patterns, fractal dimension
288 (indicating geometrical property of the boundary) and so on. All these feature extraction algorithms
289 are relatively mature in implementation.

Table 2: Results

Model	Sensitivity	Specificity	AUC
Kernel SVM	54.5%	84.6%	0.642
Naïve Bayes	58.7%	94.8%	0.660

290 After the extraction, an image is actually transformed from a 3-D tensor (RGB image) to a linear
 291 vector reflecting internal properties. PCA is performed to further reduce the dimension to 4 (achieve
 292 the weight of PCs larger than 75%) and extract the important features. For the classification part,
 293 the Kernel support vector machine (Kernel-SVM) and naïve Bayes are selected for the binary
 294 classification task in order to gain a acceptable explanability. In detail, NuSVC (Nu-Support Vector
 295 Classification) in Scikit-learn is chosen for the task, with kernel RBF to realize the non-linear
 296 classification. Deep learning algorithms based models are not selected for this task because of
 297 their low interpretability. The sensitivity, specificity and AUC are utilized to evaluate the model
 298 performance because this is a imbalanced task.

299 **A.1.3 Results and Discussion**

300 The evaluation results of the classification are shown in Table 1. As we can see in Table 1, the
 301 sensitivity and AUC of both models are bad, close to fifty percent. These results indicates that the
 302 algorithms can not get a acceptable prediction of the positive samples, or in practice, the significant
 303 frames in the CE video. The relatively high specificity results from the imbalanced data instead of
 304 that the model suits the structure of the negative samples from this point of view. And these models
 305 can not be applied to CE-video important frame detection with an acceptable result.

306 One of the possible reasons is that the data pre-processing procedure is not enough. Frames acquired
 307 from endoscopy normally suffer from kinds of noises and variations such as illumination invariance,
 308 scale invariance, rotation invariance and so on. Other noises like poor cleansing and bubbles
 309 commonly exist, disturbing the analysis. And the normalization and contrast enhancement can not
 310 deal with above issues.

311 Another important reason may be that the extracted features can not generalize to all polyps well.
 312 That is, not all polyps share similar mathematical or statistical structures in their image. Then the
 313 features extracted naturally show distinct distributions. However, this needs further investigation
 314 because there is not enough computational sources in this research. And there exist multiple steps
 315 which could be enhanced. Deep learning based models like CNN-based or transformer-based models
 316 can be utilized to investigate the effectiveness of extracted features and to discover the internal nature
 317 of the mathematical and statistical structure of polyps CE-frames in the future.

318 **A.2 YOLO model**

319 **A.2.1 Data Preparation**

320 Our target is to check the performance of YOLO-v5[2] model using 2-fold cross validation of
 321 the Kvasir-Capsule dataset. The YOLO-v5 model is the latest version of YOLO model, which is
 322 developed by the team of a company called Ultralytics. In this experiment, we used YOLO-v5-v6.1.
 323 The data set is a part of the whole Kvasir-Capsule dataset. This specific set of images is named
 324 "Labelled images" with 47284 pictures in total. We merely focus on the Luminal images, consisting
 325 of 11 classes: Angiectasia, Blood - fresh, Blood - hematin, Erosion, Erythema, Foreign Body,
 326 Lymphangiectasia, Normal clean mucosa, Polyp, Reduced Mucosal View, Ulcer. We first exclude
 327 the "Normal clean mucosa" and "Reduced Mucosal View" class, because the normal class without
 328 labelling the lesions are useless in YOLO model. Also, because the insufficiency of the image
 329 number, "Polyp" and "Blood - hematin" class are removed. The image distribution of the remaining
 330 classes is shown in Figure 7.

331

332 Another important point of data preparation lays on the dataset partition. For many images are from
 333 the same video, some of which are pictures taken in a continuous time slot. This series of images
 334 assembles each other and will leads to a over-fitting model in terms of this specific class. This will
 335 influence the evaluation of the performance of the YOLO model as well, on account of the existence

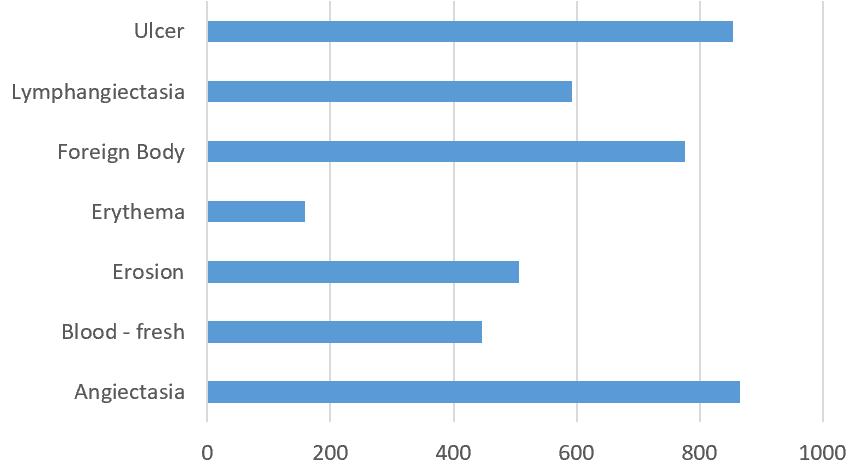


Figure 7: YOLO model classes.

Table 3: The partition of dataset based on video

Class	Video numbers	Video of the first validation set
Erosion	9	'0728084c8da942d9', 'dac1e27f7e4d4ef5', 'eb0203196e284797', '5e59c7fdb16c4228', '8ebf0e483cac48d6'
Angiectasia	6	'eb0203196e284797', 'd369e4f163df4aba', '64440803f87b4843'
Erythema	3	'5e59c7fdb16c4228'
Blood - fresh	2	'd369e4f163df4aba'
Foreign Body	4	'8885668afb844852', '7a47e8eacea04e64'
Ulcer	3	'7a47e8eacea04e64'
Lymphangiectasia	3	'7ad22d50ebaf4596'

336 of similar images appears in both training dataset and validation dataset. The following is the table
 337 showing how we divide the dataset into 2 parts based on the video code 3.

338 A.2.2 Model training and validation

339 The Code for YOLO-v5 is available from is available from:

340 <https://github.com/ultralytics/yolov5>

341 We are supposed to have a deblurred image set based the DeblurGAN model, and use the Yolo-v5
 342 model as a criteria to see how this deblurring model performs. However, because the labelled images
 343 do not contain blurred images, it is impossible to realize this idea in this dataset.
 344

345 We trained the model based on the partition introduced. The result are shown in figure 8. Using the
 346 weights file, we developed a lesion detecting app. Figure 9 is a sketch of the app.
 347

348 References

- 349 [1] Hussam Ali, Muhammad Sharif, Mussarat Yasmin, Mubashir Husain Rehmani, and Farhan
 350 Riaz. A survey of feature extraction and fusion of deep learning for detection of abnormalities
 351 in video endoscopy of gastrointestinal-tract. *Artificial Intelligence Review*, 53(4):2635–2707,
 352 2020.
- 353 [2] Yiming Fang, Xianxin Guo, Kun Chen, Zhu Zhou, and Qing Ye. Accurate and automated
 354 detection of surface knots on sawn timbers using yolo-v5 model. *BioResources*, 16(3), 2021.

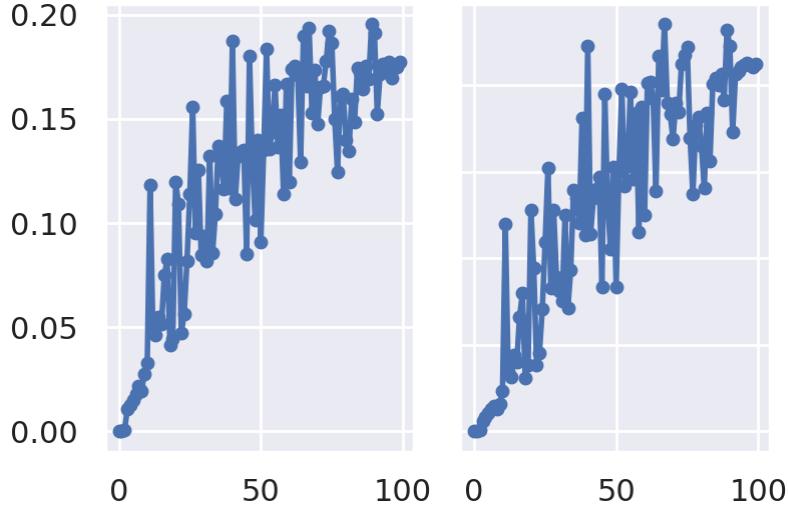


Figure 8: mAP_0.5 for the 2-fold cross validation



Figure 9: YOLO detection on a unlabeled image from Kvasir-Capsule dataset

- 355 [3] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation.
356 In *International conference on machine learning*, pages 1180–1189. PMLR, 2015.
- 357 [4] Kazuhiro Gono, Takashi Obi, Masahiro Yamaguchi, Nagaaki Oyama, Hirohisa Machida, Yasushi
358 Sano, Shigeaki Yoshida, Yasuo Hamamoto, and Takao Endo. Appearance of enhanced tissue
359 features in narrow-band endoscopic imaging. *Journal of biomedical optics*, 9(3):568–577, 2004.
- 360 [5] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros,
361 and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *International
362 conference on machine learning*, pages 1989–1998. PMLR, 2018.
- 363 [6] Debesh Jha, Pia H Smedsrud, Michael A Riegler, Pål Halvorsen, Thomas de Lange, Dag
364 Johansen, and Håvard D Johansen. Kvasir-seg: A segmented polyp dataset. In *International
365 Conference on Multimedia Modeling*, pages 451–462. Springer, 2020.
- 366 [7] Debesh Jha, Pia H Smedsrud, Michael A Riegler, Dag Johansen, Thomas De Lange, Pål
367 Halvorsen, and Håvard D Johansen. Resunet++: An advanced architecture for medical image
368 segmentation. In *2019 IEEE International Symposium on Multimedia (ISM)*, pages 225–2255.
369 IEEE, 2019.
- 370 [8] SM Krishnan, Xin Yang, Kap Luk Chan, and PMY Goh. Region labeling of colonoscopic
371 images using fuzzy logic. In *Proceedings of the First Joint BMES/EMBS Conference*. 1999

- 372 *IEEE Engineering in Medicine and Biology 21st Annual Conference and the 1999 Annual Fall*
373 *Meeting of the Biomedical Engineering Society (Cat. N, volume 2, pages 1149–vol. IEEE, 1999.*
- 374 [9] Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. *arXiv preprint*
375 *arXiv:1610.02242*, 2016.
- 376 [10] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for se-
377 mantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern*
378 *recognition*, pages 3431–3440, 2015.
- 379 [11] Saeid Motian, Marco Piccirilli, Donald A Adjerooh, and Gianfranco Doretto. Unified deep
380 supervised domain adaptation and generalization. In *Proceedings of the IEEE international*
381 *conference on computer vision*, pages 5715–5725, 2017.
- 382 [12] Sonu Sainju, Francis M Bui, and Khan A Wahid. Automated bleeding detection in capsule
383 endoscopy videos using statistical features and region growing. *Journal of medical systems*,
384 38(4):1–11, 2014.
- 385 [13] Pia H Smedsrød, Vajira Thambawita, Steven A Hicks, Henrik Gjestang, Oda Olsen Nedrejord,
386 Espen Naess, Hanna Borgli, Debesh Jha, Tor Jan Derek Berstad, Sigrun L Eskeland, et al.
387 Kvasir-capsule, a video capsule endoscopy dataset. *Scientific Data*, 8(1):1–10, 2021.
- 388 [14] Tuan-Hung Vu, Himalaya Jain, Maxime Bucher, Matthieu Cord, and Patrick Pérez. Advent: Ad-
389 versarial entropy minimization for domain adaptation in semantic segmentation. In *Proceedings*
390 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2517–2526,
391 2019.
- 392 [15] Yanchao Yang and Stefano Soatto. Fda: Fourier domain adaptation for semantic segmentation.
393 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
394 pages 4085–4095, 2020.
- 395 [16] Yang Zhang, Philip David, and Boqing Gong. Curriculum domain adaptation for semantic
396 segmentation of urban scenes. In *Proceedings of the IEEE international conference on computer*
397 *vision*, pages 2020–2030, 2017.
- 398 [17] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image
399 translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international*
400 *conference on computer vision*, pages 2223–2232, 2017.
- 401 [18] Yang Zou, Zhiding Yu, BVK Kumar, and Jinsong Wang. Unsupervised domain adaptation
402 for semantic segmentation via class-balanced self-training. In *Proceedings of the European*
403 *conference on computer vision (ECCV)*, pages 289–305, 2018.