# Mask-conditioned latent diffusion for generating gastrointestinal polyp images

Roman Macháček[*]
University of Oslo
Oslo, Norway
ro.machacek0@gmail.com

Leila Mozaffari[*]
Oslo Metropolitan University
Oslo, Norway
s372060@oslomet.no

Zahra Sepasdar
Monash University
Melbourne, Australia
zahra.sepasdar@monash.edu

Sravanthi Parasa
Swedish Medical Group
Seattle, USA
vaidhya209@gmail.com

Pål Halvorsen[†]
paalh@simula.no
SimulaMet
Oslo, Norway

Michael A. Riegler[†]
michael@simula.no
SimulaMet
Oslo, Norway

Vajira Thambawita[*]
vajira@simula.no
SimulaMet
Oslo, Norway

## ABSTRACT

In order to take advantage of artificial intelligence (AI) solutions in endoscopy diagnostics, we must overcome the issue of limited annotations. These limitations are caused by the high privacy concerns in the medical field and the requirement of getting aid from experts for the time-consuming and costly medical data annotation process. In computer vision, image synthesis has made a significant contribution in recent years, as a result of the progress of generative adversarial networks (GANs) and diffusion probabilistic models (DPMs). Novel DPMs have outperformed GANs in text, image, and video generation tasks. Therefore, this study proposes a conditional DPM framework to generate synthetic gastrointestinal (GI) polyp images conditioned on given generated segmentation masks. Our experimental results show that our system can generate an unlimited number of high-fidelity synthetic polyp images with the corresponding ground truth masks of polyps. To test the usefulness of the generated data we trained binary image segmentation models to study the effect of using synthetic data. Results show that the best micro-imagewise intersection over union (IOU) of 0.7751 was achieved from DeepLabv3+ when the training data consists of both real data and synthetic data. However, the results reflect that achieving good segmentation performance with synthetic data heavily depends on model architectures.

[*]Contributed equally to this research.

[†]Also affiliated with Oslo Metropolitan University, Norway.

## CCS CONCEPTS

• **Computing methodologies → Machine learning algorithms**; **Computer vision tasks**; **Supervised learning**; **Image segmentation**; **Neural networks**.

## KEYWORDS

diffusion model, polyp generative model, polyp segmentation, generating synthetic data

## 1 INTRODUCTION

The human digestive system can experience a range of abnormal tissue changes, from minor discomforts to severe, life-threatening illnesses [18]. Endoscopy, colonoscopy, and pilcams (wireless capsule endoscopy) [15] are the most common methods for examining the gastrointestinal (GI) tract for diagnosis. However, its effectiveness is greatly impacted by the variability in the performance of the operator (inter-rater reliability) [4]. In this regard, artificial intelligence (AI) techniques are researched to build computer-aided diagnosis (CAD) systems to aid gastroenterologists [21, 29, 40, 43].

Supervised machine learning models have become popular in many applications, such as image classification, image detection, and image segmentation. However, AI models require vast amounts of data to train. Especially supervised machine learning techniques need annotated datasets to train. In the medical field, however, acquiring a large annotated dataset is challenging. The challenges include not only privacy concerns but also costly and timely medical data labeling and annotation. In comparison with other applications of machine learning in the health area, we have limited annotated datasets to train machine learning (ML) models. GI-tract datasets are

also typically small and mostly limited to polyps [24]. To overcome this issue, one solution is to expand training datasets by generating synthetic data [17, 37].

Endoscopy diagnostics are being enhanced by AI solutions. Especially, image synthesis has made a significant contribution to overcome the issue of the limited dataset [31]. It is now common to use generative adversarial networks (GANs) to generate synthetic images because GANs produce realistic images and achieve impressive results in a wide range of applications [2, 7]. Thus, a GAN is a powerful generative model, however, it suffers from convergence instability.

To overcome the convergence issue in GANs, in recent years, diffusion models [13] have gained attention as a potential method for their ability to synthesize natural images. In this study, we introduce a framework consisting of two different diffusion models to create synthetic GI-tract images and corresponding masks. Our contributions are listed as follows:

- We introduce a fully synthetic polyp generation system.
- Our system is able to generate realistic-looking synthetic polyp masks using an improved diffusion model.
- Based on the generated masks, we are able to generate high-fidelity synthetic polyp images conditioned on pre-generated synthetic polyp masks using a conditional latent diffusion model.
- We provide a comprehensive evaluation of using synthetic polyp and mask data to train polyp segmentation models and overall results.

The source code of all the experiments is available at https://github.com/simulamet-host/conditional-polyp-diffusion and the pre-generated synthetic masks and the corresponding conditional synthetic polyp images are available at https://huggingface.co/datasets/deepsynthbody/conditional-polyp-diffusion.

## 2 RELATED WORK

There are many GI image analysis datasets available for machine learning tasks. Some of the commonly used datasets in human GI tract are: ETIS-Larib [33], CVC-ClinicDB [3], ASU-Mayo Clinic Polyp database [36], Kvasir [27], Kvasir-SEG [16] and Hyperkvasir [4]. A few datasets containing manually annotated segmentation masks for polyps. However, these real-world datasets (not limited to GI-tract data) have some limitations. The limitations include:

- Size: medical image datasets, including those for polyp detection and segmentation, are often smaller in size compared to other image datasets, such as ImageNet [20], Microsoft COCO [23] which can limit their ability to train complex machine learning models.
- Annotation quality: the accuracy and consistency of the annotations of the dataset can impact the performance of machine learning algorithms. Annotations are dependent on annotator and normally high inter-rater variability is there.
- Diversity: the diversity of the images in the dataset is important for the generalization of machine learning algorithms. If the dataset is limited to a narrow range of images, the algorithm may not perform well on new, unseen images.

- Accessibility: legal and privacy constraints can limit the accessibility of medical image datasets, making it difficult to obtain large and diverse datasets for machine learning tasks.

These limitations highlight the need for ongoing development and improvement of medical image datasets to support the advancement of machine learning in medical imaging. To overcome these limitations of real-world datasets, synthetic datasets[6, 38, 39, 41, 42] have been increasingly used in medical image analysis. For instance, to generate synthetic polyps, a GAN framework has been proposed to generate a polyp mask and then synthesize the generated polyp with the real polyp image without the use of additional datasets and processes [28]. There has also been research on the augmenting of colonoscopy images with polyps by using synthetic samples [1]. Fagereng et al. [10] present a solution to overcome the challenge of a lack of annotated data when building CAD systems for detecting polyps in the GI-tract. The authors propose a pipeline called PolypConnect, which can convert non-polyp images into polyp images to increase the size of training datasets for machine learning models. The results show 5.1% improvement in mean intersection over union (IOU) when using synthetic data in addition to real data for training. Dhariwal et al. [8] compare the performance of diffusion models and GANs for image synthesis tasks. As a result, the authors found that diffusion models outperformed GANs in terms of image quality and stability of the generated images. The results of the paper indicate that diffusion models are a promising alternative to GANs in the field of image synthesis. This work provides valuable insights into the strengths and limitations of both diffusion models and GANs.

## 3 METHODOLOGY

In this section, first, we describe the improved diffusion model which is used to generate realistic synthetic polyp mask images (the blue box of Figure 1). Then, we present the latent diffusion model which is used for generating synthetic polyp images conditioned on the input masks generated from our aforementioned mask generator (the green box of Figure 1). We evaluate the quality of the generated synthetic data and quantify similarity between generated and real data, representing using the last section of Figure 1. Finally, we present our methods used to check the quality of synthetic data using image segmentation models.

### 3.1 Improved diffusion model

In our pipeline, we use an *improved diffusion model* [25] to generate synthetic polyp masks which looks realistic to capture the distribution of the masks of the Kvasir-SEG dataset. The *improved diffusion model* is a type of generative model that uses a gradual, multi-step process to match a data distribution and generate synthetic images. In the context of generating synthetic mask images for the GI-tract, improved diffusion models can be used to generate synthetic mask images that closely resemble real images. Therefore, these models overcome the issue of limited annotated data [25, 34].

To achieve our goal, the first step is to obtain a training set for real mask images of the GI tract indicating the location of polyps. Then, the mask dataset is used to train the improved diffusion model to generate synthetic polyp images that closely resemble the real images. The improved diffusion model generates synthetic mask
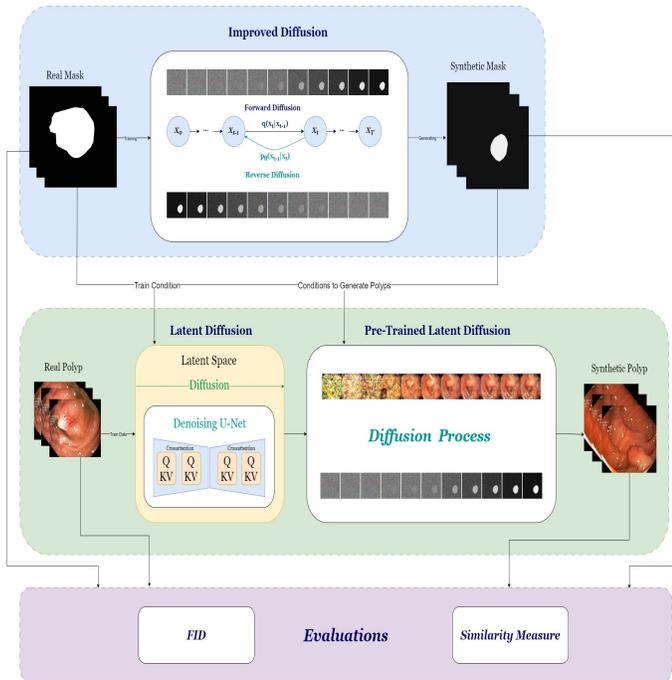
**Figure 1: The whole pipeline of generating synthetic polyps and mask. The blue box represents the diffusion model trained to generate realistic synthetic polyp masks. The green box represents the conditional latent diffusion model which is used to generate synthetic polyp conditioned on input masks. The bottom box represents the evaluation matrices.**

images by first adding noise to a randomly selected mask image from the training set. This noise would then be gradually reversed through multiple steps until a synthetic mask image is generated.

The advantage of using *improved diffusion models* to generate synthetic images is that we can overcome the issue of limited annotated data and train machine learning models more effectively. This can lead on to improvements in the accuracy and efficiency of CAD systems for detecting polyps in the GI-tract [10]

## 3.2 Latent diffusion model

The Latent Diffusion model [31], developed by CompVis and trained on the LAION-400M dataset [32], operates through a series of denoising autoencoders and diffusion models. This model has been utilized to generate images based on text prompts, and has shown exceptional results in tasks related to image inpainting [9] and various other applications, surpassing the current state of the art [30].

Latent diffusion models are a suitable choice for generating synthetic images of the GI tract for several reasons. Firstly, they possess the ability to model intricate and non-linear patterns in the data, crucial for producing convincing images of the GI tract. Secondly, they are capable of generating a large diversity of high-quality synthetic images, which enhances the generalizability of machine learning models. Lastly, they can be trained with a limited amount

of real data points, which is important in medical imaging where annotated data is often scarce.

## 3.3 Mask similarity

To assess the quality of generated images for generative models, the Fréchet inception distance (FID) [12] metric is typically used, which compares the distribution of the generated images compared to real images. Because the improved diffusion model are generating binary masks of polyps, we can also introduce the similarity measure (SIM) metric that is analogous to accuracy. Consider real image $r$ together with generated image $g$ of the same size, then the similarity $sim$ is defined as number of pixels that are same for both images divided by the total area of image:

$$sim(r, g) = \frac{\#(r == g)}{width * height}$$

To measure the similarity $SIM$ of generated images to our training real images $R$, we simply take average of the closest pairs (high similarities). If a generated image is $g$ and associated closest real image is $g*$, we can calculate largest similarity using,

$$SIM(R, G) = \frac{1}{|G|} \sum_{r \in R} sim(r, g*)$$

The idea here is that even if the generated images are highly similar to some of the training images (same size, position), they should differ in another aspects, such as rotation.

## 3.4 Segmentation models

We have used three different well-known image segmentation models, namely UNet++ [44], feature pyramid network (FPN) [22], and DeepLabv3+ [5] for evaluating the effect of using synthetic data for training polyp segmentation tasks. Initially, we trained these three models using three different approaches, i.e., we trained the system using i) 700 of real polyp images; ii) using 1000 synthetic polyp images; and iii) a combination of 700 real and 1000 synthetic polyp images. To further analyze the effect of synthetic data, we trained these three models with another set of real and synthetic data combinations. In these combinations, we fixed the number of real images to 100 samples, and we increased the number of synthetic samples from 0 to 1000 sequentially in steps of 100. The main objective of this experiment is to identify the effect of a number of synthetic samples included in the training data. However, we limited this experiment to using only 100 real images because of the time limitation, but in the future we will test with different number of real images from 200 to 700 to find the optimal combination to get better performance.

We tested these models with 300 real images and masks (from the segmentation data of HyperKvasir dataset [4]) which were not used to train either the diffusion model or the segmentation models. Then, we measured micro and micro-imagewise IOU, F1, Accuracy, and Precision from the test dataset for all the segmentation models. Micro values were calculated by summing true positive (TP), false positive (FP), false negative (FN), and true negative (TN) pixels over all images and all classes and then computing scores. In contrast, the micro-imagewise matrices were calculated by summing TP, FP, FN, and TN pixels for each image and then computing scores for

| Iter | 0 | 50k | 100k | 150k | 200k | 230k |
|------|------|------|------|------|------|------|
| FID | 140.14 | 128.95 | 117.14 | 105.63 | 88.41 | 141.44 |
| SIM | 88.22 | 89.46 | 90.81 | 91.31 | 92.49 | 88.38 |

**Table 1: Comparison of mask models based on FID, SIM**

| Epoch | 88 | 103 | 135 | 892 | 913 | 922 |
|-------|------|------|------|------|------|------|
| FID | 119.34 | 113.83 | 104.78 | 112.66 | 150.97 | 150.85 |

**Table 2: Comparison of polyp models based on FID**

each image. Finally, average scores over the dataset were calculated. In the micro-imagewise calculations, all images contributed equally to the final score. However, the second method takes into account class imbalance for each image.

## 4 RESULTS AND DISCUSSION

In this section, we discuss experiment setup and the result collected from generative models and segmentation models. A server with Nvidia A100 $80GB$ graphic processing units (GPUs), AMD EPYC 7763 64-cores processor with $2TB$ RAM were used for all the experiments of this study. Additionally, we used Pytorch [26], the Pytorch-lightning libraries, and the Pytorch segmentation library [14] as development frameworks.

### 4.1 Diffusion experiments and results

For mask generator, that is improved diffusion model we have used FID and SIM values to quantify and select appropriate model. We have generated 1000 masks for each of our saved model, and we compare them with 1000 real training masks in Table 1.

We selected the model from iteration $200,000$ based on the results from Table 1. The reason is that the model achieves lowest FID value together with high SIM values, indicating diverse and quality masks. We also inspected the generated masks visually to confirm this conclusion.

Examples of generated masks in comparison with real masks can be seen in Figure 2. We can see different masks with different shapes and numer of polyps indicating capability to generate diverse synthetic masks. Further discussion with medical professional would be required in order to determine if masks are correct.

Interestingly, high $SIM$ score doesn't necessary imply that model is producing identical masks, as can be seen in the Figure 3. For instance, masks may be located in similar positions but have different, smaller shapes, therefore achieving higher similarities.

We have used the generated masks made by the selected model as conditions to our latent polyp diffusion model and produced 1000 generated images which we used for further evaluation in Table 2.

We can see from Table 2 the model which achieved lowest FID score is at $Epoch = 135$. We inspected the generated images, similarly as in Figure 4. It can be seen that the quality of generated images deteriorates at later stages of training, reason may be overfitting. This may lead to problems while generating different samples with same condition which would be more similar.

Therefore, we selected the model from earlier stages that achieved lowest FID. We conditioned the model on one mask and generated multiple samples to see if the model generalizes well, results of this experiment can be seen in Figure 5.

### 4.2 Segmentation experiments and results

We used a learning rate of 0.0001 with the Adam optimizer [19] to train the three segmenation models, UNet++, FPN and DeepLabv3+. DiceLoss [35] was used in the training process as the loss function to update the weights. The encoder model of *resnet34* was input as the encoder network for all three models (for more details of these encoder networks, please refer to the documentation [14]). Micro metrics and micro-imagewise metrics (as discussed in the Pytorch segmentation library) were calculated from the best checkpoints and the test dataset after training 50 epochs for all the models. The calculated micro metrics are tabulated in Table 3, and micro-imagewise values are tabulated in Table 4.

According to the results in Tables 3 and 4, it is clear that adding synthetic data can improve the results of segmentation models. However, it is not always true because some models like FPN and UNet++ show the best IOU, F1, and accuracy when the training data consists of only the real data. In contrast, DeepLabv3 shows the best performance when some synthetic data is included in the training data. Overall, the best micro-imagewise IOU of 0.7751 is achieved from DeepLabv3+ when the training data contains the maximum number of images from both the real data and the synthetic data. Therefore, it is a clear evidence that synthetic data has a direct influence on the final performance of segmentation models. Moreover, we noticed that precision is always better when the synthetic data is in the training data than using only the real data. This implies that synthesized TP can improve the TP predictions which is more important in the medical domain. More visual comparisons are presented in Figure 6. This figure compares the model predictions from the three segmentation models between baseline performance and improved performance (marked using * in Figure 6) using synthetic data. The improved versions of the models were selected using the metrics of Tables 3 and 4.

Another interesting finding of these segmentation experiments is that we get the best values of precision, accuracy, F1 and IOU when we use a smaller number of real images and synthetic images. For example, Unet++ and FPN shows best precision values (micro and micro-imagewise) when the training data consist of 100 real samples and 200 synthetic samples . However, this implies that there is a direct correlation between synthetic data, models, and the number of parameters. Therefore, researchers should not conclude performance gain or degrade of using synthetic data to train segmentation models just by evaluating a single model.

## 5 CONCLUSION AND FUTURE WORKS

In this study, we used a probabilistic diffusion model-based method to generate synthetic polyp images conditioned on synthetic polyp masks. Our visual and quantitative comparisons show that the generated synthetic data are unique and realistic and not a simple copy of the training data used for the training. Our further analysis of using synthetic data to train polyp segmentation models shows that synthetic data can be used to improve the performance
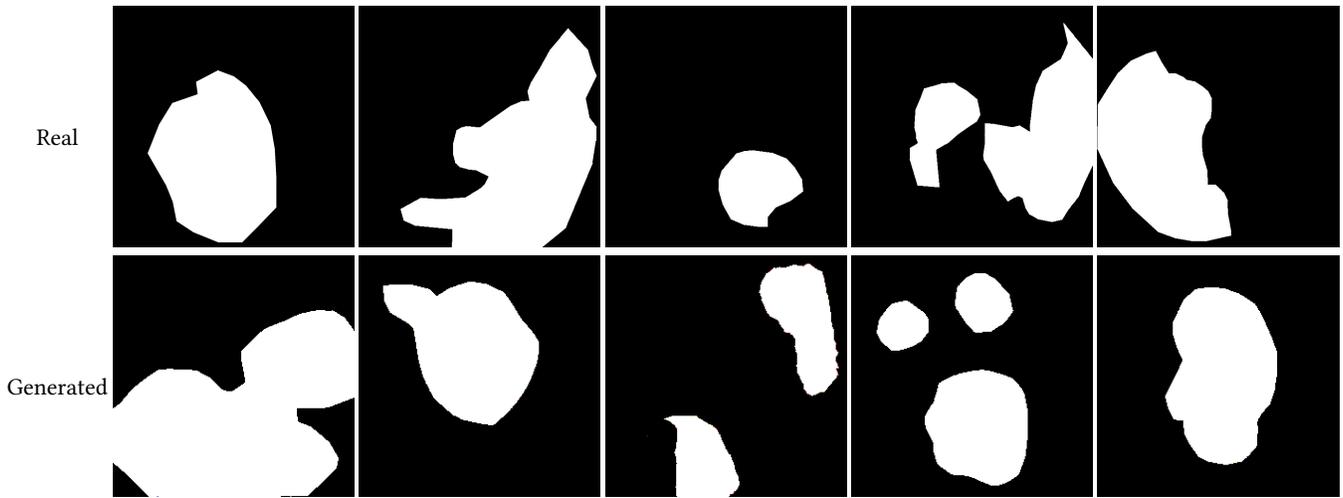
Figure 2: Examples of real masks in the first row, generated masks on the second. Note the variability of shapes and amount of polyps in the generated masks.



Figure 3: Examples of comparison of generated masks $g$ to real masks $r$ based on similarity measure $sim(r, g)$.



Figure 4: Generated synthetic polyps conditioned on the same mask illustrating changes in quality during training stages.

of segmentation models while these improvements are correlated with model architectures. In this regard, we can clearly conclude that synthetic data help to improve the performance of segmentation models. However, deep evaluations should be performed with

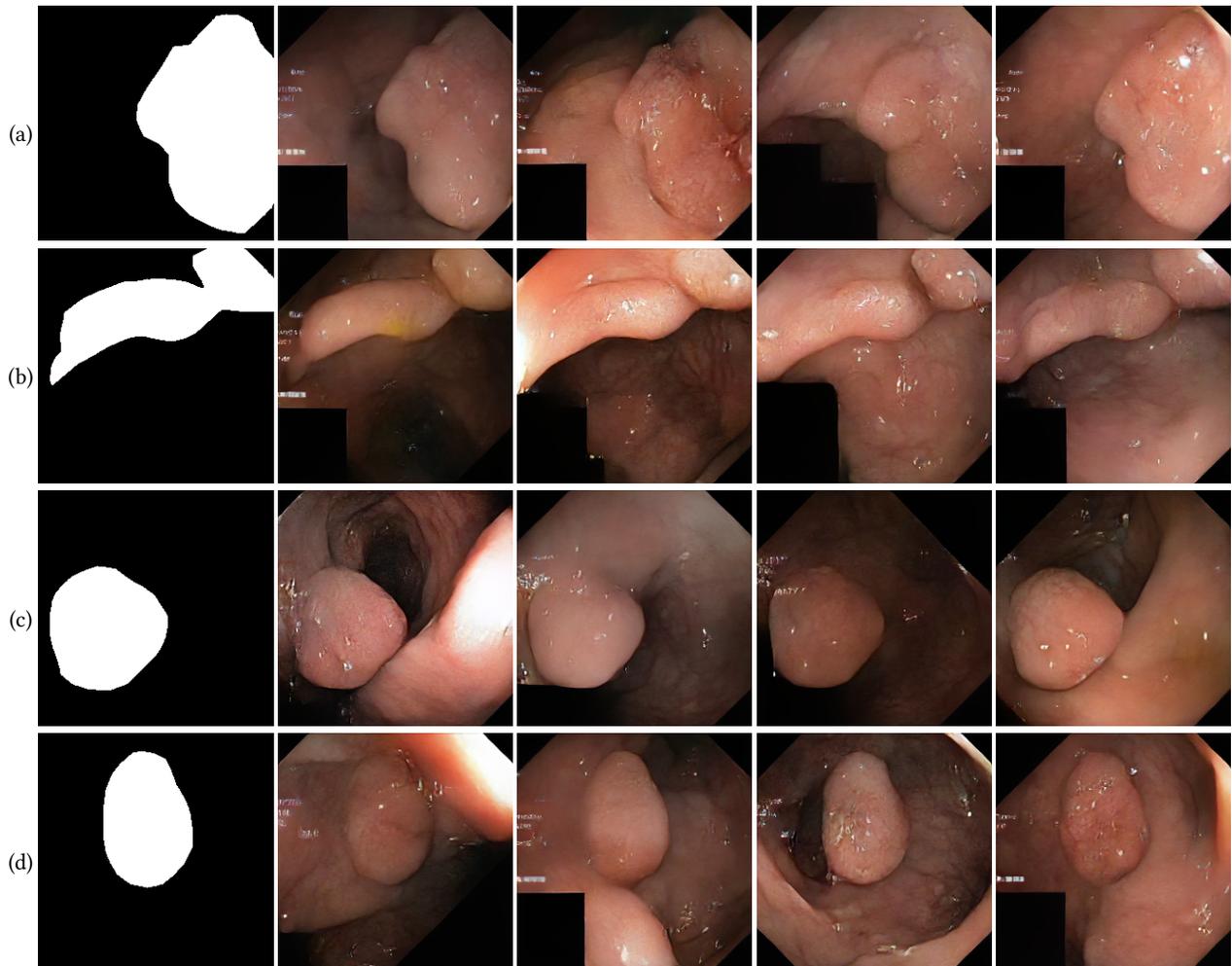**Figure 5: Generated synthetic polyp images from our conditional probabilistic diffusion model. The first column shows input conditions to the latent diffusion model. All other columns show the corresponding stochastic polyps generations with different input noises.**

multiple model architecture to see the real gain of using synthetic data.

In future studies, we will perform more segmentation experiments to get a complete result set for Tables 3 and 4, for example, increasing the synthetic training data gradually with the full real dataset. Moreover, we will generate more synthetic data using our model to train the segmentation models with large synthetic datasets to evaluate the effect of using synthetic data deeply. Generating multiple images conditioned on the same input to train the segmentation models are an another limitation of the presented segmentation experiments. Furthermore, the quality of generated images can be improved using the style-transfer technique [11] as used in the SinGAN-Seg study [41]. Cross-dataset evaluations should be performed to measure the effect of using synthetic data to train segmentation models to improve robustness and generalizability.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Prince Ebenezer Adjei, Zenebe Markos Lonseko, Wenju Du, Han Zhang, and Nini Rao. 2022. Examining the effect of synthetic data augmentation in polyp detection and segmentation. *International Journal of Computer Assisted Radiology and Surgery* 17, 7 (2022), 1289–1302.

[2] Hamed Alqahtani, Manolya Kavakli-Thorne, and Gulshan Kumar. 2021. Applications of generative adversarial networks (gans): An updated review. *Archives of Computational Methods in Engineering* 28 (2021), 525–552.

[3] Jorge Bernal, F Javier Sánchez, Gloria Fernández-Esparrach, Debora Gil, Cristina Rodríguez, and Fernando Vilariño. 2015. WM-DOVA maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from physicians. *Computerized medical imaging and graphics* 43 (2015), 99–111.

[4] Hanna Borgli, Vajira Thambawita, Pia H Smedsrud, Steven Hicks, Debesh Jha, Sigrun L Eskeland, Kristin Ranheim Randel, Konstantin Pogorelov, Mathias Lux,

**Table 3: Micro metrics calculated on the test dataset (300 real images and masks). The best value of each column is highlighted using <u>underlined</u> text. R = Real, Syn = Synthetic, Acc: = Accuracy and Prec: = Precision.**

| Model: | | Unet++ (26.1M) | | | | FPN (23.2M) | | | | DeepLabv3plus (22.4M) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| # R | # Syn | IOU | F1 | Acc: | Prec: | IOU | F1 | Acc: | Prec: | IOU | F1 | Acc: | Prec: |
| 700 | 0 | 0.7471 | 0.8552 | 0.9509 | 0.8535 | 0.7663 | 0.8677 | 0.9571 | 0.8623 | 0.7457 | 0.8543 | 0.9492 | 0.8699 |
| 0 | 1000 | 0.6852 | 0.8132 | 0.9375 | 0.8009 | 0.6784 | 0.8084 | 0.9276 | 0.7685 | 0.6580 | 0.7938 | 0.9301 | 0.7658 |
| 700 | 1000 | 0.7151 | 0.8339 | 0.9421 | 0.8235 | 0.7300 | 0.8439 | 0.9481 | 0.8323 | 0.7401 | 0.8506 | 0.9492 | 0.8252 |
| 100 | 0 | 0.6970 | 0.8215 | 0.9400 | 0.7912 | 0.6840 | 0.8123 | 0.9371 | 0.8209 | 0.6983 | 0.8224 | 0.9404 | 0.8304 |
| 100 | 100 | 0.6937 | 0.8192 | 0.9382 | 0.7692 | 0.7304 | 0.8442 | 0.9501 | 0.8509 | 0.7200 | 0.8372 | 0.9466 | 0.8305 |
| 100 | 200 | 0.7066 | 0.8281 | 0.9418 | 0.8804 | 0.7382 | 0.8494 | 0.9466 | 0.8763 | 0.7040 | 0.8263 | 0.9429 | 0.8383 |
| 100 | 300 | 0.7309 | 0.8445 | 0.9488 | 0.8536 | 0.7269 | 0.8419 | 0.9459 | 0.8219 | 0.7556 | 0.8608 | 0.9500 | 0.8521 |
| 100 | 400 | 0.6830 | 0.8116 | 0.9386 | 0.8333 | 0.7304 | 0.8442 | 0.9459 | 0.8375 | 0.7298 | 0.8438 | 0.9450 | 0.8342 |
| 100 | 500 | 0.6815 | 0.8106 | 0.9366 | 0.8152 | 0.7244 | 0.8402 | 0.9421 | 0.8209 | 0.7212 | 0.8380 | 0.9454 | 0.8427 |
| 100 | 600 | 0.7083 | 0.8292 | 0.9432 | 0.8287 | 0.7284 | 0.8429 | 0.9491 | 0.8668 | 0.7037 | 0.8261 | 0.9392 | 0.8405 |
| 100 | 700 | 0.7195 | 0.8369 | 0.9460 | 0.8420 | 0.7436 | 0.8530 | 0.9498 | 0.8107 | 0.7083 | 0.8292 | 0.9457 | 0.8347 |
| 100 | 800 | 0.6752 | 0.8061 | 0.9402 | 0.8495 | 0.7387 | 0.8497 | 0.9462 | 0.8278 | 0.7338 | 0.8465 | 0.9476 | 0.8770 |
| 100 | 900 | 0.7069 | 0.8283 | 0.9441 | 0.8319 | 0.7290 | 0.8432 | 0.9463 | 0.8234 | 0.7116 | 0.8315 | 0.9413 | 0.8171 |
| 100 | 1000 | 0.7513 | 0.8580 | 0.9506 | 0.8468 | 0.7214 | 0.8382 | 0.9457 | 0.8126 | 0.7154 | 0.8341 | 0.9401 | 0.8337 |

**Table 4: Micro-imagewise metrics calculated on the test dataset (300 real images and masks). These metrics take into account class imbalance. The best value of each column is highlighted using <u>underlined</u> text. R = Real, Syn = Synthetic, Acc: = Accuracy and Prec: = Precision.**

| Model: | | Unet++ (26.1M) | | | | FPN (23.2M) | | | | DeepLabv3plus (22.4M) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| # R | # Syn | IOU | F1 | Acc: | Prec: | IOU | F1 | Acc: | Prec: | IOU | F1 | Acc: | Prec: |
| 700 | 0 | 0.7551 | 0.8222 | 0.9509 | 0.8742 | 0.7681 | 0.8429 | 0.9571 | 0.8761 | 0.7528 | 0.8317 | 0.9492 | 0.8678 |
| 0 | 1000 | 0.7232 | 0.8013 | 0.9375 | 0.8128 | 0.6977 | 0.7903 | 0.9276 | 0.7757 | 0.7018 | 0.7896 | 0.9301 | 0.8062 |
| 700 | 1000 | 0.7442 | 0.8185 | 0.9421 | 0.8517 | 0.7371 | 0.8128 | 0.9481 | 0.8504 | 0.7751 | 0.8465 | 0.9492 | 0.8628 |
| 100 | 0 | 0.7136 | 0.8005 | 0.9400 | 0.7985 | 0.6587 | 0.7613 | 0.9371 | 0.8039 | 0.7116 | 0.8040 | 0.9404 | 0.8506 |
| 100 | 100 | 0.7146 | 0.7976 | 0.9382 | 0.7931 | 0.7357 | 0.8212 | 0.9501 | 0.8483 | 0.7183 | 0.8048 | 0.9466 | 0.8420 |
| 100 | 200 | 0.7433 | 0.8193 | 0.9418 | 0.8768 | 0.7000 | 0.7842 | 0.9466 | 0.8856 | 0.7197 | 0.8031 | 0.9429 | 0.8554 |
| 100 | 300 | 0.7392 | 0.8168 | 0.9488 | 0.8570 | 0.7302 | 0.8126 | 0.9459 | 0.8357 | 0.7337 | 0.8097 | 0.9500 | 0.8503 |
| 100 | 400 | 0.7097 | 0.7867 | 0.9386 | 0.8444 | 0.7512 | 0.8268 | 0.9459 | 0.8683 | 0.7366 | 0.8153 | 0.9450 | 0.8604 |
| 100 | 500 | 0.7088 | 0.7874 | 0.9366 | 0.8551 | 0.7376 | 0.8200 | 0.9421 | 0.8410 | 0.7264 | 0.8085 | 0.9454 | 0.8369 |
| 100 | 600 | 0.7238 | 0.7987 | 0.9432 | 0.8584 | 0.7348 | 0.8147 | 0.9491 | 0.8639 | 0.7054 | 0.7944 | 0.9392 | 0.8287 |
| 100 | 700 | 0.7230 | 0.7988 | 0.9460 | 0.8471 | 0.7319 | 0.8147 | 0.9498 | 0.8010 | 0.7317 | 0.8135 | 0.9457 | 0.8502 |
| 100 | 800 | 0.7081 | 0.7884 | 0.9402 | 0.8692 | 0.7502 | 0.8274 | 0.9462 | 0.8622 | 0.7548 | 0.8322 | 0.9476 | 0.8725 |
| 100 | 900 | 0.7244 | 0.8025 | 0.9441 | 0.8329 | 0.7441 | 0.8242 | 0.9463 | 0.8534 | 0.7314 | 0.8115 | 0.9413 | 0.8570 |
| 100 | 1000 | 0.7385 | 0.8145 | 0.9506 | 0.8495 | 0.7302 | 0.8086 | 0.9457 | 0.8386 | 0.7343 | 0.8125 | 0.9401 | 0.8574 |

Duc Tien Dang Nguyen, et al. 2020. HyperKvasir, a comprehensive multi-class image and video dataset for gastrointestinal endoscopy. *Scientific data* 7, 1 (2020), 283.

[5] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. 2017. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587* (2017).

[6] Richard J Chen, Ming Y Lu, Tiffany Y Chen, Drew FK Williamson, and Faisal Mahmood. 2021. Synthetic data in machine learning for medicine and healthcare. *Nature Biomedical Engineering* 5, 6 (2021), 493–497.

[7] Antonia Creswell, Tom White, Vincent Dumoulin, Kai Arulkumaran, Biswa Sengupta, and Anil A Bharath. 2018. Generative adversarial networks: An overview. *IEEE signal processing magazine* 35, 1 (2018), 53–65.

[8] Prafulla Dhariwal and Alexander Nichol. 2021. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems* 34 (2021), 8780–8794.

[9] Omar Elharrouss, Noor Almaadeed, Somaya Al-Maadeed, and Younes Akbari. 2020. Image inpainting: A review. *Neural Processing Letters* 51 (2020), 2007–2028.

[10] Jan Andre Fagereng, Vajira Thambawita, Andrea M Storås, Sravanthi Parasa, Thomas de Lange, Pål Halvorsen, and Michael A Riegler. 2022. PolypConnect: Image inpainting for generating realistic gastrointestinal tract images with polyps. In *2022 IEEE 35th International Symposium on Computer-Based Medical Systems (CBMS)*. IEEE, 66–71.

[11] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. 2016. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2414–2423.

[12] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems* 30 (2017).

[13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems* 33 (2020), 6840–6851.

[14] Pavel Iakubovskii. 2019. Segmentation Models Pytorch. https://github.com/qubvel/segmentation_models.pytorch.
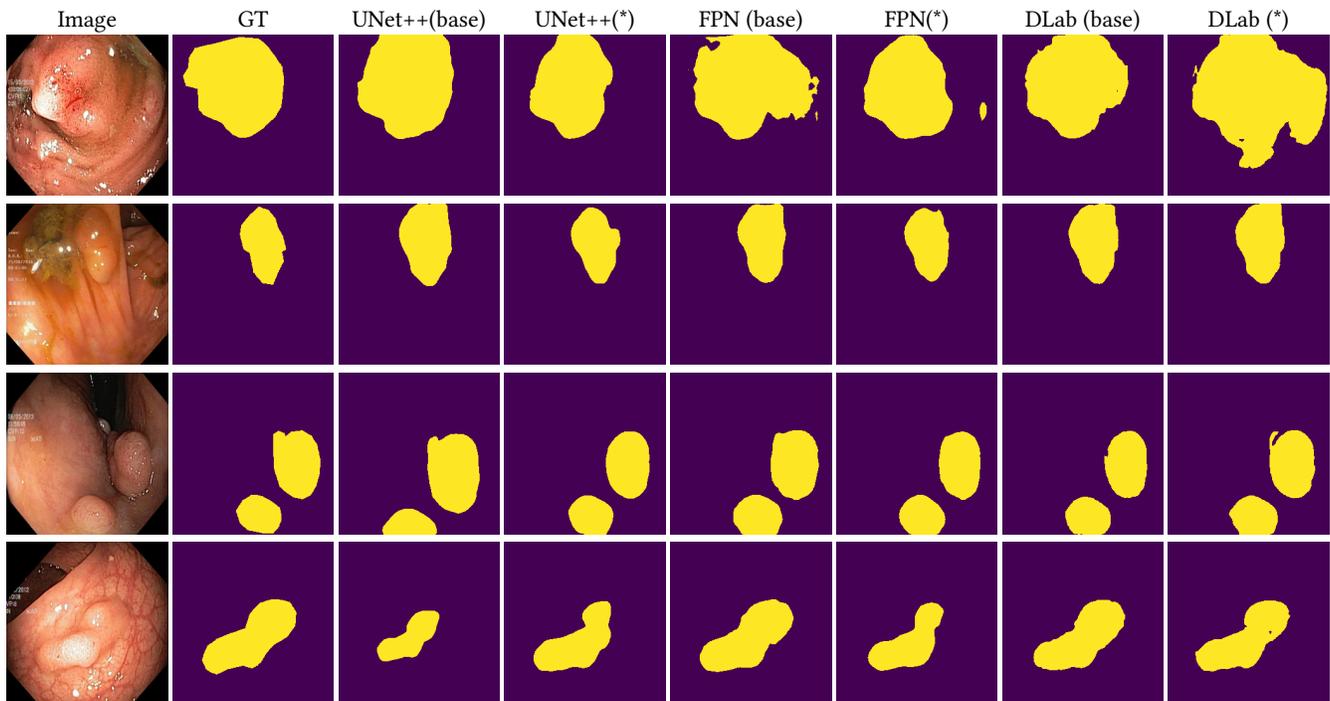
**Figure 6: Predicted masks from different segmentation models. The baseline predictions are from the model trained with only real data. Unet++(\*) is selected based on using the high IOU value in Table 3. FPN(\*) is selected using the highest Precision in Tables 3 and 4. DeepLabv3(\*)[Dlab(\*)] is selected using high IOU values in Table 4.**

[15] Gavriel Iddan, Gavriel Meron, Arkady Glukhovsky, and Paul Swain. 2000. Wireless capsule endoscopy. *Nature* 405, 6785 (2000), 417–417.

[16] Debesh Jha, Pia H Smedsrud, Michael A Riegler, Pål Halvorsen, Thomas de Lange, Dag Johansen, and Håvard D Johansen. 2020. Kvasir-seg: A segmented polyp dataset. In *MultiMedia Modeling: 26th International Conference, MMM 2020, Daejeon, South Korea, January 5–8, 2020, Proceedings, Part II 26.* Springer, 451–462.

[17] James Jordon, Lukasz Szpruch, Florimond Houssiau, Mirko Bottarelli, Giovanni Cherubin, Carsten Maple, Samuel N Cohen, and Adrian Weller. 2022. Synthetic Data–what, why and how? *arXiv preprint arXiv:2205.03257* (2022).

[18] Michal F Kaminski, Jaroslaw Regula, Ewa Kraszewska, Marcin Polkowski, Urszula Wojciechowska, Joanna Didkowska, Maria Zwierko, Maciej Rupinski, Marek P Nowacki, and Eugeniusz Butruk. 2010. Quality indicators for colonoscopy and the risk of interval cancer. *New England journal of medicine* 362, 19 (2010), 1795–1803.

[19] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[20] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2017. Imagenet classification with deep convolutional neural networks. *Commun. ACM* 60, 6 (2017), 84–90.

[21] Catherine Le Berre, William J Sandborn, Sabeur Aridhi, Marie-Dominique Devignes, Laure Fournier, Malika Smail-Tabbone, Silvio Danese, and Laurent Peyrin-Biroulet. 2020. Application of artificial intelligence to gastroenterology and hepatology. *Gastroenterology* 158, 1 (2020), 76–94.

[22] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. 2017. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition.* 2117–2125.

[23] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13.* Springer, 740–755.

[24] Min Min, Song Su, Wenrui He, Yiliang Bi, Zhanyu Ma, and Yan Liu. 2019. Computer-aided diagnosis of colorectal polyps using linked color imaging colonoscopy to predict histology. *Scientific reports* 9, 1 (2019), 1–8.

[25] Alexander Quinn Nichol and Prafulla Dhariwal. 2021. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning.* PMLR, 8162–8171.

[26] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32.* Curran Associates, Inc., 8024–8035. http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf

[27] Konstantin Pogorelov, Kristin Ranheim Randel, Carsten Griwodz, Sigrun Losada Eskeland, Thomas de Lange, Dag Johansen, Concetto Spampinato, Duc-Tien Dang-Nguyen, Mathias Lux, Peter Thelin Schmidt, et al. 2017. Kvasir: A multiclass image dataset for computer aided gastrointestinal disease detection. In *Proceedings of the 8th ACM on Multimedia Systems Conference.* 164–169.

[28] Hemin Ali Qadir, Ilangko Balasingham, and Younghak Shin. 2022. Simple U-net based synthetic polyp image generation: Polyp to negative and negative to polyp. *Biomedical Signal Processing and Control* 74 (2022), 103491.

[29] Michael Riegler, Konstantin Pogorelov, Pål Halvorsen, Thomas de Lange, Carsten Griwodz, Peter Thelin Schmidt, Sigrun Losada Eskeland, and Dag Johansen. 2016. Eir—efficient computer aided diagnosis framework for gastrointestinal endoscopies. In *2016 14th International Workshop on Content-Based Multimedia Indexing (CBMI).* IEEE, 1–6.

[30] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2021. High-Resolution Image Synthesis with Latent Diffusion Models. arXiv:2112.10752 [cs.CV]

[31] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 10684–10695.

[32] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. 2021. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114* (2021).

[33] Juan Silva, Aymeric Histace, Olivier Romain, Xavier Dray, and Bertrand Granado. 2014. Toward embedded detection of polyps in wce images for early diagnosis of colorectal cancer. *International journal of computer assisted radiology and surgery* 9 (2014), 283–293.

[34] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. 2015. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*. PMLR, 2256–2265.

[35] Carole H Sudre, Wenqi Li, Tom Vercauteren, Sebastien Ourselin, and M Jorge Cardoso. 2017. Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations. In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: Third International Workshop, DLMIA 2017, and 7th International Workshop, ML-CDS 2017, Held in Conjunction with MICCAI 2017, Québec City, QC, Canada, September 14, Proceedings 3*. Springer, 240–248.

[36] Nima Tajbakhsh, Suryakanth R Gurudu, and Jianming Liang. 2015. Automated polyp detection in colonoscopy videos using shape and context information. *IEEE transactions on medical imaging* 35, 2 (2015), 630–644.

[37] Vajira Thambawita, Steven A Hicks, Jonas Isaksen, Mette Haug Stensen, Trine B Haugen, JØrgen Kanters, Sravanthi Parasa, Thomas de Lange, Håvard D Johansen, Dag Johansen, et al. 2021. DeepSynthBody: the beginning of the end for data deficiency in medicine. In *2021 International Conference on Applied Artificial Intelligence (ICAPAI)*. IEEE, 1–8.

[38] Vajira Thambawita, Steven A. Hicks, Jonas Isaksen, Mette Haug Stensen, Trine B. Haugen, JØrgen Kanters, Sravanthi Parasa, Thomas de Lange, Håvard D. Johansen, Dag Johansen, Hugo L. Hammer, Pål Halvorsen, and Michael A. Riegler. 2021. DeepSynthBody: the beginning of the end for data deficiency in medicine. In *2021 International Conference on Applied Artificial Intelligence (ICAPAI)*. 1–8. https://doi.org/10.1109/ICAPAI49758.2021.9462062

[39] Vajira Thambawita, Jonas L Isaksen, Steven A Hicks, Jonas Ghouse, Gustav Ahlberg, Allan Linneberg, Niels Grarup, Christina Ellervik, Morten Salling Olesen, Torben Hansen, et al. 2021. DeepFake electrocardiograms using generative adversarial networks are the beginning of the end for privacy issues in medicine. *Scientific reports* 11, 1 (2021), 21896.

[40] Vajira Thambawita, Debesh Jha, Hugo Lewi Hammer, Håvard D Johansen, Dag Johansen, Pål Halvorsen, and Michael A Riegler. 2020. An extensive study on cross-dataset bias and evaluation metrics interpretation for machine learning applied to gastrointestinal tract abnormality classification. *ACM Transactions on Computing for Healthcare* 1, 3 (2020), 1–29.

[41] Vajira Thambawita, Pegah Salehi, Sajad Amouei Sheshkal, Steven A Hicks, Hugo L Hammer, Sravanthi Parasa, Thomas de Lange, Pål Halvorsen, and Michael A Riegler. 2022. SinGAN-Seg: Synthetic training data generation for medical image segmentation. *PloS one* 17, 5 (2022), e0267976.

[42] Vajira L Thambawita, Inga Strümke, Steven Hicks, Michael A Riegler, Pål Halvorsen, and Sravanthi Parasa. 2021. ID: 3523524 Data augmentation using generative adversarial networks for creating realistic artificial colon polyp images: validation study by endoscopists. *Gastrointestinal Endoscopy* 93, 6 (2021), AB190.

[43] Daniela Guerrero Vinsard, Yuichi Mori, Masashi Misawa, Shin-ei Kudo, Amit Rastogi, Ulas Bagci, Douglas K Rex, and Michael B Wallace. 2019. Quality assurance of computer-aided detection and diagnosis in colonoscopy. *Gastrointestinal endoscopy* 90, 1 (2019), 55–63.

[44] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. 2018. Unet++: A nested u-net architecture for medical image segmentation. In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: 4th International Workshop, DLMIA 2018, and 8th International Workshop, ML-CDS 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 20, 2018, Proceedings 4*. Springer, 3–11.