# Unsupervised Reverse Domain Adaptation for Synthetic Medical Images via Adversarial Training

Faisal Mahmood[1]     Richard Chen[2]     Nicholas J. Durr[1]

[1]Department of Biomedical Engineering     [2]Department of Computer Science

Johns Hopkins University (JHU)

{faisalm, rchen40, ndurr}@jhu.edu

## Abstract

*To realize the full potential of deep learning for medical imaging, large annotated datasets are required for training. Such datasets are difficult to acquire because labeled medical images are not usually available due to privacy issues, lack of experts available for annotation, underrepresentation of rare conditions and poor standardization. Lack of annotated data has been addressed in conventional vision applications using synthetic images refined via unsupervised adversarial training to look like real images. However, this approach is difficult to extend to general medical imaging because of the complex and diverse set of features found in real human tissues. We propose an alternative framework that uses a reverse flow, where adversarial training is used to make real medical images more like synthetic images, and hypothesize that clinically-relevant features can be preserved via self-regularization. These domain-adapted images can then be accurately interpreted by networks trained on large datasets of synthetic medical images. We test this approach for the notoriously difficult task of depth-estimation from endoscopy. We train a depth estimator on a large dataset of synthetic images generated using an accurate forward model of an endoscope and an anatomically-realistic colon. This network predicts significantly better depths when using synthetic-like domain-adapted images compared to the real images, confirming that the clinically-relevant features of depth are preserved.*

## 1. Introduction

Deep Learning offers great promise for the reconstruction and interpretation of medical images [27, 6]. Countless applications in clinical diagnostics, disease screening, interventional planning, and therapeutic surveillance rely on the subjective interpretation of medical images from healthcare providers. This approach is costly, time-intensive, and has well-known accuracy and precision limitations—all of
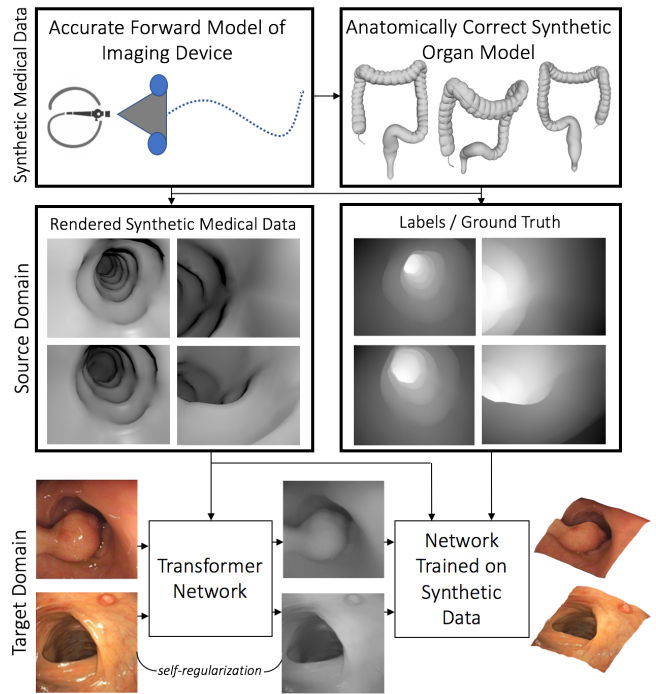


Figure 1. Unsupervised reverse domain adaption for endoscopy images. We use an accurate forward model of an endoscope and an anatomically correct colon model to generate synthetic endoscopy images with ground truth depth. This large synthetic dataset can be used to train a deep network for depth estimation. An adversarial network transforms input endoscopy images to a synthetic-like representation while preserving clinically relevant features via self-regularization. These synthetic-like images can be directly used for depth estimation from the network trained on synthetic images.

which could be mitigated by objective, automatic image analysis.

For conventional images, deep learning has achieved remarkable performance for a variety of computer vision tasks, typically by utilizing large sets of real-world im-

ages for training, such as ImageNet [12], COCO [13] and Pascal VOC [4]. Unfortunately, the potential benefits of deep learning have yet to transfer to the most critical needs in medical imaging because there are no large, annotated dataset of medical images available. Despite the compelling need for such a dataset, there are practical concerns that impede its development, including the cost, time, expertise, privacy, and regulatory issues associated with medical data collection, annotation, and dissemination.

The obstacles associated with developing a large dataset of real images can be circumvented by generating synthetic images [22, 26, 24]. Considerable effort has been devoted to adapting models generated with synthetic data as the source domain to real data as the target domain [2]. Advances in adversarial training have sparked interest in making synthetic data look more realistic via unsupervised adversarial learning (SimGAN) [28]. In the medical imaging domain, there has been recent success in generating realistic synthetic data for the relatively constrained problem of 2D retinal imaging using standard GANs [7]. In more complex applications, it is challenging to generate an appropriate span of synthetic medical images for training, because few models exist that accurately simulate the anatomical complexity and diversity found in healthy to pathologic tissues. Moreover, the forward models for medical imaging devices are more complex than those used in many conventional vision applications. Consequently, models trained on synthetic medical data may fail to generalize to real medical images, where accurate interpretation may be critically important.

Cross-patient network usage is a well-known challenge to learning-based medical imaging methods. Often a network trained on data from one patient fails to generalize to other patients. This is commonly observed for optical imaging methods, such as endoscopy, which capture both low- and high-level texture details of the patient. Low-level texture details are patient-specific and not diagnostic, such as vascular patterns. High-level texture, on the other hand, contains clinically-relevant features that should be generalized across patients. This complication makes it difficult for methods like SimGAN [28] to work both accurately and generally because the span of realistic images produced will be similar to the real images used for training.

In this work, we propose to reverse the flow of traditional adversarial training-based domain adaption. Instead of changing synthetic images to appear realistic [28], we transform real images to look more synthetic (Fig. 1). We train an adversarial transformation network which transforms real medical images to a synthetic-like representation while preserving clinically-relevant information. In summary, we can train solely on synthetic medical data as the source domain and transform real data in the target domain to a more synthetic interpretation, thus bridging the gap be-

tween the source and target domains in a reverse manner.

To transform real images to a synthetic-like representation, we train a transformer with an adversarial loss similar to GANs [5] and SimGAN [28]. However, unlike SimGAN that trains for inducing realism to synthetic data, we train for a synthetic-like representation of real data. With the roles of synthetic and real data reversed, the overall transformer architecture is similar to a standard GAN and is composed of a transformer network that tries to fool a discriminator network into thinking that the transformed medical image is synthetic. In addition to removing patient specific details from the data, the synthetic image should preserve enough information within the data that it could be used for the task at hand. To preserve this information a fully connected network is used and the adversarial loss is complemented with a self-regularization term which constrains the amount of deviation from the real image.

#### Contributions
1. We propose an adversarial training-based reverse domain adaptation method which uses unlabeled synthetic data to transform real data to a synthetic-like representation while maintaining clinically relevant diagnostic features via self-regularization.
2. **Synthetic Endoscopy Data Generation:** We generate a large dataset of perfectly-annotated synthetic endoscopy images from an endoscope forward model and an anatomically correct colon model.
3. **Reverse Domain Adaptation:** We train a transformer network via adversarial training composed of a generator which generates the a synthetic-like representation of real endoscopy images. The loss function in the generator contains a discriminator to classify the endoscopy images as real or synthetic and a self-regularization term that penalizes large deviations from the real image.
4. **Qualitative and Quantitative Study:** We validate our domain adaptation approach by using synthetically generated endoscopy data to train a monocular endoscopy depth estimation network and quantitatively testing it with real endoscopy data from: a) Colon Phantom b) Porcine Colon and qualitatively testing it with real human endoscopy data. We further show that the depth obtained from training on synthetic data can be used to improve state-of-the-art results on polyp segmentation.

## 2. Related Work

**Navigating Limited Medical Imaging Data:** Improving the performance of deep learning methods with limited data is an active research area. Standard data augmentation has been used for medical imaging for the past years. Ronneberger *et al.* [23] demonstrated success with using elastic

augmentation with U-Net architectures for medical image segmentation. Payer *et al.* [19] have demonstrated incorporating application specific *a priori* information can train better deep networks. There is a growing interest in transferring knowledge from networks trained for conventional vision to the medical imaging domain [36]. However, the major limiting factor with all these approaches is the fact that there is very limited data to train from.

**Generative Adversarial Networks:** The GAN framework was first presented by Goodfellow *et. al.* in [5] and was based on the idea of training two networks, a generator and a discriminator simultaneously with competing losses. While the generator learns to generate realistic data from a random vector, the discriminator classifies the generated image as real or fake and gives feedback to the generator. Once the training reaches equilibrium the generator is able to fool the discriminator every time it generates a new image. Initially GANs were applied to the MINST dataset [5] but recently the framework has been refined and used for a variety of applications [25]. Models with adversarial losses have been used for synthesis of 3D shapes, image-to-image translation, for generating radiation patterns etc. Recently, Zhu *et al.* [36] proposed iGAN which enables interactive image manipulation on a natural image manifold. Shrivistava *et al.* [28] have proposed an unsupervised method for refining synthetic images to look more realistic using a modified adversarial training framework.

**Adversarial Training for Biomedical Imaging:** Various kinds of adversarial training has recently been used for a variety of medical imaging tasks including noise reduction [33], segmentation [35, 16], detection [11], reconstruction [15], classification [34] and image synthesis [18]. Osokin *et al.* [18] use GANs for synthesizing biological cells imaged by fluorescence microscopy. Costa *et al.* [3] and Guibas *et al.* [7] synthesize retinal images using adversarial training. All current adversarial training-based image synthesis methods attempt to generate realistic images from a random noise vector or refine synthetic images to create more realistic images, in contrast our method, transforms real images to a synthetic-like representation allowing the desired network to be trained only on synthetic images.

## 3. Generating Synthetic Medical Data

Despite the widespread use of synthetic data for training deep networks for real world images [29, 8, 31, 20], its use for medical imaging applications has been relatively limited. Unlike conventional real-world images that may contain a constrained span of object diversity, medical images capture information of biological tissues which contain unique patient-specific texture that is difficult to model. We therefore propose a frame work where we generate a large dataset of medical images with this patient-specific detail removed so that a network can be trained on universal diag-
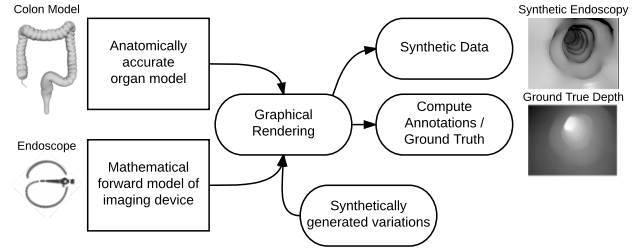


Figure 2. General framework for generating synthetic medical imaging data with endoscopy as an example.

nostic features. In general, this synthetic data can be generated by (Fig. 2):

1. Developing an accurate forward model for the medical imaging device.

2. Generating an anatomically accurate model of the organ being imaged.

3. Rendering images from a variety of positions, angles and parameters.

Typically, forward models for medical imaging devices is more complicated as compared to typical cameras, and anatomically accurate models need to represent a high degree of variation and rare conditions to cater for a diverse set of patients.

**Synthetic Endoscopy Data with Ground Truth Depth:** For the purpose of demonstration of our proposed methods we focus on the task of depth estimation from monocular endoscopy images. This is a notoriously difficult problem because of the lack of clinical images with available ground truth data, since it is difficult to include a depth sensor on an endoscope. We generate synthetic data to overcome this issue. We develop a forward model of an endoscope with a wide-angle monocular camera and two to three light sources that exhibit realistic inverse square law intensity fall-off. We use a synthetically generated and anatomically accurate colon model and image it using the virtual endoscope placed at a variety of angles and varying conditions to mimic the movement of an actual endoscope. We also generate pixel-wise ground truth depth for each rendered image. We finally create a dataset with 260,000 images with ground truth depth. Although this large dataset of images is able to train efficient deep networks these networks are not effectively generalizable to real world images.

## 4. Reverse Domain Adaptation

**Transformer Loss:** Formally, the goal of our proposed reverse domain adaptation method is to use a set of synthetic images $g_i \in \mathcal{G}$ to learn a transformer $x^{'} = \mathcal{T}_{\gamma_t}(x)$

that can transform real images $\boldsymbol{x}$ to a synthetic-like representation $\boldsymbol{x}^{'}$. The transformer should be able to fool a discriminator $\mathcal{D}_{\gamma_d}$ where $\gamma_t$ and $\gamma_d$ are the learning parameters. There are three key requirements for this setup: a) The transformer output should only remove the patient specific details in the image, while preserving diagnostic features. b) The adversarial training should not introduce artifacts in the transformed image. c) The adversarial training should be stable. The transformer loss function can be defined as,

$$\mathcal{L}_{\mathcal{T}}(\gamma_t) = \sum_i \psi(\boldsymbol{x}_i, \mathcal{G}; \gamma_t) + \lambda \phi(\boldsymbol{x}_i; \gamma_t), \qquad (1)$$

where, $\psi$ forces the real image to a synthetic-like representation and $\phi$ penalizes large variations to preserve specific properties of the real image. $\lambda$ controls the amount of self-regularization enforced by $\phi$.

**Discriminator Loss:** In order to transform a real image to its synthetic-like counterpart, the gap between the representations of the real and synthetic image needs to be minimized. An ideal transformer should be able to produce an indistinguishable synthetic representation of a real image every time, which is possible if a discriminator is embedded within the transformers loss function (Fig. 3). As explained in [5, 25, 28], a discriminator is essentially a classifier that classifies the output of another network as real or fake. However, unlike [28], in our case the role of the discriminator is reversed—instead of enforcing the transformer to produce more realistic images the role of the discriminator is to enforce the transformer to produce synthetic images. The discriminator loss can be defined as follows

$$\mathcal{L}_{\mathcal{D}}(\gamma_d) = -\sum_i \log(\mathcal{D}_{\gamma_d}(\boldsymbol{x}^{'})) - \sum_j \log(1 - \mathcal{D}_{\gamma_d}(g_j)). \qquad (2)$$

This is essentially a two class classification problem with cross-entropy error where the first term represents the probability of the input being a synthetic image and the second term represents the probability of the input image being synthetic-like representation of a real image. The discriminator works on a patch level rather than the entire image to prevent artifacts.

To train our network, we randomly sample mini-batches of synthetic images and images transformed to be synthetic by the transformer. Instead of using individual outputs of the transformer we use randomly sampled, buffered outputs and a set of randomly sampled synthetic images. This increases the stability of the adversarial training since the lack of memory can diverge the adversarial training and introduce artifacts [28]. At each step the discriminator trains using this mini-batch and parameters $\gamma_d$ are updated using stochastic gradient decent (SGD). The transformer loss is then updated with the trained discriminator, the $\psi$ term in Eq. 1 can be defined as,
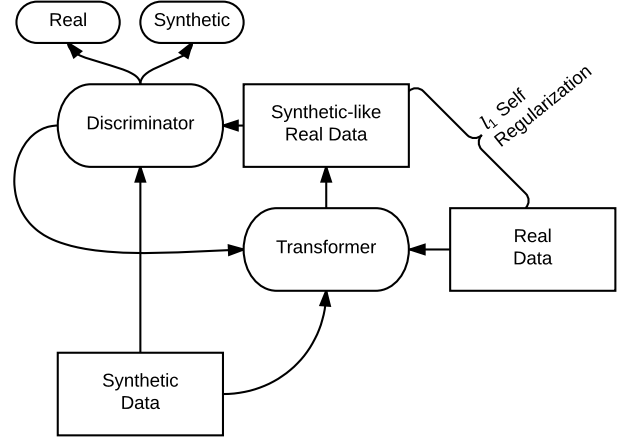


Figure 3. An overview of our proposed adversarial training architecture. Real data is transformed into a synthetic-like representation using a transformer network that minimizes an adversarial loss term and a self-regularization term. The discriminator acts as a classifier to identify the image as real or synthetic and gives feedback to the transformer via the adversarial loss. The total loss essentially defines how well the discriminator is tricked into believing that the transformed real image is synthetic and how close the transformed image is to the real image.

$$\psi(\boldsymbol{x}_i, \mathcal{G}; \gamma_t) = -\log(1 - \mathcal{D}_{\gamma_d}(\mathcal{T}_{\gamma_t}(\boldsymbol{x}))). \qquad (3)$$

As the training shuffling between the transformer and discriminator reaches equilibrium the transformer is able to fool the discriminator every time. The loss in Eq. 3 forces the discriminator to fail to classify transformed images as synthetic-like real.

**Self-Regularization:** As mentioned earlier, a key requirement for the transformer is that it should only remove patient specific data and should preserve other features such as shape. For the proof-of-concept proposed in this work, we utilize a simple, per-pixel loss term between the real image and the synthetic-like real representation of the image to penalize the transformed image from deviating significantly from the real image. The self regularization term $\phi$ can be defined as,

$$\phi(\boldsymbol{x}_i; \gamma_t) = || \Phi(\mathcal{T}_{\gamma_t}(\boldsymbol{x})) - \Phi(\boldsymbol{x}) ||_1, \qquad (4)$$

where $\Phi$ represents the feature transform and $|| \cdot ||_1$ represents the $\ell_1$ norm.

The transformer loss term can be rewritten as,

$$\begin{aligned} \mathcal{L}_{\mathcal{T}}(\gamma_t) = &-\sum_i \log(1 - \mathcal{D}_{\gamma_d}(\mathcal{T}_{\gamma_t}(\boldsymbol{x}))) \\ &+\lambda || \Phi(\mathcal{T}_{\gamma_t}(\boldsymbol{x})) - \Phi(\boldsymbol{x}) ||_1, \end{aligned} \qquad (5)$$

In summary, the total loss measures how well the discriminator is tricked into believing that the transformed real

**Algorithm 1** Adversarial training of a Transformer $x^{'} = \mathcal{T}_\gamma(x)$

---

**INPUT:** Synthetic Data: $\boldsymbol{g}_i \in \mathcal{G}$, Real Data: $\boldsymbol{x}_i \in \mathcal{X}$, Transformer Updates/step: $n_t$, Discriminator Updates/step: $n_d$
1: **for** $s = 1, 2, 3...S$ **do**
2:     **for** $n_t = 1, 2, 3...N_t$ **do**
3:         Sample a mini-batch $\{x_1, x_2, ...x_k\}$ of $k$ real images.
4:         Update the transformer network parameters $\gamma_t$ by taking an SGD step:
    $\nabla_{\gamma_t} \frac{1}{k} \sum_i \psi(\boldsymbol{x}_i, \mathcal{G}; \gamma_t) + \lambda\phi(\boldsymbol{x}_i; \gamma_t)$
5:     **end for**
6:     **for** $n_d = 1, 2, 3...N_d$ **do**
7:         Sample a mini-batchs of $k$ synthetic images $\{g_1, g_2, ...g_k\}$ and transformed real images $\{x_1, x_2, ...x_k\}$.
8:         $\boldsymbol{x}_i^{'} \leftarrow \mathcal{T}_\gamma(\boldsymbol{x}_i)$
9:         Update the discriminator network parameters $\gamma_d$ by taking an SGD step:
    $-\nabla_{\gamma_d} \frac{1}{k} \sum_i \log(\mathcal{D}_{\gamma_d}(\boldsymbol{x}^{'})) - \sum_j \log(1 - \mathcal{D}_{\gamma_d}(g_j))$
10:     **end for**
11: **end for**
**OUTPUT:** Trained Transformer Model $\mathcal{T}_\gamma(x)$

---

image is synthetic, and how close the transformed image is to the real image. This overall training process has been explained in detail in Algorithm 1.

## 5. Depth Estimation from Monocular Endoscopy Images

The previous sections have talked about generating synthetic medical data and adversarial training to bring real images within the domain of the synthetic data via a reverse domain adaptation pipeline. In order to evaluate the effectiveness of our proposed reverse domain adaptation pipeline we train a network from synthetically generated endoscopy data (Fig. 2) and demonstrate that it can be adapted to three different target domains. By demonstrating that distribution of the target domain can be brought closer to the source domain via adversarial training essentially showing that our depth estimation paradigm is domain independent.

Once the synthetic data with ground truth depths is generated we use a CNN-CRF based depth estimation framework described in [14]. Assuming $\boldsymbol{g} \in \mathbb{R}^{n \times m}$ is a synthetic endoscopy image which has been divided into $p$ superpixels and $\boldsymbol{y} = [y_1, y_2, ..., y_p] \in \mathbb{R}$ is the depth vector for each super-pixel. In this case, the conditional probability distribution of the synthetic data can be defined as,

$$Pr(\boldsymbol{y}|\boldsymbol{x}) = \frac{exp(E(\boldsymbol{y}, \boldsymbol{x}))}{\int_{-\infty}^{\infty} exp(E(\boldsymbol{y}, \boldsymbol{x}))d\boldsymbol{y}}. \quad (6)$$

where, $E$ is the energy function. In order to predict the depth of a new image we need to solve a maximum aposteriori (MAP) problem, $\widehat{\boldsymbol{y}} = \text{argmax}_y Pr(\boldsymbol{y}|\boldsymbol{x})$.

Let $\xi$ and $\eta$ be unary and pairwise potentials over nodes $\mathcal{N}$ and edges $\mathcal{S}$ of $\boldsymbol{x}$, then the energy function can be formulated as,

$$E(\boldsymbol{y}, \boldsymbol{x}) = \sum_{i \in \mathcal{N}} \xi(y_i, \boldsymbol{x}; \boldsymbol{\theta}) + \sum_{(i,j) \in \mathcal{S}} \eta(y_i, y_j, \boldsymbol{x}; \boldsymbol{\beta}), \quad (7)$$

where, $\xi$ regresses the depth from a single superpixel and $\eta$ encourages smoothness between neighboring superpixels. The objective is to learn the two potentials in a unified CNN framework. The unary part takes a single image superpixel patch as an input and feeds it to a CNN which outputs a regressed depth of that superpixel. Based on [14] the unary potential can be defined as,

$$\xi(y_i, \boldsymbol{x}; \boldsymbol{\theta}) = -(y_i - h_i(\boldsymbol{\theta}))^2 \quad (8)$$

where $h_i$ is the regressed depth of superpixel and $\theta$ represents CNN parameters.

The pairwise potential function is based on standard CRF vertex and edge feature functions studied extensively in [21] and other works. Let $\boldsymbol{\beta}$ be the network parameters and $\boldsymbol{S}$ be the similarity matrix where $S_{i,j}^k$ represents a similarity metric between i the $i^{th}$ and $j^{th}$ superpixel. Since inverse of intensity is a very valuable cue for depth estimation in endoscopy settings, we use intensity difference and greyscale histogram as pairwise similarities expressed in the general $\ell_2$ form. The pairwise potential can then be defined as,

$$\eta(y_i, y_j; \boldsymbol{\beta}) = -\frac{1}{2} \sum_{k=1}^{K} \beta_k S_{i,j}^k (y_i - y_j)^2. \quad (9)$$

The overall energy function can now be written as,

$$E = -\sum_{i \in \mathcal{N}} (y_i - h_i(\boldsymbol{\theta}))^2 - \frac{1}{2} \sum_{(i,j) \in \mathcal{S}} \sum_{k=1}^{K} \beta_k S_{i,j}^k (y_i - y_j)^2. \quad (10)$$

For training the negative log likelihood of the probability density function which can be calculated from Eq. 6 is minimized with respect to the two learning parameters. Two regularization terms are added to the objective function to penalize heavily weighted vectors $(\lambda_\theta, \lambda_\beta)$. Assuming $N$ is the number of images in the training data,

$$\min_{\theta, \beta \geq 0} -\sum_1^N \log Pr(\boldsymbol{y}|\boldsymbol{x}; \boldsymbol{\theta}, \boldsymbol{\beta}) + \frac{\lambda_\theta}{2} \|\theta\|_2^2 + \frac{\lambda_\beta}{2} \|\beta\|_2^2. \quad (11)$$

The optimization problem is solved using stochastic gradient decent-based back propagation.

# 6. Experiments

## 6.1. Evaluation Datasets

We use three kinds of datasets in our quantitative and qualitative study of the proposed methods. Since there are no publicly available endoscopy datasets with ground true depth, we generate two kinds of datasets for quantitative evaluation: a) images from a virtual endoscope in a colon phantom, and b) CT-registered optical endoscopy data collected from a real porcine colon (Fig. 4). We also use publicly-available human colonoscopy images to qualitatively assess if intuitive depth maps can be generated from real endoscopy videos.

**Colon Phantom Data:** The colon phantom data is generated from a CT-reconstructed model of a colon phantom molded from a real colon (Chamberlain Group Colonoscopy Trainer, SKU 2003[1]). A virtual endoscope is used to render images from a variety of endoscopy images with corresponding ground truth from the CT-reconstructed model. 2,160 images are generated via this procedure and are used for evaluation (Fig. 4).

**Real Porcine Colon Data:** Real endoscopy images were recorded from a pig colon fixed to a scaffold. A 3D model of the scaffold was then acquired with a CT measurement, and ground truth depth was generated for each real endoscopy image by registering virtual endoscopy views from the CT and optical endoscopy views from an endoscope (Fig. 4). 1,400 images with corresponding ground truth depth are generated using this procedure and are used for evaluation.

**Real Endoscopy Data:** We also evaluate our networks on publicly available endoscopy data[2] [1, 30]. However, these datasets do not have ground true depth and can only be used for qualitative evaluations.

## 6.2. Depth Estimation Network Trained on Synthetic Images

### Implementation Details

The architecture used for training an endoscopy depth estimation network includes training the unary and a pairwise parts of a CRF in a unified framework presented in [14]. The unary part is composed of a fully convolutional network which generates convolution maps that are fed into a superpixel pooling layer followed by three fully connected layers. The pairwise part operates on a superpixel level and is composed of a single fully connected layer. This setup was implemented using VLFeat Mat-ConvNet[3] using MATLAB 2017a and CUDA 8.0. The training data was prepared by over-segmenting each virtual endoscopy image into superpixels and corresponding ground truth depth were assigned

[1] https://www.thecgroup.com/product/colonoscopy-trainer-2003/
[2] https://polyp.grand-challenge.org/databases/
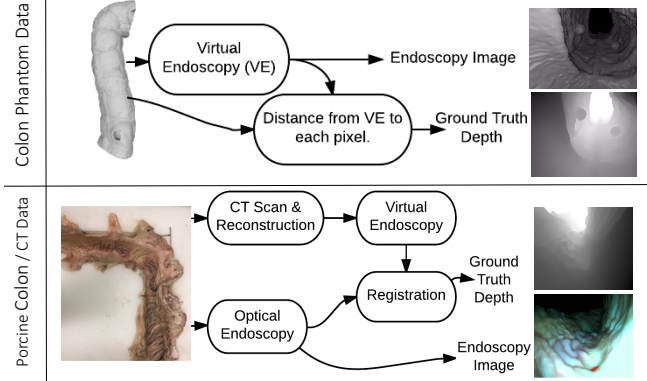[3] http://www.vlfeat.org/matconvnet/

Figure 4. The image collection and generation pipeline for colon phantom data and porcine colon data. The colon phantom data is collected from a 3D rendered colon phantom using a virtual endoscope. The ground truth depth is calculated using the 3D model. The porcine colon data is collected by imaging a porcine colon mounted on a scaffold using an optical endoscope and reconstructing a 3D model of the colon from CT measurements. The optical endoscopy and CT views are then registered to get the ground truth depth maps.

to each superpixel. Synthetic endoscopy data and its corresponding ground truth depth was generated according to the synthetic data generation pipeline presented in Section 3. The generated data was randomized to prevent the network from learning too many similar features quickly. 55% of the data was used for training and 40% for validation and 5% for testing. Training was done using K80 GPUs. Momentum was set at 0.9 as suggested in [14] and both weight decay parameters in Eq. 11 $(\lambda_\theta, \lambda_\beta)$ were set to 0.0007. The learning rate was initialized at 0.00001 and decrease by 20% every 20 epochs. These parameters were tuned to achieve best results. A total of 300 epochs were run and the epochs with least $\log 10$ error were selected to avoid the selection of an over-fitted model.

## 6.3. Adversarial Training for Reverse Domain Adaptation

### Implementation Details

Since the depth estimation network was trained solely on synthetic data all test images need to have a synthetic-like representation for the depth estimation to perform effectively. A transformer network was trained using the reverse domain adaption paradigm presented in Section 4.

The transformer and discriminator networks were implemented using tensorflow. The synthetic and real endoscopy images were down-sampled to a pixel size of $244 \times 244$ for computational efficiency. The real images were also converted to grayscale. The training between the transformer and the discriminator proceeds alternatively.

The transformer network was a standard residual network (ResNet) [9]. This is similar to [28], but for refining real data to be synthetic rather than the other way around.
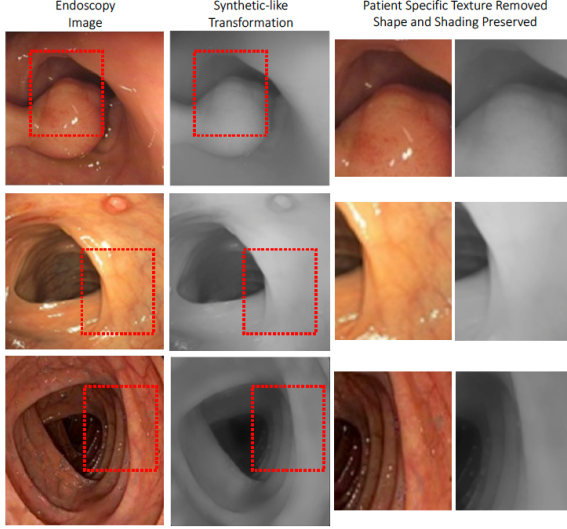
Figure 5. Examples of real endoscopy images transformed to their synthetic-like representations. Patient-specific texture is clearly removed during the transformation.
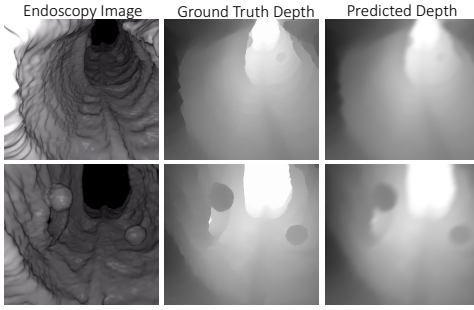


Figure 6. Examples of rendered images, corresponding ground truth depth, and depth estimates from a colon phantom.

| Test Dataset | NRMSE | HD | SSIM |
|---|---|---|---|
| Colon Phantom | 0.38 | 0.36 | 0.52 |
| Trans. Colon Phantom | **0.23** | **0.23** | **0.77** |
| Real Porcine Colon | 0.61 | 0.58 | 0.33 |
| Trans. Real Porcine Colon | **0.32** | **0.30** | **0.59** |

Table 1. A comparison between depth estimated from raw images and domain adapted images via our transformer network.

An input image of size $244 \times 244$ is convolved with a filter of $7 \times 7$ that outputs $64$ feature maps which are then passed to 10 ResNet blocks followed by a $1 \times 1$ convolution layer resulting in one feature map. The transformer is first trained with only the self-regularization term for the first 800 steps and the discriminator for 200 steps. The discriminator network is a standard classifier with five convolution layers, two max-pooling layers and softmax.

| Method | NRMSE | HD | SSIM |
|---|---|---|---|
| DiL (No Texture)[17] | 0.57 | 0.56 | 0.35 |
| DiL (Texture 1)[17] | 0.49 | 0.44 | 0.31 |
| DiL (Texture 2)[17] | 0.43 | 0.43 | 0.30 |
| DiL (Average) | 0.50 | 0.48 | 0.32 |
| Ours (No Texture) | 0.19 | 0.18 | 0.81 |
| Ours (Phantom) | 0.23 | 0.23 | 0.77 |
| Ours (Porcine) | 0.32 | 0.30 | 0.59 |
| Ours (Average) | **0.25** | **0.24** | **0.72** |

Table 2. Results of our method as compared to the state-of-the-art endoscopy depth estimation method.
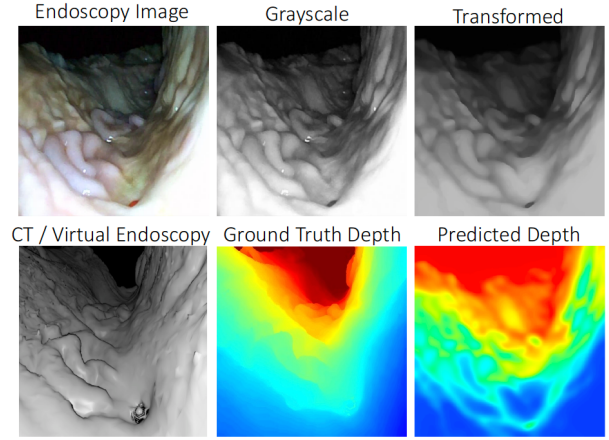


Figure 7. Depth estimates from porcine colon data. The optical endoscopy image is converted to grayscale and transformed to its synthetic-like representation using our transformer network. The optical endoscopy view is registered to it's corresponding CT view to obtain ground truth depth.

## 6.4. Results

**Transformer Network:** Fig. 5 shows examples of real endoscopy images transformed to their synthetic-like representations. It can clearly be seen that the patient specific information has been removed and clinically-relevant features have been preserved for depth estimation and polyp identification. A close-up of the images show that the vasculature has been removed while preserving the shape information. In the next subsections we demonstrate that the depth estimation network trained on synthetic data performs significantly better with images transformed to their synthetic like representations.

**Depth Evaluation Metrics:** We compare our depth estimates to corresponding ground truth values based on three metrics which have been used in previous endoscopy depth estimation work:

- Normalized root mean square error (NRMSE): NRMSE $= \frac{\sqrt{\frac{\sum_i (x_i - y_i)}{n}}}{(x_{max} - x_{min})}$, is a normalized RMS error for comparative analysis across datasets. A lower
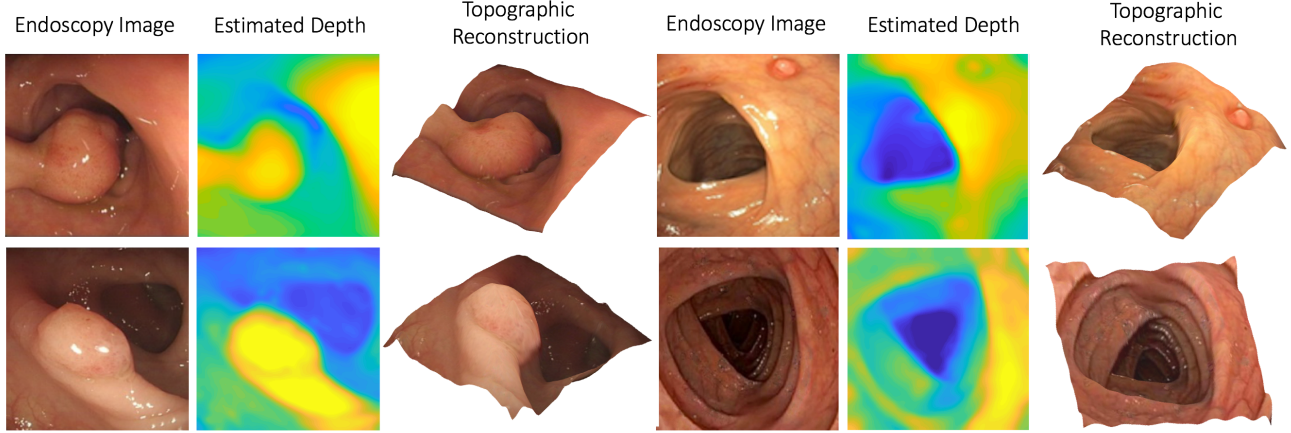
Figure 8. Depth estimates and topographical reconstructions from monocular endoscopy images. Each endoscopy image is transformed to its synthetic-like representation as shown in Fig. 5 and is fed into our depth estimation network. The depth is then used to reconstruct the surface topography.

RMSE indicates the data being compared is more similar.

- Hausdorff distance (HD): HD calculates the greatest of all the distances from a point in the ground truth data to the closest point in the calculated data [10]. It can be calculated as $H(x,y) = \max(\boldsymbol{h}(x,y), \boldsymbol{h}(y,x))$ where $\boldsymbol{h} = \max_a \min_b \|a - b\|$. A lower HD indicates the two datasets being compared are more similar.

- Structural Similarly Index (SSIM): The SSIM is an image assessment index calculated on the basis of luminance, contrast and structure. The SSIM in this paper is calculated according to the definition proposed in [32]. This index is between $-1$ and $1$ with $1$ indicating identical images.

**Quantitative Results**

Table 1 compares depth estimation results from colon phantom and real porcine colon data with and without domain transformation. It can clearly be seen that depth estimation is improved by domain transformation. As expected, the improvement in depth estimation that domain transformation provides is marginal in the colon phantom data, which has homogenous material properties, and more significant in real porcine tissue, which has natural biological variation in mucosal texture. There is a **88% improvement** in the SSIM for the porcine colon data and a **48% improvement** for the colon phantom data by transforming the input data using our proposed paradigm. Fig. 6 and 7 show representative depth estimation results for the colon phantom and porcine colon, respectively.

**Comparative Analysis**

Due to the lack of available ground truth endoscopy depth data there is currently only one learning based monocular colonoscopy depth estimation study by Nadeem *et al.*

[17]. They implement dictionary learning (DiL) and use CT colonoscopy data for training. However, unlike our data, their data does not follow optically-correct inverse square intensity fall off, which we expect to be a significant cue for absolute depth. Table 2 shows a comparative analysis of their results compared to those from our approach. We demonstrate that our depth estimation is significantly better than their method and **improves the SSIM by 125%**.

**Qualitative Results**

For the purposes of demonstration, we also show that it is possible to estimate depth from real human endoscopy data. Fig. 8 shows monocular endoscopy images, their estimated depth, and corresponding topographic reconstructions. Topographical reconstructions are reconstructed by overlaying depth on a 3D manifold. These depth estimates are qualitative and there is no corresponding ground truth depth available.

## 7. Conclusions and Future Work

In this paper, we propose a novel reverse domain adaptation method that transforms real medical images into useful synthetic representations while preserving clinically relevant features. We validated this method in the task of monocular depth estimation for endoscopy images, in which we first learned depth from a large synthetic dataset, and then demonstrated an 88% and 48% improvement in predicting depth from synthetic-like domain-adapted images over raw images for both a real porcine colon and a colon phantom respectively. Future work will focus on using the proposed reverse domain adaption paradigm for other medical imaging modalities and on using the predicted depth for improving automated polyp segmentation and classification.

# References

[1] J. Bernal, F. J. Sánchez, G. Fernández-Esparrach, D. Gil, C. Rodríguez, and F. Vilariño. Wm-dova maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from physicians. *Computerized Medical Imaging and Graphics*, 43:99–111, 2015. 6

[2] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. *arXiv preprint arXiv:1612.05424*, 2016. 2

[3] P. Costa, A. Galdran, M. I. Meyer, M. Niemeijer, M. Abràmoff, A. M. Mendonça, and A. Campilho. End-to-end adversarial retinal image synthesis. *IEEE Transactions on Medical Imaging*, 2017. 3

[4] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015. 2

[5] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. 2, 3, 4

[6] H. Greenspan, B. van Ginneken, and R. M. Summers. Guest editorial deep learning in medical imaging: Overview and future promise of an exciting new technique. *IEEE Transactions on Medical Imaging*, 35(5):1153–1159, 2016. 1

[7] J. T. Guibas, T. S. Virdi, and P. S. Li. Synthetic medical images from dual generative adversarial networks. *arXiv preprint arXiv:1709.01872*, 2017. 2, 3

[8] A. Gupta, A. Vedaldi, and A. Zisserman. Synthetic data for text localisation in natural images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2315–2324, 2016. 3

[9] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 6

[10] D. P. Huttenlocher, G. A. Klanderman, and W. J. Rucklidge. Comparing images using the hausdorff distance. *IEEE Transactions on pattern analysis and machine intelligence*, 15(9):850–863, 1993. 7

[11] S. Kohl, D. Bonekamp, H.-P. Schlemmer, K. Yaqubi, M. Hohenfellner, B. Hadaschik, J.-P. Radtke, and K. Maier-Hein. Adversarial networks for the detection of aggressive prostate cancer. *arXiv preprint arXiv:1702.08014*, 2017. 3

[12] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. 2

[13] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 2

[14] F. Liu, C. Shen, G. Lin, and I. Reid. Learning depth from single monocular images using deep convolutional neural fields. *IEEE transactions on pattern analysis and machine intelligence*, 38(10):2024–2039, 2016. 5, 6

[15] M. Mardani, E. Gong, J. Y. Cheng, S. Vasanawala, G. Zaharchuk, M. Alley, N. Thakur, S. Han, W. Dally, J. M. Pauly, et al. Deep generative adversarial networks for compressed sensing automates mri. *arXiv preprint arXiv:1706.00051*, 2017. 3

[16] P. Moeskops, M. Veta, M. W. Lafarge, K. A. Eppenhof, and J. P. Pluim. Adversarial training and dilated convolutions for brain mri segmentation. In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*, pages 56–64. Springer, 2017. 3

[17] S. Nadeem and A. Kaufman. Computer-aided detection of polyps in optical colonoscopy images. In *SPIE Medical Imaging*, pages 978525–978525. International Society for Optics and Photonics, 2016. 7, 8

[18] A. Osokin, A. Chessel, R. E. C. Salas, and F. Vaggi. GANs for biological image synthesis. 2017. 3

[19] C. Payer, D. Štern, H. Bischof, and M. Urschler. Regressing heatmaps for multiple landmark localization using cnns. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 230–238. Springer, 2016. 3

[20] B. Planche, Z. Wu, K. Ma, S. Sun, S. Kluckner, T. Chen, A. Hutter, S. Zakharov, H. Kosch, and J. Ernst. Depthsynth: Real-time realistic synthetic data generation from cad models for 2.5 d recognition. *arXiv preprint arXiv:1702.08558*, 2017. 3

[21] T. Qin, T.-Y. Liu, X.-D. Zhang, D.-S. Wang, and H. Li. Global ranking using continuous conditional random fields. In *Advances in neural information processing systems*, pages 1281–1288, 2009. 5

[22] W. Qiu and A. Yuille. Unrealcv: Connecting computer vision to unreal engine. In *Computer Vision–ECCV 2016 Workshops*, pages 909–916. Springer, 2016. 2

[23] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 234–241. Springer, 2015. 2

[24] A. A. Rusu, M. Vecerik, T. Rothörl, N. Heess, R. Pascanu, and R. Hadsell. Sim-to-real robot learning from pixels with progressive nets. *arXiv preprint arXiv:1610.04286*, 2016. 2

[25] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training gans. In *Advances in Neural Information Processing Systems*, pages 2234–2242, 2016. 3, 4

[26] A. Shafaei, J. J. Little, and M. Schmidt. Play and learn: Using video games to train computer vision models. *arXiv preprint arXiv:1608.01745*, 2016. 2

[27] D. Shen, G. Wu, and H.-I. Suk. Deep learning in medical image analysis. *Annual Review of Biomedical Engineering*, (0), 2017. 1

[28] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb. Learning from simulated and unsupervised images through adversarial training. *arXiv preprint arXiv:1612.07828*, 2016. 2, 3, 4, 6

[29] H. Su, C. R. Qi, Y. Li, and L. J. Guibas. Render for cnn: Viewpoint estimation in images using cnns trained with ren-

dered 3d model views. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2686–2694, 2015. 3

[30] N. Tajbakhsh, S. R. Gurudu, and J. Liang. Automated polyp detection in colonoscopy videos using shape and context information. *IEEE transactions on medical imaging*, 35(2):630–644, 2016. 6

[31] G. Varol, J. Romero, X. Martin, N. Mahmood, M. Black, I. Laptev, and C. Schmid. Learning from synthetic humans. *arXiv preprint arXiv:1701.01370*, 2017. 3

[32] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 8

[33] J. M. Wolterink, T. Leiner, M. A. Viergever, and I. Isgum. Generative adversarial networks for noise reduction in low-dose ct. *IEEE Transactions on Medical Imaging*, 2017. 3

[34] L. Zhang, A. Gooya, and A. F. Frangi. Semi-supervised assessment of incomplete lv coverage in cardiac mri using generative adversarial nets. In *International Workshop on Simulation and Synthesis in Medical Imaging*, pages 61–68. Springer, 2017. 3

[35] Y. Zhang, L. Yang, J. Chen, M. Fredericksen, D. P. Hughes, and D. Z. Chen. Deep adversarial networks for biomedical image segmentation utilizing unannotated images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 408–416. Springer, 2017. 3

[36] J.-Y. Zhu, P. Krähenbühl, E. Shechtman, and A. A. Efros. Generative visual manipulation on the natural image manifold. In *European Conference on Computer Vision*, pages 597–613. Springer, 2016. 3