

# Generating large labeled data sets for laparoscopic image processing tasks using unpaired image-to-image translation

Micha Pfeiffer<sup>1</sup>, Isabel Funke<sup>1</sup>, Maria R. Robu<sup>2,3</sup>, Sebastian Bodenstedt<sup>1</sup>, Leon Strenger<sup>1</sup>, Sandy Engelhardt<sup>4</sup>, Tobias Roß<sup>5</sup>, Matthew J. Clarkson<sup>2,3</sup>, Kurinchi Gurusamy<sup>6</sup>, Brian R. Davidson<sup>6</sup>, Lena Maier-Hein<sup>5</sup>, Carina Riediger<sup>7</sup>, Thilo Welsch<sup>7</sup>, Jürgen Weitz<sup>7</sup>, and Stefanie Speidel<sup>1</sup>

<sup>1</sup> Translational Surgical Oncology, National Center for Tumor Diseases, Dresden, Germany

`micha.pfeiffer@nct-dresden.de`

<sup>2</sup> Wellcome/EPSRC Centre for Interventional & Surgical Sciences, University College London, UK

<sup>3</sup> Centre for Medical Image Computing, University College London, UK

<sup>4</sup> Faculty of Computer Science, Mannheim University of Applied Sciences, Germany

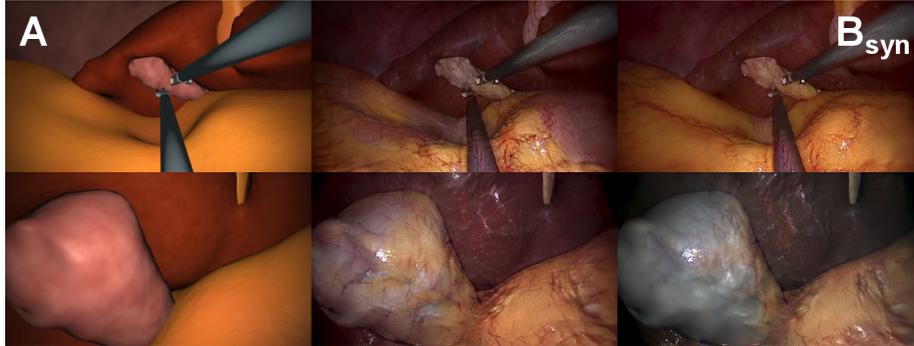
<sup>5</sup> German Cancer Research Center, Heidelberg, Germany

<sup>6</sup> Division of Surgery and Interventional Science, University College London, UK

<sup>7</sup> Department for Visceral, Thoracic and Vascular Surgery, University Hospital Dresden, Germany

**Abstract.** In the medical domain, the lack of large training data sets and benchmarks is often a limiting factor for training deep neural networks. In contrast to expensive manual labeling, computer simulations can generate large and fully labeled data sets with a minimum of manual effort. However, models that are trained on simulated data usually do not translate well to real scenarios. To bridge the domain gap between simulated and real laparoscopic images, we exploit recent advances in unpaired image-to-image translation. We extend an image-to-image translation method to generate a diverse multitude of realistically looking synthetic images based on images from a simple laparoscopy simulation. By incorporating means to ensure that the image content is preserved during the translation process, we ensure that the labels given for the simulated images remain valid for their realistically looking translations. This way, we are able to generate a large, fully labeled synthetic data set of laparoscopic images with realistic appearance. We show that this data set can be used to train models for the task of liver segmentation of laparoscopic images. We achieve average dice scores of up to 0.89 in some patients without manually labeling a single laparoscopic image and show that using our synthetic data to pre-train models can greatly improve their performance. The synthetic data set will be made publicly available, fully labeled with segmentation maps, depth maps, normal maps, and positions of tools and camera (<http://opencas.dkfz.de/image2image>).

**Keywords:** Unsupervised · Image Translation · Segmentation · Laparoscopy



**Fig. 1.** Images from simple laparoscopic computer simulation (domain  $A$ , first column) translated to look like real laparoscopic video frames (synthetic  $B_{syn}$ , second and third column) using various styles. During the unpaired training process, a multi-scale structural similarity loss ensures that structures remain similar. This enables us to use the generated images along with labels from domain  $A$  as training data for various tasks.

## 1 Introduction

With the increase in computing power, there is an obvious trend towards training larger and deeper networks. However, in the medical domain, the lack of large data sets is a strong limiting factor [10]. The difficulty of recording real patient data in an operating room, legal restrictions on sharing and the great expense of manual labeling by experts make it near impossible to generate large training benchmarks. This work focuses on the example of the segmentation of laparoscopic videos, where deep networks can achieve high accuracies, but sometimes fail to generalize to new patients due to the lack of more labeled data [4]. A solution to this problem could be the usage of synthetic training data. In computer simulations, large amounts of fully labeled data can be created automatically. The main issue here is that models trained on synthetic data usually do not generalize well to real data, due to the *domain gap* between the two.

Instead, we propose to use *image-to-image translation* techniques to translate images from the domain of simulated images in which labels are known (domain  $A$ ), to the domain of real images in which we want to train our model (domain  $B$ ). Recent advances in image translation make it possible to do this even if the data is unpaired, i.e. no direct mapping between samples in one domain to samples in the other domain exists [13]. Additionally, *multi-modal* image-to-image translation [6,8] enables us to control the style of the translation result, which can be utilized to increase the diversity in the final data set. In the present work, domain  $A$  consists of images from very simple laparoscopic 3D computer simulations while domain  $B$  is the domain of images from real laparoscopy video feeds. In order to use the translated data for training, care must be taken that a) the translated images look realistic enough to bridge the domain gap and b) the labels remain valid. This is especially difficult in laparoscopic images, since

image content can change drastically between different viewpoints and between patients. To achieve our goal, we build up on several methods:

**Unpaired translation** The CycleGAN [13] has made it possible to translate images between two unpaired domains by usage of a cycle consistency loss and adversarial losses. A generator network  $G_B$  translates images from  $A$  to  $B$  which a discriminator network  $D_B$  tries to differentiate from real images in  $B$ . At the same time, generator  $G_A$  and  $D_A$  use the same method to translate images from  $B$  to  $A$ . The cycle consistency states that an image  $a$  translated to  $B$  and back to  $A$  must match the original image, i.e.  $a = G_A(G_B(a))$  (and symmetrically for an image  $b$ ). This method can only learn a one-to-one mapping (uni-modal), meaning each input image will generate exactly one output.

**Multi-Modal translation** The key idea behind multi-modal image translation is the separation of an image’s *content* from its *style*. The assumption is that the content between domains remains the same, while the style is domain-specific (texture, lighting). An encoder  $E_A$  first extracts a *style-code*  $s_a$  and a *content-code*  $c_a$  from the source image and a generator  $G_B$  then uses this content-code together with a style-code  $s_b$  from the target domain to create the image  $b'$  in the target domain [6,8]. The opposite direction works analogously. A cycle loss and various reconstruction losses bind the networks together.

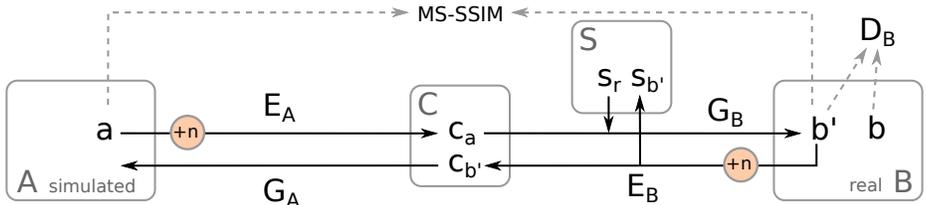
**Label-preserving translation** SPIGAN [9] proposes to train an additional network which tries to predict the depth map from the translated image, arguing that this preserves image structure. In our experience, this bears the risk of co-adaptation between the networks. AugGAN [5] and GANTruth [1] bind the generators to the image structure via weight-sharing with segmentation networks. However, AugGAN requires segmentation labels to be known for both domains and GANTruth requires a pre-trained segmentation network in the target domain. Our goal is to not use labels during the translation process, simplifying the training procedure.

**Contribution** In this work, we show how both the goal of realism as well as the preservation of label accuracy during translation can be achieved. First, we build an extension to the MUNIT framework which is asymmetrical and does not require the simulated domain to have multiple styles, speeding up the process of creating the simulated data. Next, we incorporate an additional *multi-scale structural similarity loss* [12] and show that it helps to preserve image content and structure despite large changes in camera viewpoint. Additionally, we show how the addition of noise in the encoders can help avoid *mode collapse* - where multiple images map to a similar output - and steganography. To validate the approach, we show that pre-training a segmentation model on the synthetic data can increase segmentation scores.

As part of this work, we translate 100 000 images to domain  $B$  (see Fig. 1). This data set, fully labeled with segmentation maps, depth maps and further labels as well as the code will be publicly available<sup>8</sup>, with possible applications ranging from pre-training to benchmarking.

## 2 Methods

Unpaired multi-modal image-to-image translations can output convincing results, but have mostly been tested on scenarios where the content stays similar in all images across both domains (such as faces to faces or mountains to mountains) [6,8]. In laparoscopy, viewpoints can change and structures - such as the gallbladder or abdominal wall - move into and out of the view. Incorporating this into our data set is necessary as we want it to be very diverse, however, the mismatch in domain distributions can lead to many wrongly added details, such as a gallbladder where there should only be liver and fat tissue replacing liver tissue. The following describes our extensions to the MUNIT architecture which enable us to deal with these issues, namely adding a structure-preserving loss, simplifying the encoder  $E_A$  and using noise to avoid co-adaptation of the networks. The resulting training process is outlined in Fig. 2.



**Fig. 2.** Architecture based on MUNIT [6]. Image  $a$  randomly drawn from  $A$  is translated to  $B$  and back to  $A$ , where a cycle loss ensures that  $a$  is reconstructed correctly. The same is done in the opposite direction for images drawn from  $B$ . Various reconstruction losses ensure that the generators and encoders work as expected (please see [6] for more details). During the translation process, images from  $A$  are encoded to a latent code  $c_a$ , while images from  $B$  are split into two latent codes: content  $c_b$  and style  $s_b$ . Unlike MUNIT, we do not have a style in  $A$ , which simplifies the creation of the rendered images. Furthermore, we add noise to all encoders to prevent the hiding of information and add the MS-SIM loss between source images and their translations.

### 2.1 Architecture

*Multi-Scale Structural Similarity (MS-SSIM) loss:* Unpaired translation networks often invent details in their output. This is likely due to two reasons:

<sup>8</sup> Data set and code available at: <http://opencas.dkfz.de/image2image/>

1) Some structures and some viewpoints occur more in one of the two domains than in the other. For example, domain  $A$  contains more close-ups of the liver due to the random placement of the camera. The discriminator  $D_B$  will discourage these images, resulting in the generator  $G_B$  inventing structures like an additional gallbladder. 2) Generative models are susceptible to mode collapse. We add a multi-scale structural similarity [12] loss between an image  $a$  and its translation  $G_B(a)$  (and similarly in the other direction). The loss works on the image brightness (average over the channels) which ensures that brighter regions (such as the gallbladder) remain brighter and darker regions remain dark while at the same time not penalizing style-dependent changes in hue.

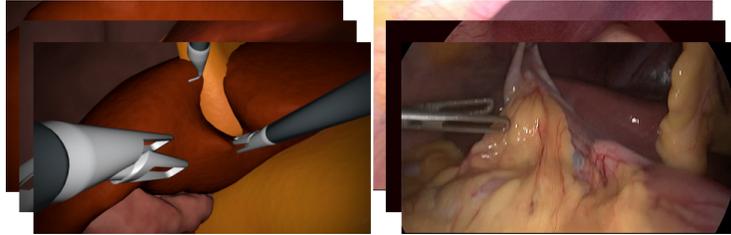
*Noise against steganography:* GANs have shown to be very effective hiding information in their output images [2]. Since the generators  $G_A$  and  $G_B$  are trained jointly to fulfill the cycle consistency,  $G_A$  learns to hide details of the image  $b$  in its translation which are useful for  $G_B$ . This is problematic when giving  $G_B$  a real image to translate, since these details are not present in this case. To circumvent this effect, we add Gaussian noise to the input of each translation network.

*Asymmetrical style:* One of our aims is to reduce the amount of manual work required to generate data. In this spirit, we want to translate from a simple and easy to set up domain  $A$  to a very complex domain  $B$  and let the computer do the bulk of the work automatically. We remove the part of encoder  $E_A$  which extracts the style and the style-injection from  $G_A$ . As a result, our setup becomes asymmetrical and we do not need to worry about creating multiple textures or lighting styles in the simulated domain  $A$ , simplifying the simulation process. During training, both the style extracted by  $E_B$  as well as randomly drawn style vectors are used when translating from  $A$  to  $B$ . In this way, the network can later translate images either using a random style or the style taken from a real image.

## 2.2 Translation data

To train our translation networks, we use two unpaired data sets, which both contain images with livers, gallbladders, tools, fat and abdominal wall (see Fig. 3).

*Rendered data set - Domain A:* We create six synthetic laparoscopic 3D-scenes using the liver and gallbladder surface meshes extracted from CT scans of six patients (3D-IRCAdB 01 data set, IRCAD, France). We add meshes which represent fat tissue, ligament and the inflated abdominal wall. Each tissue type is assigned a distinctive texture with small random details. We randomly place the camera together with a light source (representing the laparoscope) and tools. In this way, we render 2000 images from random perspectives for each patient, resulting in 12 000 synthetic images. To increase the diversity in our translated results, we repeat the process for four additional patients where no gallbladder is present, resulting in scenes similar to liver staging procedures. The images from all ten patients together make up our extended rendered data set  $A^+$ .



**Fig. 3.** Sample images from the two domains. Both contain similar objects, but no pairing information is known, and the distribution of content does not necessarily match.

*Real data set - Domain B:* The real images are taken from 80 videos of the Cholec80 data set (videos of 80 laparoscopic cholecystectomies) [11]. We first identify parts of the videos in which the gallbladder is still intact and then extract frames at five frames per second. We separate the resulting images into a training data set  $B_{tr}$  (75 patients, roughly 74 000 images) and a segmentation data set  $B_v$  (5 patients). We manually segment the liver in 196 images of  $B_v$  (at a rate of one frame every five seconds).

### 2.3 Experiments

We train the translation networks for 375 000 iterations. Afterwards, we translate all images from  $A^+$ , using five randomly drawn style vectors for each image, resulting in 100 000 images which we call the synthetic data set  $B_{syn}$ .

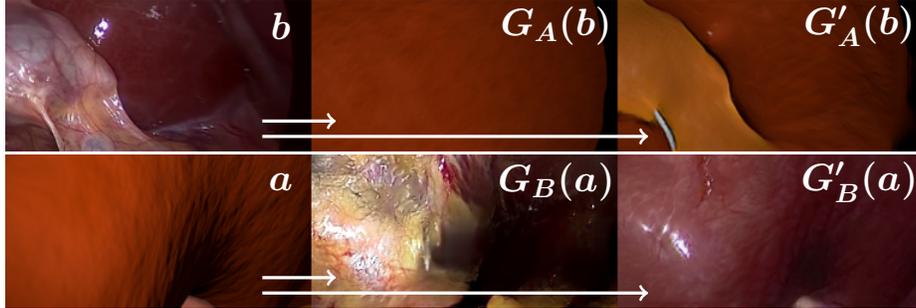
Evaluating the image quality quantitatively is difficult. Instead, we validate the usefulness of the synthetic data set by using it as training data for a segmentation task: As a baseline, we first train a TerausNet-11 [7] on the real Cholec80 validation data set  $B_v$  in a leave-one-patient-out cross-validation (five models trained, each time one patient is left out of the training data to be used for testing). We then train the same network only on the synthetic data  $B_{syn}$  and validate it on all five patients in  $B_v$ . Furthermore, we test how the performance changes if the network which is already trained on  $B_{syn}$  is fine-tuned on the real data in the same cross-validation as before. The experiments are repeated for a TerausNet which has previously been pre-trained on the ImageNet data set [3].

To see how our synthetic data helps in the adaptation to a wider diversity of images, we evaluate the pre-initialized TerausNets on images from 13 liver staging sequences, in which a total of roughly 2000 images are segmented [4].

## 3 Results

Using the MS-SSIM loss can greatly improve the preservation of image structure, as shown qualitatively in Fig. 4 and helps in the correct usage of textures: The

correct assignment of texture to the various organs can be clearly seen and close-up shots of the liver surface result in highly detailed liver texture translations (more translation results in the supplementary materials).



**Fig. 4.** Qualitative results for the MS-SSIM loss. During translation of images  $b$  and  $a$ , the networks tend to remove ( $G_A(b)$ ) or add ( $G_B(a)$ ) detail. In contrast, networks  $G'_A$  and  $G'_B$ , which are trained with an MS-SSIM loss, preserve structures in both directions.

**Table 1.** Median dice scores for  $B_v$  (Patients 75 to 80 from Cholec80 data) and for the 13 staging procedures. In all cases where  $B_v$  is part of the training data, the reported results are from a leave-one-patient-out cross-validation (except for the staging procedures, where all five patients were used). In patient P78 most of the visible liver region is covered in ligament and fat tissue. Median scores for the 13 staging procedures increase considerably by using the synthetic data  $B_{syn}$  for pre-training. An additional improvement is achieved by pre-training on the ImageNet data  $I$ .

Training data	P76	P77	P78	P79	P80	Staging Procedures
$B_v$	0.50	0.68	0.42	0.52	0.56	
$B_{syn}$	0.73	0.70	0.13	0.74	0.76	
$B_{syn} + B_v$	0.74	0.72	0.40	0.64	0.61	
$I + B_v$	0.80	0.81	0.48	0.86	0.83	0.25
$I + B_{syn}$	0.89	0.80	0.12	0.80	0.85	0.61
$I + B_{syn} + B_v$	0.92	0.83	0.64	0.89	0.91	0.77

Training on our synthetic data shows considerable improvements over training only on the real data (Table 1). When using the synthetic data for pre-training, the median dice score improved by an average of 16 percent (no ImageNet pre-training) and 11 percent (with ImageNet pre-training).

When the network was tested on the 13 staging procedures [4] containing data that had not been seen at all during training, the mean dice score using

only real training data  $B_v$  was 0.25, and improved to 0.77 when the network was pre-trained with the synthetic data  $B_{syn}$ .

## 4 Discussion

In this work, we have shown that consistent translation results can be achieved despite having a large change in content and viewpoints.

The translated results alone can be used to achieve reasonably good scores on a segmentation task without labeling a single image. When pre-training a network with our synthetic data, we can demonstrate an increase in performance, compared with only using real data. We also show that the training data can help a network in generalizing to new situations.

Unpaired image-to-image translation is proving to be a very powerful tool in the generation of training data. Since the domain of surgical data science still mostly lacks large benchmarks and open data sets, it could greatly benefit from further development in this field.

## References

1. Bujwid, S., Martí, M., Azizpour, H., Pieropan, A.: Gantruth - an unpaired image-to-image translation method for driving scenarios (11 2018)
2. Chu, C., Zhmoginov, A., Sandler, M.: CycleGAN, a master of steganography (2017)
3. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: ImageNet: A Large-Scale Hierarchical Image Database. In: CVPR09 (2009)
4. Gibson, E., Robu, M.R., Thompson, S., Edwards, P.E., Schneider, C., Gurusamy, K., Davidson, B., Hawkes, D.J., Barratt, D.C., Clarkson, M.J.: Deep residual networks for automatic segmentation of laparoscopic videos of the liver (2017)
5. Huang, S.W., Lin, C.T., Chen, S.P., Wu, Y.Y., Hsu, P.H., Lai, S.H.: Auggan: Cross domain adaptation with gan-based data augmentation. In: Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y. (eds.) Computer Vision – ECCV 2018. pp. 731–744. Springer International Publishing, Cham (2018)
6. Huang, X., Liu, M.Y., Belongie, S., Kautz, J.: Multimodal unsupervised image-to-image translation. In: The European Conference on Computer Vision (ECCV) (September 2018)
7. Iglovikov, V.I., Shvets, A.A.: Terausnet: U-net with vgg11 encoder pre-trained on imagenet for image segmentation. CoRR **abs/1801.05746** (2018)
8. Lee, H.Y., Tseng, H.Y., Huang, J.B., Singh, M., Yang, M.H.: Diverse image-to-image translation via disentangled representations. In: The European Conference on Computer Vision (ECCV) (September 2018)
9. Lee, K.H., Ros, G., Li, J., Gaidon, A.: SPIGAN: Privileged adversarial learning from simulation. In: International Conference on Learning Representations (2019)
10. Maier-Hein, L., Vedula, S.S., Speidel, S., Navab, N., Kikinis, R., Park, A., Eisenmann, M., Feussner, H., Forestier, G., Giannarou, S., et al.: Surgical data science for next-generation interventions. Nature Biomedical Engineering **1**(9), 691 (2017)
11. Twinanda, A., Shehata, S., Mutter, D., Marescaux, J., De Mathelin, M., Padoy, N.: Endonet: A deep architecture for recognition tasks on laparoscopic videos. IEEE Transactions on Medical Imaging **36** (02 2016)

12. Wang, Z., Simoncelli, E.P., Bovik, A.C.: Multiscale structural similarity for image quality assessment. In: The Thrity-Seventh Asilomar Conference on Signals, Systems Computers, 2003. vol. 2, pp. 1398–1402 Vol.2 (Nov 2003)
13. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: Computer Vision (ICCV), 2017 IEEE International Conference on (2017)