# Domain Adaptive Sim-to-Real Segmentation of Oropharyngeal Organs

Guankun Wang[1], Tian-Ao Ren[2,3], Jiewen Lai[1], Long Bai[1] and Hongliang Ren[1*]

[1]Department of Electronic Engineering, The Chinese University of Hong Kong, Hong Kong, China.
[2]College of Mechanical and Electrical Engineering, Beijing University of Chemical Technology, Beijing, China.
[3]Shenzhen Research Institute, The Chinese University of Hong Kong, Shenzhen, Guangdong, China.

*Corresponding author(s). E-mail(s): hlren@ee.cuhk.edu.hk; ren@nus.edu.sg;
Contributing authors: gkwang@link.cuhk.edu.hk; taren@buct.edu.cn; jiewen.lai@cuhk.edu.hk; b.long@link.cuhk.edu.hk;

**Abstract**

Video-assisted transoral tracheal intubation (TI) necessitates using an endoscope that helps the physician insert a tracheal tube into the glottis instead of the esophagus. The growing trend of robotic-assisted TI would require a medical robot to distinguish anatomical features like an experienced physician which can be imitated by utilizing supervised deep-learning techniques. However, the real datasets of oropharyngeal organs are often inaccessible due to limited open-source data and patient privacy. In this work, we propose a domain adaptive Sim-to-Real framework called **I**oU-**R**anking **B**lend-**A**rt**F**low (IRB-AF) for image segmentation of oropharyngeal organs. The framework includes an image blending strategy called IoU-Ranking Blend (IRB) and style-transfer method ArtFlow. Here, IRB alleviates the problem of poor segmentation performance caused by significant datasets domain differences; while Art-Flow is introduced to reduce the discrepancies between datasets further. A virtual oropharynx image dataset generated by the SOFA framework is used as the learning subject for semantic segmentation to deal

with the limited availability of actual endoscopic images. We adapted IRB-AF with the state-of-the-art domain adaptive segmentation models. The results demonstrate the superior performance of our approach in further improving the segmentation accuracy and training stability.

**Keywords:** Domain adaption, Semantic segmentation, Sim-to-Real Transfer

# 1 Introduction

Transoral tracheal intubation (TI) is the gold standard for securing a patient's airway when they require respiratory assistance. However, the success of this procedure hinges on the physician's skills to correctly insert an endotracheal tube into the patient's trachea [1]. Asphyxia, hypoxia, and pulmonary aspiration can cause severe morbidity and death if the TI is not performed in a timely manner [2]. With advances in robotics and AI technology, image-guided automation that uses visual information to plan, complete, and recognize specific tasks is becoming an emerging area in medical robotics [3–5] and rehabilitation [6]. On top of video-assisted transoral TI, robot-assisted transoral TI makes intubation even more effective through automation. Yet, robotization requires the robot to distinguish and understand the contour of the anatomical features like an experienced physician. Therefore, it is expected that the robot should segment the endoscopic vision with sufficient fidelity. The above-mentioned initiative motivates the work herewith.

Segmenting oropharyngeal organs is one of the most critical steps in robot-assisted intubation. In practice, semantic segmentation divides each pixel into a label, and each label is assigned to a class, which can be applied to radiation, image-guided therapies, and enhanced radiological diagnostics [7]. However, obtaining real datasets of oropharyngeal organs is challenging due to patient privacy and the need for sufficient open-source data. Deploying supervised learning techniques for segmentation tasks with insufficient data would result in poor performance of the segmentation network model. To alleviate this problem, one can exploit synthetic or simulation data for Sim-to-Real transfer [8]. The method generally trains the network using rich simulation data while the tests are performed in real data that lacks abundance. Current simulation frameworks in the medical field, such as Blender [9], SOFA [10], and so on, provide users with increasingly realistic 3D scenarios with more sensations like textures and interaction-responsive physical properties. All anatomical structures can be rotated, scaled, and moved at any angle, and the hierarchy and adjacency of organs can be viewed from the inside out. They not only aid in providing a potentially abundant data source [11] but also relieve privacy concerns with real datasets.

Nonetheless, there are some challenges in employing virtual data for Sim-to-Real in deep learning. For instance, the open-source project Surgical Blender
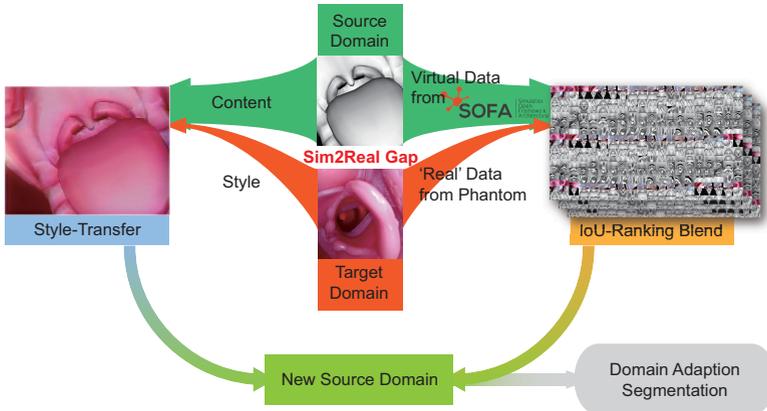
**Figure 1** Overview of the proposed IRB-AF framework. The style-transfer module first narrows the difference (i.e., the sim-to-real gap) between two domains. Then the source domain is reconstructed based on IoU-Ranking Blend. The new source domain will be used as the training set for segmentation.

can build virtual organ textures and morphologies based on advanced computer graphic techniques. Still, the pre-processing of this framework can be dramatically complex and time-consuming, which leads to limited accessibility of the data [12, 13]. While Surgical Blender is designed explicitly for virtual surgical simulation, SOFA is a relatively mature and high-performance library for general-purpose physical simulations. Although SOFA can easily handle general 3D models with a broader range of accessibility, their textures and morphologies are usually missing. In this work, we choose SOFA-generated images as the source domain, since data generation is more accessible with higher efficiency. For the problem of lower-quality data, if our work can improve the segment performance of models in Sim-to-Real, a better demonstration can be expected for domain adaption of other general data. The model trained with simulation images may have difficulty in segmentation compared with the model trained with real images, as the oropharyngeal organs in simulation images are constrained in representing natural textures like color and reflection. Deep learning models that are sensitive to data have yet to guarantee an effective generalization to unrestricted test cases. Therefore, segmentation performance will be greatly degraded in real datasets (i.e., target domains) if the model training relies on simulated datasets only, making the Sim-to-Real deployment challenging.

Given the abovementioned problems, reducing the differences between datasets is the most direct and effective way. Among all, domain randomization and unsupervised domain adaptation are often used. Domain randomization, as a commonly used technique in Sim-to-Real transfer learning, is often used to increase the robustness of a model to real-world variability [14]. It is generally done by randomizing various aspects of the simulation environment, such as the texture, lighting, and even physical property of objects, followed by

model training on the randomized simulation data to adapt to different environment variations [15]. The training process shall be repeated until the model is sufficiently robust after being evaluated by real-world data.

Unsupervised Domain Adaptation (UDA) refers to the process of adapting a model that has been trained on annotated samples from one distribution (the source domain) to function on another distribution (the target domain) for which no annotations are provided [16]. Most recent UDA works focus on lessening the difference between the two datasets during training on the source domain. Maximum Mean Disparity (MMD) and its kernel variants are proposed by [17, 18] as a common measure of discrepancy. Central Moment Discrepancy (CMD) [19] extended the measure to higher-order statistics. In addition, some researchers used pseudo labels provided by self-training [20], or target data [21, 22] produced by generative networks. Traditional domain adaptation methods [16, 23, 24] demonstrate excellent domain adaptation performances in the datasets with minor differences like GTA5 [25], SYN-THIA [26], and CityScapes [27]. However, the model's performance will decline rapidly when the differences become more significant.

Therefore, semi-supervised learning (SSL) techniques are required. During SSL, the model is adapted to the dataset with one of its subsets annotated. By properly aligning domains, SSL can become unsupervised domain adaptation. There are some common methods between SSL and UDA, such as self-training [20, 28], class balancing [29], and generative model. As one of the successful approaches, entropy minimization is also used in semi-supervised learning [30]. Both the adversarial loss of the entropy maps and the entropy of the pixel-wise prediction are minimized by [23].

In this work, we use the SOFA framework to reconstruct a virtual scene where a steerable robotic endoscope navigates through a 3D oropharynx model with its endoscopic vision recorded. Besides, real images are captured on a real-world phantom. To improve the performance of UDA models, we propose a domain adaptive Sim-to-Real framework called **I**oU-**R**anking **B**lend-**A**rt**F**low (IRB-AF) that includes a novel image-blending strategy and the style-transfer method. Firstly, we blend a small batch of real images into the simulation domain with the Intersection over Union (IoU)-Ranking Blend (IRB) mechanism. The mechanism sorts the resultant IoU among classes after training and refines the blending proportion for the next training iteration according to the sorting. In this regard, the potential of a limited number of mixed images can be fully utilized in the training process to improve the segmentation performance of real domains. Then, the style-transfer technique ArtFlow [31] is used to reduce the differences between the source and target domains. In practice, IRB-AF combines these two methods by first modifying the style of the source domain images and then blending target domain images via IRB. The overview of this framework is depicted in Fig. 1.

This work contributes the following:

- An image segmentation method targeting oropharyngeal organs with domain adaptive Sim-to-Real transfer;

- A novel IRB approach aiming at reducing the domain gap between virtual and real datasets for segmentation accuracy improvement;
- A style-transfer-based domain adaptive segmentation strategy to improve the network's training stability.
- An open-sourced dataset of endoscopic images generated from SOFA-based oropharynx model with style transfer from phantom (EISOST).

The rest of this paper is organized as follows. Section 2 introduces the data preparation for the image segmentation task. Section 3 illustrates the domain adaptive sim-to-real convention based on the proposed IRB-AF strategy. Section 4 showcases the performance of the IRB-AF strategy based on our dataset. The final section concludes the paper.

## 2 SOFA-Generated Virtual Dataset

As shown in Fig. 2A, robot-assisted TI employs a steerable flexible endoscope that works as a stylet to navigate to the repository tract instead of the digestive tract with the aid of endoscopic vision. The automation of such a process requires the robot to recognize or even segment the oropharyngeal organs by learning from a dataset with a great number of medical images. However, real-world endoscopic images of the oropharynx are often inaccessible due to ethical and privacy issues. While using different phantoms for image collection can be tedious and time-consuming, we propose to employ virtual endoscopic images as the dataset to train the segmentation network. In the virtual environment, one can conveniently generate an abundance of datasets [11] of different types for deep learning. Several virtual environments are capable of establishing 3D anatomical scenes, such as Blender [9], SOFA [10], etc. In previous work [32], we built an interactive environment of a soft robotic endoscope and oropharynx in the simulation (SOFA v22.06.99) as shown in Fig. 2B and 2D. The rich SOFA image data also motivated us to develop a virtual dataset-based segmentation approach, which could be useful in real-world soft robotic applications. However, the over-animated scenes pose a major challenge to the Sim-to-Real transfer.

### 2.1 Generation of Virtual Data

In recent years, researchers have been trying to develop large-scale virtual datasets to supplement the limited real datasets. For instance, [34, 35] release datasets and challenges to evaluate the state-of-the-art in surgical image segmentation and improve surgeon capabilities. With excellent virtual engines and rendering, modern computer games can also be used to generate virtual data sets to replace realistic scenes that require high costs to obtain [25]. In [32], a 3D oropharyngeal phantom modified based on [36] was imported into the virtual scene. The phantom includes three necessary oropharyngeal organs, namely, the uvula, epiglottis, and glottis. To reduce the expensive finite element computation, we trimmed the insignificant entities from the phantom,
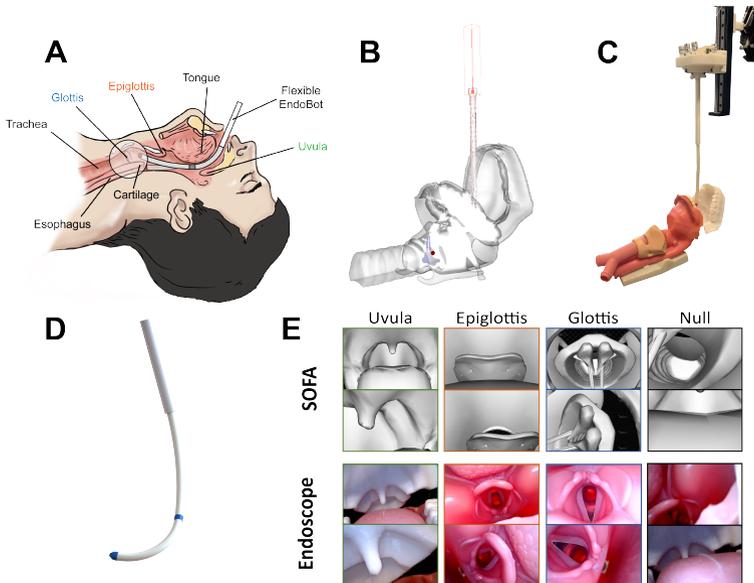
**Figure 2** (A) Using a flexible robotic endoscope (EndoBot) as a stylet to guide the next step's tracheal intubation. (B) SOFA scene. (C) Real-world scene. (D) CAD design of our EndoBot (adapted from [33]) that was used in the virtual and reality setup. (E) Dataset examples. "Null" indicates there is no target in the scene.

such as teeth and miscellaneous muscles. Using virtual images allows one to generate a large dataset with customized differentiation like the field of view, color, texture, and model variety, which significantly reduces the dependence on limited real images.

## 2.2  Preparation of Training Data

The training dataset is composed of two parts. The first part consists of the virtual images obtained by a flexible endoscopic robot navigating to the desired spot near the glottis in SOFA, as shown in Fig. 2B. The second part includes real images obtained by a similar setup and manner in the phantom environment in the real world, as shown in Fig. 2C. Some sample images are demonstrated in Fig. 2E. The size of the initial dataset labeled with bounding boxes is given in Table 1. As some images contain multiple organs, the available images in virtual and real sets are 1194 and 203, respectively. For the annotations, we provide coarse and fine annotations at the pixel level, including instance-level labels for oropharyngeal organs.

# 3  Domain Adaptive Sim-to-Real with IRB-AF

In this section, we introduce our domain adaption segmentation in two aspects. The proposed IRB approach is a compelling dataset blending strategy used for the Sim-to-Real training, while the image style-transfer is used to

**Table 1** Size of the Blended Dataset for Oropharyngeal Organ Recognition (Unit: Frame)

|                          | Uvula | Epiglottis | Glottis |
|--------------------------|-------|------------|---------|
| SOFA's Virtual Images    | 601   | 395        | 507     |
| Colored Real Images      | 119   | 69         | 82      |

further reduce the differences between the source domain and the target domain, thereby improving the domain adaption performances. We integrate the above two methods and propose IRB-AF that aligns the image distributions of different datasets in terms of content and style.

## 3.1 IRB: IoU-Ranking Blend

Conventional domain adaption methods focus on modifying the image features (e.g., color features, texture features, shape features, and spatial relationship features of an image) to reduce the distribution between the source and target domains. However, it is noticeable that only minor image differences occur in their dataset selection in the first place. This is often seen in applications like urban scenes and city traffic images. Yet, in the medical image field, obtaining virtual images that are significantly recognizable (i.e., being photo-realistic) to the real organs or anatomical features would be impractical. To deal with this problem, we propose to mix a small batch of images from the target domain during the training whenever using virtual or synthetic medical images, which can increase the content similarity between datasets. Since the distribution of features differs among images, the segmentation performance based on randomly-mixed training will be unstable. To maximize the usability and potential of a designated small batch of blended images, we propose the IRB based on the model training performance. The IoU-Ranking Blend method mixes images based on each class's segmentation testing results after training the model.

Since our dataset (see Table 1) contains three classes (excluding the background/null), we devise the following criteria to formulate the blending:
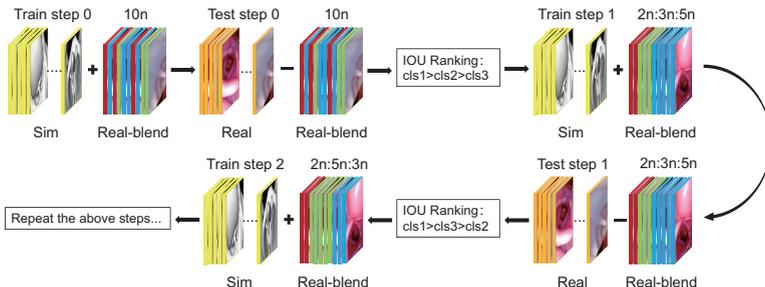
1. The number of blending images is a multiple of 10 to facilitate statistics and differentiate the outcomes of blending;
2. The proportion of blends is increased among classes to ensure lower-ranked classes have larger weights;
3. The least number of classes is a multiple of 2 in case the segmentation performance drops rapidly due to too insufficient proportions.

Therefore, a ratio of 5:3:2 for each class that fulfills the above criteria was chosen to investigate the IRB strategy, while the number of blends is 10 to 40. The resultant percentage of blending images in each subset is shown in Table 2. A maximum blending of 40 real images would ensure a 'virtual purity' of over 90% for each organ. Fig. 3 shows the flow chart of the blending sequence.

In the first step of training, randomly distributed images $I_b = \{x_t\}_{i=1}^{10n}$ where $x_t \in \mathbb{R}^{H \times W \times 3}$ are mixed in source domain. When the first training is

**Table 2** Percentage of blending images in each subset

| # of Blends | Uvula | Epiglottis | Glottis |
|---|---|---|---|
| 10 | 1.637% | 2.469% | 1.934% |
| 20 | 3.221% | 4.819% | 3.795% |
| 30 | 4.754% | 7.059% | 5.587% |
| 40 | 6.240% | 9.195% | 7.313% |



**Figure 3** Flow chart of IoU-Ranking Blend. *10n* and *cls* represent the number of blended images in the class, respectively. The colors of three *cls* are red, green, and blue. The above steps will be repeated until there is no new IoU ranking. The best mIoU will be selected from all test results.

completed, the segmentation results of the model on the remaining images in the target domain reflect the model's performance for each class. During the training process, the model's sensitivity on different oropharyngeal organs reflects differently, so the number of mixed images required for different classes to achieve the same testing performance varies. For some classes, if the differences between their source and target domains are minor, or the features are sharper compared to the background, they would require a lighter blend to improve the Sim-to-Real performance. However, other classes may require more to achieve the same effect.

After training, the resultant IoU of the test dataset is ranked in descending order. Then, the proportion of blends can be adjusted as deemed. For example, if the IoU ranking ends up with $\{\mathrm{IoU}_A > \mathrm{IoU}_B > \mathrm{IoU}_C\}$ where $\{A, B, C\}$ represents the classes, it indicates that class $C$ demonstrates relatively poor performance in the segmentation and the blending proportion of class $C$ needs to be increased. In contrast, the blending proportion of class $A$ can be reduced appropriately due to its better performance. For the next step's training set, the dataset is proportionally partitioned as $I_b = \{\{x_t^C\}_{i=1}^{5n}, \{x_t^B\}_{i=1}^{3n}, \{x_t^A\}_{i=1}^{2n}\}$. To avoid the model learning features outside the blending images, we randomly initialize the model parameters at the beginning of each training step. Although the performance of some classes could slightly reduce in the testing after adjusting the blending ratios, the rest could perform better, thus improving the training's mean of IoU (mIoU). Moreover, with the continuous adjustment of each training iteration, the ranking of IoUs will be gradually fixed. After that, we can opt for their optimal ratio so that the potential for a limited amount of mixed images can be exploited to the greatest extent.

The experiments in the next section will demonstrate the effectiveness of our strategy by comparing random and optimal proportions.

## 3.2 Image Style-Transfer for Domain Adaption

The style-transfer provides a new viewpoint to narrow the differences between the source and target domain in terms of image style. Here, the style essentially refers to the textures, colors, and visual patterns in images at various spatial scales, and they are considered low-level image features. A typical style-transfer strategy can be described as optimizing the weighted sum of content loss $\mathcal{L}_c$ and style loss $\mathcal{L}_s$ [37], which can be defined as follows:

$$\mathcal{L}_s = \|\mathcal{G}\left[\mathcal{A}(t)\right] - \mathcal{G}\left[\mathcal{A}(s)\right]\|_F^2 \tag{1}$$

$$\mathcal{L}_c = \|\mathcal{A}(t) - \mathcal{A}(c)\|_2^2. \tag{2}$$

Here, $\|...\|_F^2$ and $\|...\|_2^2$ are the Frobenius Norm and the Euclidean Norm, respectively. $s$ and $c$ indicate lower layers and higher layers in an image classification network. $\mathcal{A}(t)$ is the network activations, and $\mathcal{G}(t)$ denotes the Gram matrix of the network activations. The variables $s$ and $c$ represent the lower and higher layers, respectively. Neural style-transfer (NST) aims to transfer the style of the target domain into the image of the source domain to enhance the similarity of the features from different domains. However, due to the information loss caused by pooling [38], the training bias from the loss function [39], and the biased style-transfer module [40], the content leak problem occurs—the content information of the source domain may be lost during the style-transfer training process [31]. Therefore, we employ ArtFlow [31], which includes an unbiased feature transfer module and reversible neural flows to prevent content leaks during the style-transfer. ArtFlow establishes the Projection Flow Network (PFN) following the Glow model [41] and replaces the traditional encoder-decoder structure with a projection-reversion strategy. The components of PFN (additive coupling [42], invertible 1×1 convolution [41], and Actnorm [41]) are completely reversible, which also ensures that information is lossless when transmitted through PFN.

We have made numerous efforts to generate synthetic images. However, they still have a certain appearance gap with real images and cannot be directly used to train oropharyngeal organ segmentation. An example is shown in Fig. 4. With the help of ArtFlow, we try to convert the appearance of virtual images into real oropharyngeal organs' appearance, thereby enhancing the sense of photo-realistic of virtual data while preserving useful anatomical features for model training. A schematic of employing style-transfer to the virtual data is shown in Fig. 5. Image style-transfer helps us to reduce the differences between datasets from the low-level features. And the image content with high-level features is optimized by the IoU-Ranking Blend method in the previous subsection. In practice, IRB-AF first modifies the style of the source domain images and then blends a small number of images via IRB, which jointly modifies the distribution between the source domain and target

**Figure 4** Reducing the sim-to-real gap between the virtual and 'real' data using transfer learning based on Fourier Transform. There are still appearance gaps between virtual images and real images.
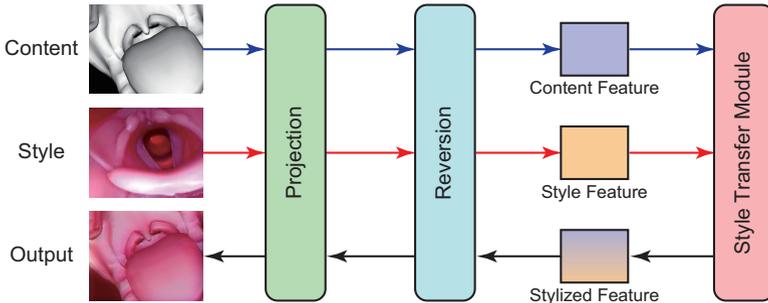


**Figure 5** The overview of employing ArtFlow in this work. Instead of the auto-encoder-based framework, ArtFlow applies a projection-transfer-reversion scheme. The projection module extracts content and style images' features which will be used to generate stylized features in Transfer Module. Afterward, the stylized feature is transferred to the stylized image via reversion inference.

domain, and improves the performance of the segmentation model from the perspective of the low-level and high-level features of the images, respectively. Therefore, our dataset contains three parts, i.e., endoscopic images generated from SOFA-based oropharynx model, endoscopic images captured from real-world phantom, and style-transferred virtual images. We have published this dataset under the name EISOST[1].

## 4  Experiments

To demonstrate the effectiveness of the proposed IRB-AF on domain adaptation segmentation, we conduct extensive experiments. We make a comparison before and after the introduction of IRB-AF on the state-of-the-art domain adaptive segmentation models. Moreover, an ablation experiment is made to qualitatively analyze the respective effects of IoU-Ranking Blend and style-transfer in domain adaptation segmentation.

---

[1]Endoscopic Images generated from SOFA-based oropharynx model with style transfer from phantom (EISOST) - https://github.com/gkw0010/EISOST-Sim2Real-Dataset-Release

**Table 3** Performances of Different SOTA Segmentation Methods Adopting our IRB-AF Strategy

| Method | Train Step | Intersection over Union | | | | | Accuracy | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | background | glottis | epiglottis | uvula | mIoU | background | glottis | epiglottis | uvula | mAcc |
| FDA [16] | 40-r | 94.990 | 72.030 | 54.540 | 65.660 | 71.805 | 97.73 | 85.19 | 59.44 | 87.34 | 82.425 |
| | 40-253 | 96.620 | 77.280 | 74.260 | 63.890 | 78.013 | 98.21 | 86.55 | 88.16 | 79.78 | 88.175 |
| | 40-235 | 96.490 | 78.440 | 73.980 | 66.590 | **78.875** | 98.15 | 86.4 | 88.47 | 80.18 | **88.300** |
| Advent [23] | 40-r | 96.147 | 80.790 | 65.060 | 67.900 | 77.474 | 98.71 | 87.79 | 70.17 | 80.78 | 84.363 |
| | 40-253 | 97.100 | 84.060 | 79.330 | 64.790 | **81.320** | 99.21 | 83.33 | 92.29 | 71.85 | 86.670 |
| | 40-235 | 96.600 | 78.290 | 76.630 | 67.460 | 79.745 | 98.69 | 82.41 | 89.35 | 79.22 | **87.418** |
| Cycada [24] | 40-r | 95.100 | 72.930 | 62.830 | 59.080 | 72.485 | 97.51 | 77.95 | 89.42 | 66.41 | 82.823 |
| | 40-253 | 94.640 | 75.320 | 72.380 | 64.890 | **76.808** | 97.02 | 87.31 | 72.6 | 76.16 | 83.273 |
| | 40-235 | 96.040 | 69.790 | 62.290 | 63.620 | 72.935 | 97.89 | 85.36 | 87.62 | 74.88 | **86.438** |

## 4.1  Evaluation and comparison on different networks

{*1) Implementation details:* We train and test three domain adaptive segmentation models, namely, FDA [16], ADVENT [23] and CyCADA [24], to validate the generality of our methods. The training is conducted on an Nvidia GeForce RTX 3090 GPU, and the batch size is set to 4 in all our experiments. We take DeepLabV2 [43] with ResNet101 [44] as the baseline. The optimization algorithm is selected as stochastic gradient descent (SGD), and the learning rate is $5 \times 10^{-3}$ with a momentum of 0.9 and weight decay of $5 \times 10^{-4}$. For the training parameters of ArtFlow, we will show them in Section 4.2.2. The larger the number of blended images, the lower the improvement in Sim-to-Real performance, which is demonstrated in Section 4.2.1. As a result, we just experimented with 40 blending images. If there is an improvement with this setting, it is inevitable with less blending quantity.

*2) Results and discussion:* The results are shown in Table 3 in which IRB-AF achieves higher mIoU and mean Accuracy (mAcc) in terms of Sim-to-Real performance. Compared to the original settings of these models, although increasing the number of blended images weakens the performance improvement, our method still increases 4.96%–9.85% in mIoU and 3.62%–7.13% in mAcc, respectively.

## 4.2  Ablation Experiment

### 4.2.1  IoU-Ranking Blend in FDA

*1) Implementation:* Fourier Domain Adaptation (FDA) is one of state-of-the-art domain adaptive segmentation models. The main idea of FDA is to use Fourier transform to swap the low-frequency spectra of the source domain and target domain images to reduce their differences. The experimental parameters are the same as those in Section 4.1. We bring in the IRB strategy and adapt it to the FDA.

Since the improvement would be significantly weakened when the blending quantity exceeds 40, we set 40 as the upper limit in our quantitative experiments. Note that a blend of 40 images accounts for only 3% of the training set. For each blended set, the same training and test will be conducted three times to reduce the randomness of training, and all records will be made.

*2) Results and discussion:* The mIoU of different blending groups is shown in Fig. 6. Due to the Sim-to-Real gap between the source and target domain, the original FDA may not achieve a satisfactory domain adaptation performance with low-frequency spectrum exchange. Consequently, the mIoU of the test set has always been low at about 21%. However, when ten real images are randomly mixed in the training set (the ratio between the original and the blended images is about 120:1), the mIoU of multiple tests can sharply increase to over 50%. Moreover, after introducing our IRB strategy, the model fully uses the blended images according to the classes, making the average mIoU exceed 56% with the lower limit exceeding 55%. When further increasing the mixed quantity, although the mIoUs fluctuate, their mean of upper and lower
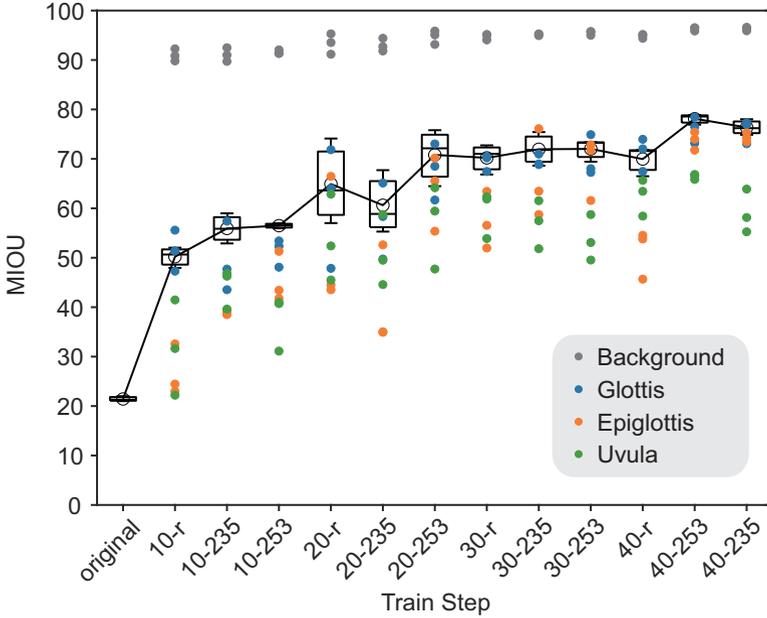
**Figure 6** The mIoU comparison between the IRB-FDA and original FDA. The horizontal axis represents the training sequence. Each element of the boxplot contains the results of three experiments, which are also shown by scatter points. The data of the line chart comes from the mean of mIoU.

**Table 4** Segmentation results on FDA Style-Transfer

| Train Step | mIoU | * | Stability | * |
|---|---|---|---|---|
| 10-r | 46.495 | ↓ 3.712 | 2.048 | ↓ 2.055 |
| 10-235 | 50.861 | ↓ 5.061 | 4.178 | ↓ 1.915 |
| 10-253 | 53.493 | ↓ 2.998 | 5.810 | ↑ 4.793 |
| 20-r | 61.432 | ↓ 3.486 | 3.480 | ↓ 13.620 |
| 20-235 | 62.694 | ↑ 2.070 | 5.558 | ↓ 6.870 |
| 20-253 | 72.005 | ↑ 1.199 | 4.313 | ↓ 7.008 |
| 30-r | 70.543 | ↑ 0.340 | 2.387 | ↓ 3.503 |
| 30-235 | 73.198 | ↑ 2.379 | 2.758 | ↓ 5.225 |
| 30-253 | 73.143 | ↑ 0.683 | 3.895 | ↓ 0.425 |
| 40-r | 72.996 | ↑ 3.034 | 3.373 | ↓ 1.945 |
| 40-235 | 74.723 | ↓ 1.642 | 6.153 | ↓ 3.013 |
| 40-253 | 75.823 | ↓ 2.296 | 0.632 | ↓ 1.328 |

*Compared to non-style-transfer.

limits gradually ascended. It is shown that IRB can increase mIoU by more than 5% under different blending quantities. The qualitative comparison of segmentation output is shown in Fig. 7.
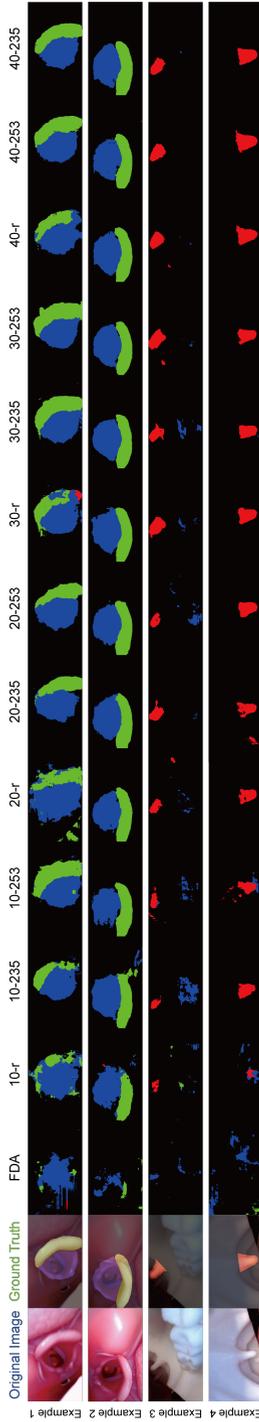
**Figure 7** Performance comparison on segmenting in the real domain of different oropharyngeal organs.

### 4.2.2 Style-Transfer in FDA

*1) Implementation details:* For the ArtFlow model, we set two blocks in the backbone network and eight flow modules in each block. The training of Art-Flow is carried out according to the [31] so that the optimization algorithm is Adam [45] instead of SGD, and the learning rate is $1 \times 10^{-4}$ with weight decay of $5 \times 10^{-5}$. The number of training iterations is set to 120,000. The configuration opts empirically.

After training the ArtFlow, each source image will be assigned to a random target image with the same oropharyngeal organs as style-transfer input to ensure that the style of the target domain can be fully utilized. Then, we set the style-transfer images as the new source domain for FDA training. Note that the experimental parameters are the same as those in Section 4.1.

*2) Results and discussion:* The experimental results are shown in Table 4. We focus on the mean of mIoU and the stability of repeated experiments with different mixed combinations to find the best training configuration for future applications. According to the results, the mean mIoU of the tests does not generally increase but increases when the blending proportion becomes larger. It shows that the improvement of style-transfer on mIoU is limited, and the contribution to the segmentation accuracy improvement of the target domain is mainly due to the blended images from the target domain. However, the data fluctuations of multiple tests of style-transfer are much lower than that of previous experiments, which shows that the style-transfer can improve the training stability and reduce the impact of randomness. The results indicate that the style-transfer could enhance the stability of semi-supervised domain adaptive segmentation training by mitigating variations between datasets, although the improvement in segmentation accuracy is limited.

## 5 Conclusion

In this paper, we propose a domain adaptive Sim-to-Real framework called IRB-AF. The framework includes an image blending strategy called IoU-Ranking Blend (IRB) and style-transfer method ArtFlow. The IRB solves the accuracy degradation problem caused by the significant difference between simulation and real datasets during unsupervised domain adaptive segmentation of oropharyngeal organs. By sorting the segmentation results between classes, a more appropriate blending strategy can be formulated to maximize the potential of employing a limited number of blended images. The style-transfer method ArtFlow is introduced to reduce the differences between datasets further. The Sim-to-Real segmentation results show that our proposed method improves the performance of the existing domain adaptive segmentation models. Furthermore, we also found that the style-transfer can enhance training stability and alleviate the impact of randomness.

Although we have made some progress in oropharyngeal organ segmentation, the accuracy of the model still has much space for improvement. Future work can be expected to enrich domain adaptation segmentation methods in

the medical field. Besides, the enrichment of real datasets is also significant to further improve the accuracy of the model.

# Declarations

- Conflict of interest/Competing interests: No benefits in any form have been or will be received from a commercial party related directly or indirectly to the subject of this manuscript.
- Ethics approval: Ethical approval was not sought for the present study because this article does not contain any studies with human or animal subjects.
- Authors' contributions: G.W., J.L., and H.R. conceived the concepts. G.W., T.R., J.L., and L.B. advised on the design and implementation of the experiments. G.W., T.R., and J.L. conducted experiments and analyzed the data. G.W., T.R., J.L., and L.B. wrote the manuscript. All authors read, edited, and discussed the manuscript and agree with the claims made in this work. H.R. coordinated and supervised the research.

# References

[1] Thomas, E.B., Moss, S.: Tracheal intubation. Anaesth. Intensiv. Care Med. **15**(1), 5–7 (2014)

[2] Caplan, R.A., Benumof, J.L., Berry, F.A., Blitt, C.D., Bode, R.H., Cheney, F.W., Connis, R.T., Guidry, O.F., Nickinovich, D.G., Ovassapian, A.: Practice guidelines for management of the difficult airway. Anesthesiology **98**(1269-1277), 2 (2003)

[3] Lu, B., Li, B., Chen, W., Jin, Y., Zhao, Z., Dou, Q., Heng, P.-A., Liu, Y.: Toward image-guided automated suture grasping under complex environments: A learning-enabled and optimization-based holistic framework. IEEE Transactions on Automation Science and Engineering **19**(4), 3794–3808 (2021)

[4] Lai, J., Lu, B., Chu, H.K.: Variable-stiffness control of a dual-segment soft robot using depth vision. IEEE/ASME Transactions on Mechatronics **27**(2), 1034–1045 (2021)

[5] Lu, B., Li, B., Dou, Q., Liu, Y.: A unified monocular camera-based and pattern-free hand-to-eye calibration algorithm for surgical robots with rcm constraints. IEEE/ASME Transactions on Mechatronics **27**(6), 5124–5135 (2022)

[6] Yu, B.X., Liu, Y., Zhang, X., Zhong, S.-h., Chan, K.C.: Mmnet: A model-based multimodal network for human action recognition in rgb-d videos. IEEE Transactions on Pattern Analysis and Machine Intelligence (2022)

[7] Asgari Taghanaki, S., Abhishek, K., Cohen, J.P., Cohen-Adad, J., Hamarneh, G.: Deep semantic segmentation of natural and medical images: a review. Artif. Intell. Rev. **54**(1), 137–178 (2021)

[8] Frangi, A.F., Tsaftaris, S.A., Prince, J.L.: Simulation and synthesis in medical imaging. IEEE Trans. Med. Image. **37**(3), 673–679 (2018)

[9] Rehman, M., Arsenault, L., Javan, R.: Organs in color: utilizing free software and emerging multi jet fusion technology to color and surface label 3d-printed anatomical models. J. Digit. Imaging **35**(6), 1611–1622 (2022)

[10] Duriez, C.: Control of elastic soft robots based on real-time finite element method. In: Proc. IEEE Int. Conf. Robot. Autom. (ICRA), pp. 3982–3987 (2013)

[11] Zhao, W., Queralta, J.P., Westerlund, T.: Sim-to-real transfer in deep reinforcement learning for robotics: a survey. In: Proc. IEEE Symp. Ser. Comput. Intell. (SSCI), pp. 737–744 (2020)

[12] Ganry, L., Hersant, B., Quilichini, J., Leyder, P., Meningaud, J.: Use of the 3d surgical modelling technique with open-source software for mandibular fibula free flap reconstruction and its surgical guides. J. Stomatol. Oral Maxillofac. Surg. **118**(3), 197–202 (2017)

[13] Pierri, R., Nogueira, L., Balan, I., Iwaki, L., *et al.*: Bimaxillary orthognatic surgery planned with the software blender, through the addon ortogonblender. Int. J. Oral Maxillofac. Surg. **48**, 254 (2019)

[14] Chen, X., Hu, J., Jin, C., Li, L., Wang, L.: Understanding domain randomization for sim-to-real transfer. arXiv preprint arXiv:2110.03239 (2021)

[15] Tobin, J., Fong, R., Ray, A., Schneider, J., Zaremba, W., Abbeel, P.: Domain randomization for transferring deep neural networks from simulation to the real world. In: IEEE/RSJ Int. Conf. Intell. Robot. Syst. (IROS), pp. 23–30 (2017). IEEE

[16] Yang, Y., Soatto, S.: Fda: Fourier domain adaptation for semantic segmentation. In: Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 4085–4095 (2020)

[17] Geng, B., Tao, D., Xu, C.: Daml: Domain adaptation metric learning. IEEE Trans. Image Process. **20**(10), 2980–2989 (2011)

[18] Long, M., Cao, Y., Wang, J., Jordan, M.: Learning transferable features with deep adaptation networks. In: Proc. Int. Conf. Mach. Learn. (ICML), pp. 97–105 (2015). PMLR

[19] Zellinger, W., Grubinger, T., Lughofer, E., Natschläger, T., Saminger-Platz, S.: Central moment discrepancy (cmd) for domain-invariant representation learning. arXiv preprint arXiv:1702.08811 (2017)

[20] Zou, Y., Yu, Z., Kumar, B., Wang, J.: Unsupervised domain adaptation for semantic segmentation via class-balanced self-training. In: Proc. Eur. Conf. Comput. Vis. (ECCV), pp. 289–305 (2018)

[21] Wu, Z., Han, X., Lin, Y.-L., Uzunbas, M.G., Goldstein, T., Lim, S.N., Davis, L.S.: Dcan: Dual channel-wise alignment networks for unsupervised scene adaptation. In: Proc. Eur. Conf. Comput. Vis. (ECCV), pp. 518–534 (2018)

[22] Sankaranarayanan, S., Balaji, Y., Jain, A., Lim, S.N., Chellappa, R.: Learning from synthetic data: Addressing domain shift for semantic segmentation. In: Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 3752–3761 (2018)

[23] Vu, T.-H., Jain, H., Bucher, M., Cord, M., Pérez, P.: Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation. In: Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 2517–2526 (2019)

[24] Hoffman, J., Tzeng, E., Park, T., Zhu, J.-Y., Isola, P., Saenko, K., Efros, A., Darrell, T.: Cycada: Cycle-consistent adversarial domain adaptation. In: Proc. Int. Conf. Mach. Learn. (ICML), pp. 1989–1998 (2018). Pmlr

[25] Richter, S.R., Vineet, V., Roth, S., Koltun, V.: Playing for data: Ground truth from computer games. In: Proc. Eur. Conf. Comput. Vis. (ECCV), pp. 102–118 (2016). Springer

[26] Ros, G., Sellart, L., Materzynska, J., Vazquez, D., Lopez, A.M.: The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In: Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 3234–3243 (2016)

[27] Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban scene understanding. In: Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 3213–3223 (2016)

[28] Li, Y., Yuan, L., Vasconcelos, N.: Bidirectional learning for domain adaptation of semantic segmentation. In: Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 6936–6945 (2019)

[29] Zhu, X.J.: Semi-supervised learning literature survey. Technical report, University of Wisconsin-Madison Department of Computer Sciences (2005)

[30] Springenberg, J.T.: Unsupervised and semi-supervised learning with categorical generative adversarial networks. arXiv preprint arXiv:1511.06390 (2015)

[31] An, J., Huang, S., Song, Y., Dou, D., Liu, W., Luo, J.: Artflow: Unbiased image style transfer via reversible neural flows. In: Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 862–871 (2021)

[32] Lai, J., Ren, T.-A., Yue, W., Su, S., Chan, J.Y.K., Ren, H.: Sim-to-real transfer of soft robotic navigation strategies that learns from the virtual eye-in-hand vision. Unpublished (2023)

[33] Lai, J., Lu, B., Zhao, Q., Chu, H.K.: Constrained motion planning of a cable-driven soft robot with compressible curvature modeling. IEEE Robot. Autom. Lett. **7**(2), 4813–4820 (2022)

[34] Allan, M., Shvets, A., Kurmann, T., Zhang, Z., Duggal, R., Su, Y.-H., Rieke, N., Laina, I., Kalavakonda, N., Bodenstedt, S., et al.: 2017 robotic instrument segmentation challenge. arXiv preprint arXiv:1902.06426 (2019)

[35] Allan, M., Kondo, S., Bodenstedt, S., Leger, S., Kadkhodamohammadi, R., Luengo, I., Fuentes, F., Flouty, E., Mohammed, A., Pedersen, M., et al.: 2018 robotic scene segmentation challenge. arXiv preprint arXiv:2001.11190 (2020)

[36] University of Dundee, School of Medicine: Pharynx and Floor of Mouth. https://skfb.ly/6QXqr. Accessed: 2022-08-01

[37] Ghiasi, G., Lee, H., Kudlur, M., Dumoulin, V., Shlens, J.: Exploring the structure of a real-time, arbitrary neural artistic stylization network. arXiv preprint arXiv:1705.06830 (2017)

[38] Li, Y., Fang, C., Yang, J., Wang, Z., Lu, X., Yang, M.-H.: Universal style

transfer via feature transforms. Adv. Neural Info. Processing Syst. **30** (2017)

[39] Huang, X., Belongie, S.: Arbitrary style transfer in real-time with adaptive instance normalization. In: Proc. IEEE Int. Conf. Compt. Vis. (ICCV), pp. 1501–1510 (2017)

[40] Liao, J., Yao, Y., Yuan, L., Hua, G., Kang, S.B.: Visual attribute transfer through deep image analogy. arXiv preprint arXiv:1705.01088 (2017)

[41] Kingma, D.P., Dhariwal, P.: Glow: Generative flow with invertible 1x1 convolutions. Adv. Neural Info. Processing Syst. **31** (2018)

[42] Dinh, L., Krueger, D., Bengio, Y.: Nice: Non-linear independent components estimation. arXiv preprint arXiv:1410.8516 (2014)

[43] Chen, L.-C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. IEEE Trans. Pattern Anal. Mach. Intell. **40**(4), 834–848 (2017)

[44] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 770–778 (2016)

[45] Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)