Improving Realism in Abdominal Ultrasound Simulation Combining a Segmentation-Guided Loss and Polar Coordinates training

Santiago Vitale^{1,2}, José Ignacio Orlando^{1,2}, Emmanuel Iarussi^{1,3}, Alejandro Díaz^{1,4}, Ignacio
 Larrabide^{1,2}

Correspondence: Santiago Vitale, Pladema, UNICEN, Tandil, Argentina. Campus Universi tario, Paraje Arroyo Seco. Email: santiago.vitale@pladema.exa.unicen.edu.ar

10

g

6

Abstract

Background: Ultrasound (US) simulation helps train physicians and medical students 11 in image acquisition and interpretation, enabling safe practice of transducer manipulation 12 and organ identification. Current simulators generate realistic images from reference scans. 13 Although physics-based simulators provide real-time images, they lack sufficient realism, 14 while recent deep learning-based models based on unpaired image-to-image translation im-15 prove realism but introduce anatomical inconsistencies. Purpose: We propose a novel 16 framework to reduce hallucinations from generative adversarial networks (GANs) used on 17 physics-based simulations, enhancing anatomical accuracy and realism in abdominal US 18 simulation. Our method aims to produce anatomically consistent images free from artifacts 19 within and outside the field of view (FoV). Methods: We introduce a segmentation-guided 20 loss to enforce anatomical consistency by using a pre-trained Unet model that segments ab-21 dominal organs from physics-based simulated scans. Penalizing segmentation discrepancies 22 before and after the translation cycle helps prevent unrealistic artifacts. Additionally, we 23 propose training GANs on images in polar coordinates to limit the field of view to non-blank 24 regions. We evaluated our approach on unpaired datasets comprising 617 real abdominal 25 US images from a SonoSite-M turbo v1.3 scanner and 971 artificial scans from a ray-casting 26 simulator. Data was partitioned at the patient level into training (70%), validation (10%), 27 and testing (20%). Performance was quantitatively assessed with Frechet and Kernel In-28 ception Distances (FID and KID), and organ-specific χ^2 histogram distances, reporting 29 95% confidence intervals. We compared our model against generative methods such as 30 CUT, UVCGANv2, and UNSB, performing statistical analyses using Wilcoxon tests (FID 31 and KID with Bonferroni-corrected $\alpha = 0.01$, χ^2 with $\alpha = 0.008$). A perceptual realism 32 study involving expert radiologists was also conducted. Results: Our method significantly 33 reduced FID and KID by 66% and 89%, respectively, compared to CycleGAN, and by 34% 34 and 59% compared to the leading alternative UVCGANv2 ($p \ll 0.01$). No significant dif-35 ferences (p > 0.008) in echogenicity distributions were found between real and simulated 36 images within liver and gallbladder regions. The user study indicated our simulated scans 37 fooled radiologists in 36.2% of cases, outperforming other methods. Conclusions: Our 38 segmentation-guided, polar-coordinates-trained CycleGAN framework significantly reduces 39

¹National Scientific and Technical Research Council (CONICET), Buenos Aires, Argentina ²Pladema Institute UNICEN Tandil Argentine

²Pladema Institute, UNICEN, Tandil, Argentina

³Laboratory of Artificial Intelligence, University Torcuato Di Tella, Buenos Aires, Argentina

⁴Facultad de Ciencias de la Salud, UNICEN, Olavarría, Argentina

hallucinations, ensuring anatomical consistency and realism in simulated abdominal US
 images, surpassing existing methods.

42 I. Introduction

Abdominal ultrasound (US) is an essential non-invasive imaging technique for diagnosing various
 abdominal conditions². Effective clinical use requires specialists skilled in both image acquisition
 and interpretation. Typically, this training involves hands-on sessions with patients or volunteers,
 limiting scalability due to the need for devices and human subjects³.

US simulation has emerged as a valuable training tool, allowing medical professionals to 47 safely develop technical skills and procedural proficiency without needing real patients or equip-48 ment^{4,5}. Simulators provide repeatable and controlled scenarios where users practice device 49 manipulation⁵, organ localization⁶, and complex procedures⁷. Hence, these risk-free platforms 50 contribute to improved clinical outcomes and increased confidence of clinicians to handle the 51 complexities of real-world medical imaging. Additionally, US simulation supports applications 52 like image registration⁸ and expands datasets for deep learning⁹, highlighting the necessity for 53 realistic simulated images. High-fidelity simulations are crucial for achieving anatomical accuracy 54 in training and clinical applications. 55

Several methods have been proposed to generate synthetic US images, such as ray-casting 56 algorithms applied to CT volumes¹⁰ or ray-tracing methods on deformable meshes^{11,12}. While 57 efficient, these physics-based approaches lack the realism needed for clinical training in image in-58 terpretation and diagnosis¹³. Recent generative models using convolutional neural networks have 59 gained considerable attention for their enhanced realism¹⁴. These models have primarily focused 60 on simulating images from specific areas of interest, such as intravascular¹⁵ or fetal examina-61 tions¹⁶, and regions like the brain⁸, ovaries¹⁷, kidneys¹⁸, and musculoskeletal structures¹⁹. More 62 complex regions, such as the abdominal cavity, have been less explored using these techniques. 63 Previously, we applied an unpaired CycleGAN-based translation model²⁰ to improve ray-casting 64 simulations¹. While this refinement enhances the overall realism of the generated images, it 65 suffers from hallucinated features typical of distribution matching losses²¹. In particular, the re-66 sulting scans include both unexpected organs in anatomically incorrect areas and distorted edges 67 of the observable area captured by the device, typically referred to as the field of view (FoV). 68



FoV Deformation Rullucinations

Figure 1: Examples of different artificial US scans obtained with a ray-casting model, our previous approach based on a standard CycleGAN model¹, and our improved method using a segmentation-guided loss and polar coordinates. Anatomical masks are provided as reference.

In this study, we propose some novel changes to our previous approach¹, with the goal of 69 eliminating hallucinations and enabling the generation of anatomically consistent abdominal US 70 scans from ray-casting-based simulations²². We achieve this by introducing a novel segmentation-71 guided loss, which leverages a pretrained Unet²³ segmentation model that penalizes differences 72 between organ segmentations in the input image and its reconstructed versions after completing 73 a full translation cycle. This information propagates through the entire cycle, compelling the 74 fake-to-realistic generator to preserve anatomical consistency in the forward cycle. Otherwise, 75 any hallucinations and unrealistic artifacts introduced will be propagated in the realistic-to-fake 76 generator, and detected by the segmentation network. This aids to eliminate one of the sources 77 of mistake, the hallucinations within organs. Additionally, we propose training our models di-78 rectly in polar coordinates to remove irrelevant blank areas outside the field of view (FoV) and 79 reduce artifacts in these regions. In summary, our key contributions with respect to our previous 80 CycleGAN approach are threefold: 81

1) We introduce an objective term that enforces consistency between organ segmentations in the input scan, and its equivalent after the realism improvement transformation. To the best of our knowledge, such an "asymmetrical" approach for backpropagating anatomical knowledge have not been applied before to reduce hallucinations.

2) We adapted the training process to be directly applied to images in polar coordinates, eliminating empty spaces outside the FoV and preventing FoV deformations.

3) We demonstrate the model's generalization capability—unlike our previous patient-specific approach, the new model can be trained on multiple subjects and effectively applied to simulate ⁹⁰ new individuals

Experimental results confirm that our approach significantly improves realism and anatomical
 accuracy over previous CycleGAN-based methods¹ and an improved ray-casting-based simulator.

⁹³ II. Related work

The proposed model builds on top of our previous approach for improving realism in US sim-94 ulations¹. Originally, the method was based on a standard CycleGAN model²⁰, which allows 95 image-to-image translation with unpaired samples. Technically, this model features two GANs, 96 each defined by its own pair of generators and discriminators. In this context, it formally trans-97 lates images from the domain $\mathcal A$ of **a**rtificially generated US images to another set $\mathcal R$ of **r**eal US 98 images (both described in Section IV.A.1.), and viceversa. Formally, let $G_{A \rightarrow R}$ be the generator 99 that translates an artificial image $a \in \mathcal{A}$ to \mathcal{R} , and $D_{\mathcal{R}}$ the discriminator that distinguishes 100 between real images r and the translated ones $G_{\mathcal{A}\to\mathcal{R}}(a)$. On the other hand, let $G_{\mathcal{R}\to\mathcal{A}}$ be the 101 generator that translates an image $r \in \mathcal{R}$ to the domain \mathcal{A} while trying to avoid being detected 102 by a discriminator $D_{\mathcal{A}}$. In the original CycleGAN definition, both pairs of networks are simul-103 taneously trained by optimizing a linear combination of losses, including a standard adversarial 104 penalty $\mathcal{L}_{GAN}\text{,}$ a cycle-consistency term $\mathcal{L}_{cyc}\text{,}$ and the identity loss $\mathcal{L}_{idt}\text{.}$ 105

 \mathcal{L}_{GAN} is defined per each pair of generator and discriminator as follows:

$$\mathcal{L}_{\mathsf{GAN}}(G_{\mathcal{A}\to\mathcal{R}}, D_{\mathcal{R}}, \mathcal{A}, \mathcal{R}) = \mathbb{E}_{r \sim p_{\mathsf{data}}(r)} [\log(D_{\mathcal{R}}(r) - 1)^{2}] + \mathbb{E}_{a \sim p_{\mathsf{data}}(a)} [\log(D_{\mathcal{R}}(G_{\mathcal{A}\to\mathcal{R}}(a))^{2}], \\ \mathcal{L}_{\mathsf{GAN}}(G_{\mathcal{R}\to\mathcal{A}}, D_{\mathcal{A}}, \mathcal{A}, \mathcal{R}) = \mathbb{E}_{a \sim p_{\mathsf{data}}(a)} [\log(D_{\mathcal{A}}(a) - 1)^{2}] + \mathbb{E}_{r \sim p_{\mathsf{data}}(r)} [\log(D_{\mathcal{A}}(G_{\mathcal{R}\to\mathcal{A}}(r))^{2}],$$

$$(1)$$

where \mathbb{E} stands for the expected value of each corresponding data distribution, and each term is based on the least-squares GAN loss (LSGAN)²⁴, which prevents vanishing gradient issues.

To allow unpaired image-to-image translation, the training scheme incorporates an additional cycle-consistency loss \mathcal{L}_{cyc} . This term enforce that translations produced by one generator are reversible and retain the original domain's characteristics (Step 1, Figure 2). Formally, a forward cycle translates an image $a \in \mathcal{A}$ previously translated to domain \mathcal{R} back to \mathcal{A} (that is, $a \rightarrow G_{\mathcal{A} \rightarrow \mathcal{R}}(a) \rightarrow G_{\mathcal{R} \rightarrow \mathcal{A}}(G_{\mathcal{A} \rightarrow \mathcal{R}}(a)) \approx a$). Similarly, a reverse cycle ensures an image $r \in \mathcal{R}$ translated to domain \mathcal{A} is brought back to \mathcal{R} (by doing $r \to G_{\mathcal{R} \to \mathcal{A}}(r) \to G_{\mathcal{A} \to \mathcal{R}}(G_{\mathcal{R} \to \mathcal{A}}(r)) \approx$ 114 r). \mathcal{L}_{cyc} can then be defined as the sum of two losses:

$$\mathcal{L}_{\mathsf{cyc}}(G_{\mathcal{A}\to\mathcal{R}}, G_{\mathcal{R}\to\mathcal{A}}) = \mathbb{E}_{r \sim p_{\mathsf{data}}(r)}[||G_{\mathcal{A}\to\mathcal{R}}(G_{\mathcal{R}\to\mathcal{A}}(r)) - r||_{1}] + \mathbb{E}_{a \sim p_{\mathsf{data}}(a)}[||G_{\mathcal{R}\to\mathcal{A}}(G_{\mathcal{A}\to\mathcal{R}}(a)) - a||_{1}].$$

$$(2)$$

The identity loss \mathcal{L}_{idt} regularizes the generators towards identity mappings, thereby biasing the models towards learning only what is needed to accurately generate realistic images:

$$\mathcal{L}_{\mathsf{idt}}(G_{\mathcal{A}\to\mathcal{R}}, G_{\mathcal{R}\to\mathcal{A}}) = \mathbb{E}_{a \sim p_{\mathsf{data}}(a)}[||G_{\mathcal{R}\to\mathcal{A}}(a) - a||_{1}] \\ + \mathbb{E}_{r \sim p_{\mathsf{data}}(r)}[||G_{\mathcal{A}\to\mathcal{R}}(r) - r||_{1}].$$
(3)

117 III. Methods

Figure 2 depicts a schematic representation of the training and test phases of our abdominal 118 US simulation model. Our approach requires two sets of unpaired images for training, one with 119 intermediate artificial US images (\mathcal{A}) and one with real US scans (\mathcal{R}). The first one is obtained 120 by applying a ray-casting-based simulation algorithm²² on cross-sectional 2D slices retrieved from 121 multiple 3D CT scans and their associated 3D segmentation masks, based on the coordinates of 122 an artificial probe. These 2D images are then transformed to polar coordinates to eliminate blank 123 spaces outside the FoV and avoid the generative model hallucinating features outside the area. 124 Images in \mathcal{A} , and their associated set of 2D segmentation masks (\mathcal{M}), are used offline to train a 125 segmentation model S, which remains fixed later on while training our SG-CycleGAN model. This 126 approach learns to map images from $\mathcal A$ to $\mathcal R$ and viceversa using two image-to-image translation 127 models $G_{\mathcal{A}\to\mathcal{R}}$ and $G_{\mathcal{R}\to\mathcal{A}}$. The optimization minimizes a combined loss: a cycle-consistency term 128 (\mathcal{L}_{cyc}) and a segmentation-guided term (\mathcal{L}_{sg}) . The latter penalizes anatomical inconsistencies 129 by comparing the predicted segmentations of the artificial scan and its reconstruction. During 130 the testing phase, we input an intermediate artificial ultrasound image into the $G_{\mathcal{A}\to\mathcal{R}}$ generator, 131 provided it was generated using the same ray-casting approach utilized during training. Doing so 132 will yield a more realistic version of the original image. 133

In this study we propose to improve the previous approach by incorporating a novel segmentation-guided term (\mathcal{L}_{sg}) that enforces consistency between segmentation predictions of



Figure 2: Schematic representation of the training (top) and testing (bottom) phases of our proposed approach for improving abdominal US simulation using a novel anatomically consistent image-to-image translation model.

images from A and their reconstructed counterparts. By penalizing discrepancies between the
segmentation maps of the original and reconstructed fake images, the model is encouraged to
maintain realistic anatomical features throughout the cycle during fake-to-real translation process.
This consistency reduces the likelihood of introducing unrealistic artifacts and hallucinations, as
any deviations from expected anatomical structures are penalized during training.

Figure 2 illustrates the proposed additional asymmetric objective, which integrates infor-141 mation about tissue locations within $a \in \mathcal{A}$ and enforces anatomical consistency between the 142 original input and its reconstructed counterpart. Let S(x) represent a deep neural network that 143 produces a pixel-wise multiclass segmentation of a given input image x. The model S is trained 144 offline using images $a \in \mathcal{A}$ and the corresponding segmentation masks, remaining fixed during 145 the CycleGAN training phase (Step 0, Figure 2). During CycleGAN training, each image $a \in \mathcal{A}$ 146 is translated into the \mathcal{R} domain by the generator $G_{\mathcal{A}\to\mathcal{R}}$. The resulting image is subsequently 147 translated back into the original domain by the generator $G_{\mathcal{R}\to\mathcal{A}}$ to obtain the reconstructed 148 image (Step 1, Figure 2). Both the original image a and its reconstruction are segmented by S, 149 yielding anatomical masks which are subsequently compared for consistency (Step 2, Figure 2). 150 Formally, our proposed loss function, \mathcal{L}_{sg} , penalizes differences between S(a) (the segmentation 151

¹⁵² map of an image $a \in \mathcal{A}$) and $S(G_{\mathcal{R}\to\mathcal{A}}(G_{\mathcal{A}\to\mathcal{R}}(a)))$ (the segmentation map of the reconstructed ¹⁵³ image after completing the full cycle):

$$\mathcal{L}_{\mathsf{sg}}(G_{\mathcal{A}\to\mathcal{R}}, G_{\mathcal{R}\to\mathcal{A}}, S) = -\sum S(a) \log \left(S(G_{\mathcal{R}\to\mathcal{A}}(G_{\mathcal{A}\to\mathcal{R}}(a))) \right), \tag{4}$$

By means of this term, anatomical knowledge is transferred between generators, forcing $G_{A\to R}$ to preserve organs shape so that the reverse cycle through $G_{R\to A}$ does not produce an inconsistent sample.

In summary, the proposed training scheme is defined as a linear combination of the CycleGAN
 loss terms and the novel objective introduced above, namely:

$$\mathcal{L}(G_{\mathcal{A}\to\mathcal{R}}, G_{\mathcal{R}\to\mathcal{A}}, D_{\mathcal{A}}, D_{\mathcal{R}}, S) = \mathcal{L}_{\mathsf{GAN}}(G_{\mathcal{A}\to\mathcal{R}}, D_{\mathcal{R}}, \mathcal{A}, \mathcal{R}) + \mathcal{L}_{\mathsf{GAN}}(G_{\mathcal{R}\to\mathcal{A}}, D_{\mathcal{A}}, \mathcal{R}, \mathcal{A}) + \lambda_{\mathsf{cyc}} \cdot \mathcal{L}_{\mathsf{cyc}}(G_{\mathcal{A}\to\mathcal{R}}, G_{\mathcal{R}\to\mathcal{A}}) + \lambda_{\mathsf{idt}} \cdot \mathcal{L}_{\mathsf{idt}}(G_{\mathcal{A}\to\mathcal{R}}, G_{\mathcal{R}\to\mathcal{A}}) + \lambda_{\mathsf{sg}} \cdot \mathcal{L}_{\mathsf{sg}}(G_{\mathcal{A}\to\mathcal{R}}, G_{\mathcal{R}\to\mathcal{A}}, S),$$

$$(5)$$

where λ_{cyc} , λ_{idt} and λ_{sg} are hyperparameters that control the relative importance of each term in the final loss. Supplementary materials provide a flow chart with a visual representation of the calculation of the global loss throughout the training process.

Notice that the identity loss and the segmentation-guided loss serve different purposes in the model. The identity term enforces that each generator maintains features from the target domain that are already present in the source domain. On the other hand, our segmentation-guided loss focuses on preserving anatomical structure when transitioning from one domain to another.

¹⁶⁶ IV. Experimental setup

¹⁶⁷ IV.A. Materials

¹⁶⁸ IV.A.1. Artificial US dataset

We generated a set of simulated images using 13 contrast-enhanced CT volumes (60% male) from the VISCERAL's Anatomy3 Challenge dataset²⁵. To standardize the images, we manually cropped them to retain only the abdominal cavity, from the thoracic diaphragm to the pelvic inlet. Hounsfield Units (HUs) were then normalized to [0, 1] using histogram equalization. A 2D Gaussian smoothing kernel of size 50×50 pixels (ranging from 34×34 mm to 37.5×37.5 mm, depending on voxel size) with a standard deviation of 2.5 pixels (approximately 1.7–1.875 mm) was applied to reduce high-frequency noise and improve uniformity.

For intermediate simulation, an artificial probe was placed at various abdominal locations to extract clinically relevant cross-sectional slices from both the CT scans and their segmentation masks (see Segmentation masks dataset). These slices served as inputs for a modified version of the ray-casting simulation algorithm by Rubí *et al.*²² (see supplementary materials for further details). This process generated 926 artificial scans.

¹⁸¹ IV.A.2. Segmentation masks dataset

The anatomical masks used correspond to the cross-sectional slices extracted from the silver corpus segmentations of the 13 CT volumes in Artificial US dataset. The original dataset included segmentations of the spleen, liver, gallbladder, aorta, and kidneys. To provide additional anatomical references, we manually segmented the rib cage and spine.

186 IV.A.3. Real US scan dataset

¹⁸⁷ Our real US dataset comprised 617 prospectively collected images from 11 volunteers (60% male, ¹⁸⁸ age = 27 ± 3 years) with no known abdominal conditions. A specialist acquired these scans ¹⁸⁹ during routine abdominal exams using a SonoSite-M turbo v1.3 US Scanner (FUJIFILM, Bothell, ¹⁹⁰ USA). The scanning parameters differed from those used in the ray-casting model, as there is no ¹⁹¹ direct correspondence between the device and the algorithm. All images were exported in JPEG ¹⁹² format at 640×480 pixels.

¹⁹³ IV.A.4. Dataset preprocessing and partition

To standardize spatial dimensions and align with the transducer's curvature, we applied a 194 Cartesian-to-Polar transformation to both artificial and real ultrasound scans. This process in-195 volved calculating the center, inner and outer radii, and angular range (θ) for each image. For 196 simulated images, these parameters were derived from the ray-casting algorithm, while for real 197 scans, they were manually extracted using FoV masks. This transformation corrected the trans-198 ducer's curvature and removed non-informative areas (see supplementary materials for a graphical 199 explanation). The final images were resized to 256×256 pixels and randomly partitioned at the 200 patient level into training (70%, 8 patients), validation (10%, 2 patients), and test (20%, 3 201 patients) subsets. 202

²⁰³ IV.B. Architectures

²⁰⁴ IV.B.1. Generator architecture

We evaluated three generator architectures, all based on a Unet encoder-decoder network. The 205 first was a standard Unet²³ (Unet in our experiments), where the decoder was replaced with 206 nearest-neighbor upsampling followed by a convolutional layer to prevent checkerboard artifacts¹. 207 The second was a modified Unet with bottleneck layers and residual connections²⁶ (ResUnet in our 208 experiments), implemented in two width variations. Lastly, we included the densely connected 209 image-to-image translation generator by Dangi et al.²⁷ (DenseUnet in our experiments). All 210 generators used a tanh activation function. Further architectural details are provided in the 211 supplementary materials. 212

²¹³ IV.B.2. Discriminator architecture

Following previous studies^{16,17,28}, we employed a 70 × 70 patchGAN²⁹ as the discriminator. The network consists of four convolutional blocks, each with a 4 × 4 kernel and a stride of 2. Instance normalization was used instead of batch normalization, as it has been shown to enhance diversity and prevent mode collapse^{30,31}. Each block applies Leaky-ReLU activation, as commonly done in patchGANs²⁹, progressively reducing spatial dimensions while increasing feature maps to 64, 128, 256, and 512, respectively. A final 1-filter convolution, followed by a sigmoid activation function, produces the output probability for each patch.

²²¹ IV.B.3. Segmentation model

The segmentation network S is based on a Unet architecture. The encoder consists of four 222 convolutional blocks with 64, 128, 256, and 512 filters, each followed by 2×2 max-pooling 223 for downsampling. Each block comprised a sequence of a 3×3 convolutional layer, a batch 224 normalization operation, and a ReLU activation, repeated twice. A bottleneck layer with 1024 225 filters precedes the decoder, which uses nearest-neighbor upsampling followed by convolutional 226 layers with progressively fewer filters, from 512 down to 64. A final 1×1 convolutional layer 227 produces class logits, converted into probabilities using softmax activation. The network was 228 trained to segment the liver, spleen, gallbladder, aorta, and kidneys. Since the kidney consists of 229 two ultrasound-differentiable structures—the hyperechoic renal pelvis and the hypoechoic renal 230 cortex—we treated them as separate classes, using weak annotations for each (see supplementary 231 for further details). 232

²³³ IV.C. Model configuration

Hyperparameters were empirically selected based on validation set performance using Fréchet 234 Inception Distance (FID). In tied cases, we visually inspected the results and chose parameters 235 that produced more realistic and anatomically consistent images. Coefficients λ_{cyc} , λ_{idt} , and λ_{sg} 236 were experimentally fixed to 10, 0.5, and 0.5, respectively. We found that a higher λ_{cyc} improved 237 cycle consistency in image translations. We trained the model for 200 epochs using Adam³² 238 optimization with an initial learning rate of 2×10^{-4} and a batch size of 4. After 100 epochs, the 239 learning rate was reduced linearly by $\frac{1}{101}$. The segmentation network S was trained offline using 240 a multiclass cross-entropy objective, Adam optimization with an initial learning rate of 1×10^{-4} , 241 and a batch size of 16 for 150 epochs. The learning rate was decreased by a factor of 0.5 every 242 time that the performance plateaued for 20 epochs, measured by the average Dice coefficient. 243 Hyperparameters were selected to maximize the average Dice score for all organs in the validation 244 set. 245

All CNNs, including the segmentation network, were implemented in Pytorch 1.10.0 and trained on a desktop workstation with an AMD Ryzen 9 5900X CPU and an NVIDIA GeForce RTX 3060 GPU (12GB RAM).

²⁴⁹ IV.D. Baselines for comparison

We compared SG-CycleGAN with the ray-casting-based method¹⁰ used to generate the input 250 images and four state-of-the-art image-to-image translation models. Given the limited number of 251 models available for unpaired datasets in this task, we focused on CycleGAN-based approaches, 252 which have shown promise in US simulation. First, we compared SG-CycleGAN to our previously 253 published CycleGAN¹, trained with images in Cartesian coordinates. Second, we included the 254 Contrastive Unpaired Translation $(CUT)^{33}$ model, which has been used as a baseline for obstetric 255 US simulation³⁴. To incorporate recent advances, we tested the UNet Vision Transformer cycle-256 consistent GAN (UVCGANv2)³⁵, which integrates a U-Net with a Vision Transformer encoder. 257 Finally, we included the Unpaired Neural Schrödinger Bridge (UNSB)³⁶, a diffusion-based model 258 that provides an alternative to GANs and has been applied to US simulation³⁷. This selection 259 covers both standard approaches and recent innovations in generative modeling for US simulation. 260

²⁶¹ IV.E. Evaluation metrics & statistical analysis

Assessing the quality and realism of simulated US scans is challenging, as in any image genera-262 tion task^{38,39}. The most widely used metrics are Fréchet Inception Distance (FID)⁴⁰ and Kernel 263 Inception Distance (KID)⁴¹, which have been applied in various US studies^{16,17,34,42}. Both met-264 rics quantify the statistical distance between feature distributions of real and artificial images. 265 extracted from an Inception v3⁴³ network pretrained on ImageNet. This comparison captures 266 macro-level differences in speckle noise texture. A lower FID score indicates that the generated 267 images better resemble real US scans in terms of noise and echogenicity. We used the intermedi-268 ate 768-feature layer to avoid highly specialized low-level descriptors³⁴. For evaluation, we used 269 the validated TorchFidelity implementation⁴⁴. Statistical significance was assessed using one-270 tailed Wilcoxon signed-rank tests, with Bonferroni correction⁴⁵ adjusting the significance level 271 from 0.05% to 0.01% (5 comparisons). Effect sizes were evaluated using Cohen's d^{46} , which 272 measures differences relative to the pooled standard deviation. According to Cohen's criteria, 273 0.2 represents a small effect size and indicates that the difference between groups is noticeable 274 but not substantial; 0.5 represents a medium effect size, suggesting a moderate difference that is 275 likely to be meaningful in most contexts; and 0.8 represents a large effect size, indicating a sub-276 stantial difference between groups, which is often considered to be practically significant. Very 277 small effects (below 0.2) indicate negligible differences that may not have practical relevance. 278

Values greater than 1, on the other hand, are considered very large, and highlight a difference
 that is both statistically and practically significant.

The χ^2 distance⁴⁷, commonly used in US simulation¹⁶, quantifies dissimilarities between image histograms:

$$\chi^{2}(h_{A}||h_{B}) = \frac{1}{2} \sum_{l=1..d} \frac{(h_{A}[l] - h_{B}[l])^{2}}{h_{A}[l] + h_{B}[l]},$$
(6)

where d is the number of histogram bins (50 in our case). While alternatives like Jensen-Shannon (JS) divergence⁴⁸ compare entire histograms, we opted for χ^2 as it is more sensitive to relative differences in individual bins.

Histogram-based methods are affected by intensity shifts and contrast variations⁴⁹. To 286 evaluate potential mismatches in echogenicity, we compared intensities locally within segmented 287 gallbladder, liver, and kidney regions. Segmentation masks were slightly eroded using a 5 imes 5288 structuring element to reduce edge irregularities. Pairwise χ^2 distances from real US images 289 were used as reference values. To ensure fair comparisons, scans with minimal tissue representa-290 tion were excluded, and histograms were normalized by the number of pixels within each mask. 291 Statistical significance was tested using a one-tailed Wilcoxon rank-sum test with a Bonferroni-292 corrected threshold of 0.0083 (6 comparisons), alongside effect size analysis via Cohen's d. Notice 293 that if simulations are realistic, the χ^2 distance distribution for each organ should closely match 294 that of real scans, showing no significant differences. For all metrics, 95% confidence intervals 295 (95% CI) were computed using bootstrap resampling (N = 1000). 296

Finally, we assessed anatomical accuracy by comparing segmentation masks from our method and standard CycleGAN using mean Intersection over Union (mIoU). These masks were created by manually segmenting a set of 16 simulated images and comparing the resulting organ masks with those used as input to the physical model.

³⁰¹ IV.F. User study-based evaluation

We further evaluated our approach through a custom-made online user study, implemented using the jsPsych JavaScript library⁵⁰ (see supplementary materials for further details). The study comprised two experiments. The first experiment assessed experts' ability to distinguish real

from simulated US images. Participants were shown a US scan and asked to classify it as 305 real or simulated. The dataset included 45 images: 15 real US scans, 15 generated by our 306 approach, and 15 by the original CycleGAN model. Classification accuracy was measured as the 307 fraction of correctly identified real and fake images. The second experiment evaluated anatomical 308 preservation in the generated images. Experts were presented with two simulated scans-one 309 generated with and one without the segmentation-guided term—and asked to select the scan 310 with better anatomical preservation. The original segmentation mask was provided as a reference. 311 This test included 10 scan pairs covering typical abdominal capture windows such as intercostal, 312 subcostal margin, longitudinal, oblique, and transverse views. A total of 16 clinicians, all experts 313 in US imaging, participated in the study, most of whom were affiliated with Sociedad Argentina 314 de Ultrasonido en Medicina y Biología (SAUMB). 315

³¹⁶ V. Results

We conducted a comprehensive evaluation of the proposed approach, using both qualitative and quantitative approaches. Our method was compared to state-of-the-art techniques outlined in Subsection IV.D., elaborated upon in Subsection V.A.2.. Additionally, an ablation study was carried out to evaluate the impact of each design choice on the final results, detailed in Subsection V.B..

³²² V.A. Simulation performance

³²³ V.A.1. Qualitative evaluation

Figure 3 presents example simulations generated using the original CycleGAN in Cartesian coordinates¹, the same model in polar coordinates, and our SG-CycleGAN. The samples correspond to different abdominal windows commonly used in clinical analyses.

The Cartesian CycleGAN results exhibit FoV deformations in all scans, mainly as irregular edges (Figures 3 (a) and (d)). In some cases, these distortions remove anatomical structures, such as part of the liver (Figures 3 (a) and (c)), the aorta (Figures 3 (c) and (d)), or the kidneys (Figures 3 (a) and (e)). Alternatively, using polar coordinates ensures images that are consistent with the input FoV, with both the standard CycleGAN and our proposed SG-CycleGAN, preserving (a)

(b)

(c)

(d)

(e)



Figure 3: Qualitative results for abdominal US simulation obtained using a standard Cycle-GAN trained in Cartesian and polar coordinates and our proposed SG-CycleGAN approach. Dotted lines indicate inconsistent organs (yellow) and their improved counterparts (green). From top to bottom: (a) right subcostal margin, (b) longitudinal, (c) oblique, and (d,e) right and left intercostal acquisition windows.

Gallbladder

Anatomical preservation

Spleen

Kidney

Inconsistencies

Liver

Aorta

Bones

³³² all the organs that are present in the images.

Figures 3 (a)–(c) show that standard CycleGANs introduce inhomogeneities in the liver, appearing as hallucinated shadows (Figures 3 (a) and (c)) or anatomically inconsistent hyperechoic structures (Figures 3 (b) and (c)). Our segmentation-guided approach preserves liver structure, maintaining homogeneous echogenicity (green contours).

Figures 3 (d) and (e) present results for windows that include part of the kidney. Training with 337 Cartesian coordinates produces unrealistic kidneys, with artifacts such as hyperechoic reflections 338 that are inconsistent with this anatomical area (Figure 3 (d)), or intensities of the renal pelvis 339 below the usual echogenicities (Figure 3 (e)). Similarly, Figures 3 (c) and (d) show poor aorta 340 representations, which disappear into larger anechoic areas. While polar coordinates mitigate this 341 issue, they still generate anatomical inconsistencies (e.g., hyperechoic streaks in the kidney or 342 diffuse spleen edges in Figure 3 (e)). Our approach better preserves organs, yielding anatomically 343 accurate results for the gallbladder (Figure 3 (a)), aorta (Figures 3 (b) and (c)), bones (Figures 3 344 (c)–(e)), kidneys (Figures 3 (a), (d), and (e)), and spleen (Figure 3 (e)). On this last area, a 345 better scattering effect can be observed on top of the artifact generated by the skin (top green 346 arrow), as well as more defined interfaces at the bottom (bottom green arrow). 347

Figure 4 visually compares our method to other baselines. Further qualitative results are pro-348 vided in the supplementary materials. The previous CycleGAN model reduces the FoV, removing 349 image regions (e.g., the missing backbone in Figure 4 (c) or the truncated kidney in Figure 4 350 (e)). CUT better preserves anatomical structures but still producing hallucinations such as a 351 hyperechoic artifact in the liver (Figure 4 (a)) and an anechoic tubular formation in the kidney 352 (Figure 4 (c)). It also fails to maintain spleen integrity (Figure 4 (e)). UVCGANv2 struggles 353 to maintain structures, reducing gallbladder size (Figure 4 (b)) and distorting kidneys (Figures 4 354 (c) and (e)). The UNSB model preserves structures like the liver, gallbladder, and vessels (see 355 Figure 4 (a), (b) and (d), respectively), but struggles with kidney structures, where it halluci-356 nates anechoic formations (Figure 4 (c) and (e)). Additionally, it fails to simulate the skin layer 357 artifacts, which are captured in the other models. Finally, our model corrects the FoV limitations 358 observed in our previous version, while also preserving all the anatomical structures provided in 359 the ray-casting based input. 360



Figure 4: Qualitative examples for each model and their associated segmentations as reference. Yellow arrows indicate inconsistencies.

³⁶¹ V.A.2. Quantitative evaluation

Table 1 compares our approach to all baselines detailed in Section IV.D.. While all generative models outperform the physics-based simulator, our SG-CycleGAN achieves statistically significant reductions in FID (80%) and KID (97%) (p < 0.01). Very large effect sizes (Cohen d = 84.24and 79.97) further support these findings. Among baselines, UVCGAN performed best, but still lags behind our method, with substantial Cohen effect sizes (d = 10.43 for FID and d = 9.20 for KID).

Our SG-CycleGAN also exhibits χ^2 distances within the liver and gallbladder that closely 368 resemble those observed between real images (Table 1). Figure 5 (A) provides a detailed analysis 369 of this metric for each tissue, with colored boxplots representing the distribution of pairwise χ^2 370 distances between simulated and real US images, and gray boxplots representing the reference 371 distribution between real scans. Although these cannot be compared directly one other for being 372 calculated using different samples, it can be observed that methods incorporating generative 373 approaches achieve χ^2 distances that distribute approximately similar as in real images, for all 374 organs. All generative models produce echogenicities in the gallbladder that are statistically 375 indistinguishable from those in real US images, with p values greater than 0.021. However, it 376 should be noted that our model, like CUT and UVCGANv2, presents closer mean values and 377 a very low Cohen's d value (< 0.09), indicating a very small effect size compared to the rest 378

models, which have values close to 0.2. Within the liver, our SG-CycleGAN and UNSB model 379 achieved distances comparable to the distances observed between real images. In this case, 380 the statistical tests performed between these models and real US images showed no statistically 381 significant differences, with p > 0.127 for all comparisons. On the contrary, performing the same 382 comparison between CUT, CycleGAN and UVCGANv2 exhibited statistically significant differences 383 (p < 0.008). Nonetheless, all models exhibit a small effect size (Cohen's d < 0.16), with the 384 UNSB model standing out with a Cohen's d of 0.01. In the the kidney, the CycleGAN and the 385 UNSB did not exhibit statistically significant differences when compared to real US images, with 386 p > 0.183, showing a very small effect size (Cohen d < 0,09). 387

To further illustrate echogenicity similarities, Figure 5 (b) presents histograms of cumulative intensity distributions for each organ. These histograms differ from those used for organ-specific χ^2 comparisons in Table 1 and Figure 5 (A). Consistent with previous observations, our model produces intensities that closely resemble real images, particularly in the liver and gallbladder. For the kidney, UNSB outputs are more similar to real images.

We also report training and inference time comparisons in the supplementary material. SG-CycleGAN increased training time from 95 seconds (standard CycleGAN) to 127 seconds per epoch, similar to CUT and notably faster than UVCGANv2 and UNSB. For inference, SG-CycleGAN and CycleGAN were the fastest at 0.0813 seconds per scan, while other models required 2–3 times longer.

³⁹⁸ V.B. Ablation analysis

³⁹⁹ V.B.1. Quantitative evaluation

Table 2 presents results from CycleGAN models trained with different strategies. Models using 400 polar coordinates (rows 2 to 4) achieved better FID and KID scores than the Cartesian-based 401 model (row 1). However, improvements in χ^2 distances appeared only in the gallbladder and 402 liver when combined with the segmentation-guided loss and LSGAN objective. Regarding adver-403 sarial loss, LSGAN outperformed the vanilla loss (Table 2 rows 2 and 3). The best results were 404 achieved by incorporating the segmentation-guided loss (row 4), which further improved FID and 405 KID scores. In terms of anatomical preservation relative to ground truth label maps, our model 406 achieved a higher overall mIoU (0.68) than the standard CycleGAN (0.59). For individual organs 407



Figure 5: Organ-wise quantitative evaluation. (A) Box plots illustrating the distribution of pairwise χ^2 distances between pairs of simulated and real US images for each organ of interest (colored), and between pairs of real US images (gray).p-values are included for all comparison where no statistical differences observed. (B) Histograms representing the distribution of echogenicity values for each organ, for simulated (colored) and real (gray) images.



Figure 6: FID and KID results for different architectures of generator models. Each network was trained with (right) and without (left) our proposed loss term. The bubble size is proportional to the number of parameters of each model, indicated in millions (M) on top of each one.

Table 1: Quantitative comparison of the proposed model with respect to other alternatives in terms of FID and KID distances (lower value, marked as \downarrow , is better), and mean χ^2 distances for different organs of interest. Asterisks (*) next to FID and KID values indicate statistically significant differences, when compared to our approach (p < 0.01). χ^2 distances between pairs of real scans are included as a reference (closer to this reference is better). Daggers (†) in χ^2 distances indicate no statistical differences with the real scans (p > 0.008). Sub-indices indicate Cohen's d values. Best values are indicated in bold. Last row corresponds to the number of real and simulated US used to calculate each metric.

Model	$\mathrm{FID}\downarrow$	$\mathrm{KID} \downarrow (\times 10^{-3})$	$\chi^2 \; [95\% \; {\rm CI}]$		
Model	[95% CI]	[95% CI]	Liver	Kidney	Gallbladder
Ray-casting ²²	1.73 [1.69 - 1.77] $*_{84.24}$	$5.02 \ [4.85 - 5.19] *_{79.97}$	$0.23 \ [0.02 - 0.54]_{0.21}$	$0.17 \ [0.03 - 0.34] \dagger_{0.06}$	$0.09 \ [0.0 - 0.50]_{0.90}$
$CycleGAN^1$	$0.99 \ [0.96 - 1.03]_{48.63}^{*}$	2.61 [2.48 - 2.74] $*_{46.83}$	$0.21 \ [0.07 - 0.41]_{0.13}$	$0.19 \ [0.03 - 0.47]_{0.09}$	$0.28 \ [0.02 - 0.65] \dagger_{0.22}$
CUT ³³	$0.80 \ [0.76 - 0.84] *_{27.25}$	1.90 [1.74 - 2.06] $*_{26.79}$	$0.21 \ [0.06 - 0.42]_{0.12}$	$0.28 \ [0.05 - 0.55]_{0.74}$	$0.26 \ [0.0 - 0.66] \dagger_{0.09}$
UVCGANv2 ³⁵	$0.48 \ [0.45 - 0.51] *_{10.43}$	$0.69 \ [0.59 - 0.79]_{9.20}^{*}$	$0.17 \ [0.03 - 0.43]_{0.16}$	$0.23 \ [0.03 - 0.52]_{0.41}$	$0.25 \ [0.00 - 0.54] \dagger_{0.03}$
UNSB ³⁶	$0.95 \ [0.90 - 0.99] *_{36.23}$	2.42 [2.26 - 2.58] $*_{37.31}$	$0.19 \ [0.04 - 0.46] \dagger_{0.01}$	0.18 [0.02 - 0.50] $\dagger_{0.03}$	$0.22 \ [0.05 - 0.45] \dagger_{0.23}$
SG-CycleGAN (ours)	$0.33 \ [0.32 - 0.35]$	$0.28 \ [0.25 - 0.31]$	0.18 [0.05 - 0.40] $\dagger_{0.13}$	$0.22 \ [0.03 - 0.48]_{0.33}$	$0.25 \ [0.00 - 0.53] \dagger_{0.07}$
Real US	-	-	$0.19 \ [0.00 - 0.51]$	$0.18 \ [0.00 - 0.45]$	$0.24 \ [0.00 - 0.48]$
Number of scans $\mathcal{R} \mathcal{A}$	213 213	213 213	40 90	16 48	12 28

⁴⁰⁸ (liver, kidney, gallbladder), our model outperformed CycleGAN with IoU values of 0.84, 0.93, and ⁴⁰⁹ 0.86, respectively, compared to 0.75, 0.89, and 0.81, demonstrating superior anatomical fidelity. ⁴¹⁰ We also analyzed the impact of different generator architectures by comparing FID and KID ⁴¹¹ metrics across network types and backbone sizes (Figure 6). The standard Unet consistently out-⁴¹² performed ResUnets and DenseUnets in FID and KID scores. Additionally, adding \mathcal{L}_{sg} improved ⁴¹³ performance across all networks, except for the DenseUnet with the smallest capacity (0.2 million ⁴¹⁴ parameters).

⁴¹⁵ V.B.2. Qualitative effect of using polar coordinates

To assess the impact of using polar instead of Cartesian coordinates for training, Figure 7 com-416 pares input simulations from the ray-casting algorithm with their improved versions using both 417 alternatives. All scans share the same FoV, outlined in green. With Cartesian coordinates, the 418 model either restricts the original FoV (left edge of image (a)) or introduces organs outside of it 419 (bottom of both scans). In Figure 7 (b), the network hallucinates large shadowed areas near the 420 contours while partially preserving original image details (yellow arrow, left side), creating false 421 tissue reflections beyond the incorrect FoV. In contrast, images generated in polar coordinates 422 remain confined to the pre-defined FoV, free of deformations or hallucinated artifacts. Figure 7 423

Table 2: Evaluation of the ablation test in terms of FID, KID (lower value, marked as \downarrow , is better) and mean χ^2 distances for different organs. Asterisks (*) next to FID and KID values indicate statistically significant differences (p < 0.016), when compared to our approach. χ^2 distances between pairs of real scans are included as a reference (closer to this reference is better).Daggers (†) in χ^2 distances indicate no statistical differences with the real scans (p > 0.012). Sub-indices indicate Cohen's *d* values. The best values are indicated in bolds. The last row corresponds to the number of real and simulated US images used to calculate each metric respectively.

Model	Adversarial	Coordinate	$\mathrm{FID}\downarrow$	$\mathrm{KID}\downarrow(\times10^{-3})$		χ^2	
	loss	space	[95% CI]	[95% CI]	Liver	Kidney	Gallbladder
CG	Vanilla	С	$0.99 \ [0.96-1.03]_{48.63}^{*}$	2.61 [2.48 - 2.74] $*_{46.83}$	$0.21 \ [0.07 - 0.41]_{0.13}$	$0.19 \ [0.03 - 0.46]_{0.09}$	$0.28 \ [0.02 - 0.65]_{0.22}$
CG	Vanilla	Р	$0.73 \ [0.71 - 0.76] *_{36.12}$	1.82 [1.72 - 1.92] $*_{38.93}$	$0.26 \ [0.09 - 0.52]_{0.49}$	$0.29 \ [0.03 - 0.56]_{0.85}$	$0.23 \ [0.00 - 0.53]_{0.08}$
CG	LSGAN	Р	$0.42 \ [0.40 - 0.44] *_{7.48}$	$0.38 \ [0.33 - 0.43]_{9.64}^{*}$	0.21 [0.05 - 0.44] $\dagger_{0.05}$	$0.22 \ [0.04 - 0.47]_{0.25}$	$0.27 \ [0.00 - 0.54]_{0.04}$
SG	LSGAN	Р	$0.33 \ [0.32 - 0.35]$	$0.28 \ [0.25 - 0.31]$	$0.18 [0.05 \text{-} 0.40] \dagger_{0.13}$	$0.22 \ [0.03 - 0.48]_{0.34}$	$0.25 \ [0.00 - 0.53]$ † $_{0.07}$
Real US		-	-	$0.19 \ [0.00 - 0.51]$	$0.18 \ [0.00 - 0.45]$	$0.24 \ [0.00 - 0.48]$	
Number of scans $(\mathcal{R} \mathcal{A})$		213 293	213 293	40 90	16 48	6 27	

Abbreviations: CG, Standard CycleGAN; SG, SG-CycleGAN; Vanilla, Jensen-Shannon divergence loss; LSGAN, least squares GAN loss; C, cartesian; P, Polar

also includes patches illustrating speckle noise patterns. Unlike input simulated scans, Cartesian based outputs exhibit randomly oriented patterns, misaligned with the US transducer. Polar
 coordinates mitigate this issue, producing more realistic lateral speckle orientations consistent
 with the convex transducer's azimuthal angle.

$_{\tt 428}$ V.B.3. Qualitative effect of the generator architecture

Figure 8 compares results from SG-CycleGAN using different generator architectures. All gen-429 erative models enhance overall brightness, but the ResUnet introduces bright artifacts that are 430 anatomically inconsistent, such as in the renal pelvis (Figure 8 (a), yellow arrow) and an un-431 segmented region (Figure 8 (b), green arrow). Additionally, ResUnet produces an overly blurred 432 and poorly defined speckle pattern. In contrast, the Unet and DenseUnet backbones yield better 433 intensity distributions while preserving organ shapes and boundaries (e.g., the aorta in Figure 8 434 (a), red arrow). The kidney (Figure 8 (a), yellow arrows) also shows well-defined interfaces both 435 externally and within the renal pelvis. These networks generate more realistic speckle noise pat-436 terns (e.g., in the liver, Figure 8 (b)), though DenseUnet hallucinates interfaces in unsegmented 437 areas compared to Unet (green arrow). 438



Figure 7: Comparison of simulated images with CycleGANs trained on different coordinate systems. Green boundaries indicate the original FoV.

439 V.B.4. User study

Figure 9 presents the user survey results. Figure 9.A) shows bar charts of user accuracy in clas-440 sifying images—generated by CycleGAN, SG-CycleGAN, or real US, as fake or real. The average 441 and standard deviation for each type are also included. Lower accuracy indicates more frequent 442 misclassification of fake images as real and vice versa. Most participants correctly identified 443 CycleGAN-generated images as fake with high accuracy (98%), reflecting their lower realism. 444 However, for SG-CycleGAN images, accuracy averaged 63.75%, meaning 36.25% were mistaken 445 for real. This trend is also evident in real US scan classification, where expert accuracy averaged 446 below 80%. Figure 9.B) presents a pie chart summarizing radiologists' responses on anatomi-447 cal preservation. When asked about the preservation of the anatomy in fake images generated 448 with both synthetic methods, 81.6% of cases favored SG-CycleGAN to be more anatomically 449 consistent over CycleGAN. 450



Figure 8: Comparison of simulation results obtained using an SG-CycleGAN with Unet, ResUnet, or DenseUnet based generator.



Figure 9: User study results. A) Classification accuracy for each simulation model and real scans as a bar per participant. Additionally, a bar plot with average accuracy per method. B) Pie chart comparing responses about which generative model performs better in terms of anatomy conservation.

⁴⁵¹ VI. Discussion

⁴⁵² VI.A. Effect of our segmentation-guided loss

Simulating abdominal US images is challenging. While physics-based approaches generate anatomically plausible images, their echogenicities remain unrealistic. In contrast, CycleGANs enhance visual quality but introduce hallucinated artifacts that distort the underlying anatomy¹. These inconsistencies appear as non-uniform echogenicity patterns within organs (yellow dotted lines in Figure 3), a common issue in unpaired models relying on distribution-matching losses²¹.

To alleviate this issue, we proposed a segmentation-guided loss, penalizing segmentation mismatches before and after completing the cycle. This term prevents the generator $\mathcal{G}_{\mathcal{A}\to\mathcal{R}}$ to introduce artifacts that cannot be removed through the reversed cycle $\mathcal{G}_{\mathcal{R}\to\mathcal{A}}$, without any extra annotation. The anatomical labels from ray-casting simulations suffice for training. As seen in
 Figure 3 (green lines), our approach produces well-defined organ interfaces and homogeneous
 speckle noise patterns. Compared to existing methods (Figure 4), our loss function preserves
 anatomical structures while preventing hallucinated patterns within them.

Quantitatively, our model significantly reduces FID and KID scores by 66% and 89%, re-465 spectively $(p \ll 0.01)$, as shown in Table 2. Our model not only presents the lowest FID and 466 KID values, but when comparing with the others, we obtain high Cohen's d values (> 9.20), 467 which imply a very large effect size between the simulations of our model and the others. Lower 468 FID scores suggest improved statistical similarity to real images, resulting from the reduction in 469 hallucinations and unusual artifacts in the constrained areas. This ensures that simulated images 470 closely resemble real ones, making them more valuable for medical training. Furthermore, our 471 segmentation-guided loss enhances anatomical accuracy, improving mIoU by up to 15.3% over 472 standard CycleGAN. This advantage is reinforced by our user study, where SG-CycleGAN was 473 rated as more anatomically consistent in 81.9% of cases compared to the standard CycleGAN. 474

475 VI.B. Impact of training in polar coordinates

Another key contribution of our work is migrating CycleGAN training from Euclidean to polar coordinates. As illustrated in Figure 3 and highlighted by the yellow arrows in the intermediate column of Figure 7, CycleGANs trained in Euclidean coordinates produce jagged edges, distort the FoV, or introduce warped regions. This occurs because the network lacks prior knowledge of the region of interest, making it difficult to distinguish between acoustic shadows and empty areas outside the FoV.

Training in polar coordinates addresses this issue by constraining the network's focus to 482 the relevant area while excluding blank spaces. This prevents the model from having to learn 483 the FoV shape itself, allowing for better utilization of its capacity. As a result, the model 484 more accurately mimics speckle noise patterns (see zoomed patches in Figure 7) and better 485 leverages the segmentation-guided loss, as evidenced by improvements in FID and KID values 486 (Table 2). Additionally, since areas outside the FoV are absent in the input, the network naturally 487 avoids generating artifacts in those regions. This is evident in Figure 7, where all images exhibit 488 consistent FoVs without irregularities or hallucinations beyond the designated area. 489

⁴⁹⁰ VI.C. Influence of the generator architecture

Our approach proves effective across different generator architectures and network sizes, consis-491 tently improving FID and KID values when using the segmentation-guided loss (Figure 6). Among 492 the tested architectures, the standard Unet outperformed ResUnet and DenseUnet, aligning with 493 previous findings¹. As illustrated in Figure 8, Unet generates anatomically more coherent outputs 494 than ResUnet. This discrepancy is likely due to the absence of skip connections in ResUnet's 495 bottleneck layers. Without these connections, the decoder must reconstruct anatomical struc-496 tures using only low-level features from earlier layers, leading to information loss. The bottleneck 497 acts as a lossy compression of the input, making it difficult for the decoder to reconstruct organs 498 without introducing unrealistic artifacts. 499

500 VI.D. Advantages of SG-CycleGAN

Integrating all our proposed modifications into the standard CycleGAN framework resulted in a 501 robust generative model that outperforms several state-of-the-art approaches in realism. We 502 compared SG-CycleGAN against recent deep learning models, including Vision Transformers 503 (UVCGANv2) and conditional diffusion models (UNSB). As shown in Table 1, these methods 504 reduced FID and KID scores relative to the ray-casting model, with Vision Transformers achiev-505 ing the largest improvement. However, SG-CycleGAN achieved the lowest FID and KID values 506 $(p = 0.33 \times 10^{-3} \text{ and } p = 0.25 \times 10^{-3}$, respectively), with a very large effect size (Cohen's 507 d > 9.20). Our model also closely matches real ultrasound (US) echogenicity distributions. As 508 shown in Table 2, χ^2 tests indicate no statistically significant differences in liver and gallbladder 509 echogenicities between SG-CycleGAN-generated images and real scans (p > 0.008). The effect 510 size is minimal (Cohen's d = 0.07 for the gallbladder and d = 0.13 for the liver), suggesting 511 that our model generates tissue echogenicities within the natural variability of real US images. 512 While UNSB achieves a slightly better match for the liver (d = 0.01), our approach still per-513 forms competitively, as showed in Figure 5 (B). From a qualitative perspective, SG-CycleGAN 514 produces more realistic scans. If the generated images were easily distinguishable from real ones, 515 expert classification accuracy would approach 100%. While this was true for standard CycleGAN, 516 experts misclassified 36% of SG-CycleGAN images as real (Figure 9). This suggests that our 517 model generates anatomically consistent and realistic US scans, making it a promising tool for 518 improving ultrasound training applications. 519

⁵²⁰ VI.E. Limitations

The primary limitation of this study is its focus on healthy subjects, as all experiments were conducted on individuals without pathologies or lesions. While we have demonstrated that our approach reduces hallucinations in simulated scans, we cannot guarantee the same for pathological cases or lesions. Future work should extend the evaluation to pathological cases to assess the method's robustness in simulating complex anatomical variations. Nevertheless, preventing hallucinations in healthy cases is already a promising step forwards, as it avoids introducing unrealistic artifacts that could be interpreted as pathologies.

It should be pointed out also that, despite the model exhibiting a substantial reduction in 528 hallucinations compared to its original counterpart, we still observed unrealistic features occurring 529 outside the segmented areas (e.g., around organ interfaces in Figure 3 (d)). In our current setup, 530 we utilized masks for six different tissues available in our set of volumetric segmentations, so 531 anatomical inconsistencies outside these regions are to be expected. In particular, we observed 532 this phenomenon to occur in areas such as the stomach or the pancreas, which are not segmented 533 in our training set. Clinically, these inaccuracies could affect the usefulness of the simulations 534 in training scenarios where detailed anatomy of these regions is critical, such as as in surgical 535 planning or procedural training, where a precise understanding of the anatomical structures is 536 crucial. 537

Nevertheless, notice that the proposed approach is general enough to include any other 538 organ without considerable modifications, should they are already available for the ray-casting 539 based simulator (e.g. by segmenting the organs from the input CT scans). While these masks 540 are essential for training the segmentation-guided CycleGAN, notice they do not increase the 541 annotation costs beyond that already incurred in the first stage of the pipeline. Furthermore, 542 these input segmentations are obtained from CT scans and not from US images, as it is needed 543 for other US simulation approaches^{17,18}. Therefore, accurate CT segmentation models such as 544 TotalSegmentator⁵¹ and Auto3DSeg⁵² might be a promising alternative to automate this step 545 and ease the incorporation of new simulation cases. 546

Notice that our image translation approach was trained and evaluated using images simulated with a single ray-casting approach with a fixed configuration, and with real scans obtained from a single US device. Consequently, it does not generalize to produce images from other probes or devices. However, notice also that our proposed model is general enough to be retrained with images from other sources. Hence, by changing \mathcal{A} and/or \mathcal{R} with sets of artificial and/or real scans generated with a different simulator or US device, respectively, or under different imaging setups, the model would adapt to produce new artificial images for other practical applications.

As with all generative models, another limitation of this study is the lack of a trustworthy 554 automated evaluation metric. The best approach for assessing the performance of US simulation 555 algorithms is to run user tests with US experts, where individual images are analyzed and ranked 556 based on their realism, without knowing their source. However, this becomes impractical for 557 ablation studies, which require a substantial number of comparisons across multiple models and 558 images. Furthermore, it is affected by subjective factors such as the level of experience of the 559 human graders and their fatigue while performing the assessment. Although we conducted a user 560 study with participants who are professionals specializing in abdominal ultrasound to add reliability 561 to our findings, we acknowledge that a larger sample size could provide additional insights into the 562 generalizability of the results. While the sample size is small, it enabled us to obtain meaningful 563 insights that allowed to complement the validation of our approach. Furthermore, it is important 564 to notice that most user studies in US simulation research use even smaller sample sizes (between 565 4 and $6^{17,18,42,53}$) than the one presented in this work (16). To the best of our knowledge, only 566 one study used more experts for the validation than ours 34 . 567

Measuring the quality of results obtained using unpaired generative models is inherently 568 complex since it cannot be done using standardized metrics, such as SSIM and SNR, which 569 require ground truth matching between real and artificial scans¹⁶. In an effort to provide a 570 quantitative evaluation, we employed several metrics commonly used in the context of US sim-571 ulation. These metrics enable the assessment of different aspects of the generated images from 572 multiple complementary perspectives^{16,17,34}. FID and KID allow to evaluate scans at a macro 573 level, characterizing their texture patterns using filters from a pre-trained convolutional neural 574 network. The χ^2 distance in particular is commonly employed for tissue characterization in paired 575 image patches¹⁶. Alternatively, we used it to characterize intensities using segmentation masks 576 to extract organ histograms (Section IV.E.). To complement this analysis, we also compared the 577 cumulative distribution of echogenicities of each organ of interest (Figure 5.B). For homogeneous 578 structures, such as the liver and the gallbladder, the histograms from SG-CycleGAN outputs were 579 more alike to the ones computed from real scans. However, some notorious differences persisted 580 in the kidney. The kidney has a complex internal anatomical structure (renal pelvis, renal cortex, 581

etc.) which might be the cause of these differences. Considering the presence or absence of these structures separately, might be a way to account for these differences.

The fact that US images obtained in DICOM format are, by default, JPEG compressed, is 584 a drawback. JPEG is a lossy compression format that introduces artifacts in the images. As our 585 models were trained to produce artificial scans that match the target distribution, it is expected 586 for them to also feature these artifacts. This does not compromise our proposed model nor its 587 evaluation, since they are compared to images presenting the same artifacts. In a more general 588 context, image data used in the training of the proposed model should be consistent in the 589 characteristics of the data where it will be applied. Failing to do so, might notoriously affect the 590 results. 591

⁵⁹² VII. Conclusions

In this paper we introduce a series of contributions to improve anatomical consistency and re-593 duce artifacts in hybrid abdominal US simulators than combine ray-casting-based methods and 594 CycleGANs. Our approach preserves anatomical structures and reduces hallucinations both inside 595 organs and outside the FoV. We demonstrated that the weakly supervised segmentation-guided 596 loss prevents significant alterations in anatomical areas, by penalizing differences in predicted 597 masks obtained from a pre-trained Unet before and after the cycle consistency term. Addition-598 ally, training with images in polar coordinates constrains the FoV, enabling the model to focus 599 on relevant content within non-blank areas. Our model demonstrated to be able to generate 600 synthetic US images with fewer unrealistic artifacts, scattering patterns that are compatible with 601 the acquisition probe's azimuthal angle, and a consistent FoV, closely resembling real scans. This 602 approach enhances the realism of simulators, aiding in training and localization of abdominal or-603 gans. We believe future research can further improve these results by incorporating more organs 604 and simulating abnormalities such as liver tumors or cysts, benefiting training for clinicians. Ad-605 ditionally, eliminating the ray-casting stage by training paired models directly from segmentation 606 masks could lead to end-to-end trainable simulators. We encourage researchers to explore these 607 promising directions to advance this field. 608

609

This work is funded by ANPCyT PICTs 2020-0045 and PIP GI 2021-2023-11220200102472CO. A Kaggle Open Data Research Grant also supported us with a financial grant to purchase the GPU used for this research. We thank all the expert radiologists who participated in the user study.

614 Data availability statement

 $_{615}$ The data that support this study was made publicly available by the authors as a Kaggle dataset 5

616 Conflict of interest statement

⁶¹⁷ The authors declare that they have no conflict of interest.

619 References

- ⁶²⁰ ¹ S. Vitale, J. I. Orlando, E. Iarussi, and I. Larrabide, Improving realism in patient-specific ⁶²¹ abdominal ultrasound simulation using CycleGANs, International journal of computer assisted ⁶²² radiology and surgery 15, 183–192 (2020).
- ⁶²³ ² T. Kameda and N. Taniguchi, Overview of point-of-care abdominal ultrasound in emergency ⁶²⁴ and critical care, Journal of Intensive Care 4, 53 (2016).
- ³ J. Urbina, S. M. Monks, and S. B. Crawford, Simulation in Ultrasound Training for Obstetrics and Gynecology: A Literature Review, Simulation 15 (2021).
- ⁶²⁷ ⁴ V. A. Dinh, J. Y. Fu, S. Lu, A. Chiem, J. C. Fox, and M. Blaivas, Integration of ultra-⁶²⁸ sound in medical education at United States medical schools: a national survey of directors' ⁶²⁹ experiences, Journal of ultrasound in medicine **35**, 413–419 (2016).
- ⁵ M. Østergaard, C. Ewertsen, L. Konge, E. Albrecht-Beste, and M. B. Nielsen, Simulation-⁶³¹ based abdominal ultrasound training–a systematic review, Ultraschall in der Medizin-⁶³² European Journal of Ultrasound **37**, 253–261 (2016).

⁵https://www.kaggle.com/datasets/ignaciorlando/ussimandsegm

⁶³³ D. J. Canty, J. A. Hayes, D. A. Story, and C. F. Royse, Ultrasound simulator-assisted teaching
 ⁶³⁴ of cardiac anatomy to preclinical anatomy students: A pilot randomized trial of a three-hour
 ⁶³⁵ learning exposure, Anatomical sciences education 8, 21–30 (2015).

⁶³⁶ ⁷ B. P. Dromey, D. M. Peebles, and D. V. Stoyanov, A systematic review and meta-analysis
 ⁶³⁷ of the use of high-fidelity simulation in obstetric ultrasound, Simulation in Healthcare 16,
 ⁶³⁸ 52–59 (2021).

- ⁸ M. Donnez, F.-X. Carton, F. Le Lann, E. De Schlichting, and M. Chabanas, Realistic
 ⁶⁴⁰ synthesis of brain tumor resection ultrasound images with a generative adversarial network,
 ⁶⁴¹ in *Medical Imaging 2021: Image-Guided Procedures, Robotic Interventions, and Modeling*,
 ⁶⁴² volume 11598, pages 637–642, SPIE, 2021.
- ⁹ L. Bargsten and A. Schlaefer, SpeckleGAN: a generative adversarial network with an adaptive
 speckle layer to augment limited training data for ultrasound image processing, International
 journal of computer assisted radiology and surgery 15, 1427–1436 (2020).
- ⁶⁴⁶ ¹⁰ R. Shams, R. Hartley, and N. Navab, Real-time simulation of medical ultrasound from CT
 ⁶⁴⁷ images, in *International Conference on Medical Image Computing and Computer-Assisted* ⁶⁴⁸ *Intervention*, pages 734–741, Springer, 2008.
- ⁶⁴⁹¹¹ B. Burger, S. Bettinghausen, M. Radle, and J. Hesser, Real-time GPU-based ultrasound
 ⁶⁵⁰ simulation using deformable mesh models, IEEE transactions on medical imaging 32, 609–
 ⁶⁵¹ 618 (2012).
- ⁶⁵² ¹² O. Mattausch and O. Goksel, Monte-carlo ray-tracing for realistic interactive ultrasound
 ⁶⁵³ simulation, in *Proceedings of the Eurographics Workshop on Visual Computing for Biology* ⁶⁵⁴ and Medicine, pages 173–181, 2016.
- D. Tomar, L. Zhang, T. Portenier, and O. Goksel, Content-preserving unpaired translation
 from simulated to realistic ultrasound images, in *Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part VIII 24*, pages 659–669, Springer, 2021.
- L. Ruthotto and E. Haber, An introduction to deep generative modeling, GAMM-Mitteilungen
 44, e202100008 (2021).

- ¹⁵ F. Tom and D. Sheet, Simulating patho-realistic ultrasound images using deep generative
 networks with adversarial learning, in *2018 IEEE 15th international symposium on biomedical imaging (ISBI 2018)*, pages 1174–1177, IEEE, 2018.
- L. Zhang, T. Portenier, and O. Goksel, Learning ultrasound rendering from cross-sectional
 model slices for simulated training, International Journal of Computer Assisted Radiology
 and Surgery 16, 721–730 (2021).
- ¹⁷ J. Liang, X. Yang, Y. Huang, H. Li, S. He, X. Hu, Z. Chen, W. Xue, J. Cheng, and D. Ni, Sketch guided and progressive growing GAN for realistic and editable ultrasound image synthesis, Medical Image Analysis **79**, 102461 (2022).
- G. Pigeau, L. Elbatarny, V. Wu, A. Schonewille, G. Fichtinger, and T. Ungi, Ultrasound
 image simulation with generative adversarial network, in *Medical Imaging 2020: Image- Guided Procedures, Robotic Interventions, and Modeling*, volume 11315, pages 54–60, SPIE,
 2020.
- ⁶⁷⁴ ¹⁹ N. J. Cronin, T. Finni, and O. Seynnes, Using deep learning to generate synthetic B-mode
 ⁶⁷⁵ musculoskeletal ultrasound images, Computer methods and programs in biomedicine 196,
 ⁶⁷⁶ 105583 (2020).
- ⁶⁷⁷²⁰ J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, Unpaired image-to-image translation using ⁶⁷⁸ cycle-consistent adversarial networks, in *Proceedings of the IEEE international conference* ⁶⁷⁹ *on computer vision*, pages 2223–2232, 2017.
- J. P. Cohen, M. Luck, and S. Honari, Distribution matching losses can hallucinate features in medical image translation, in *Medical Image Computing and Computer Assisted Intervention–MICCAI 2018: 21st International Conference, Granada, Spain, September 16-20, 2018, Proceedings, Part I*, pages 529–536, Springer, 2018.

P. Rubi, E. F. Vera, J. Larrabide, M. Calvo, J. D'Amato, and I. Larrabide, Comparison of
 real-time ultrasound simulation models using abdominal CT images, in *12th international symposium on medical information processing and analysis*, volume 10160, pages 55–63,
 SPIE, 2017.

- O. Ronneberger, P. Fischer, and T. Brox, Unet: Convolutional networks for biomedical
 image segmentation, in *International Conference on Medical image computing and computer- assisted intervention*, pages 234–241, Springer, 2015.
- ⁶⁹¹ ²⁴ X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, and S. Paul Smolley, Least squares generative
 ⁶⁹² adversarial networks, in *Proceedings of the IEEE international conference on computer vision*,
 ⁶⁹³ pages 2794–2802, 2017.
- O. Jimenez-del Toro et al., Cloud-based evaluation of anatomical structure segmentation
 and landmark detection algorithms: VISCERAL anatomy benchmarks, IEEE transactions on
 medical imaging 35, 2459–2475 (2016).
- ⁶⁹⁷²⁶ K. He, X. Zhang, S. Ren, and J. Sun, Deep residual learning for image recognition, in
 ⁶⁹⁸ *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–
 ⁶⁹⁹ 778, 2016.
- ²⁷ S. Dangi and C. Linte, DenseUNet-K: A simplified Densely Connected Fully Convolutional
 ⁷⁰⁰ Network for Image-to-Image Translation, (2019).
- ⁷⁰² ²⁸ X. Sun, H. Li, and W.-N. Lee, Constrained CycleGAN for effective generation of ultrasound
 ⁷⁰³ sector images of improved spatial resolution, Physics in Medicine and Biology (2023).
- P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, Image-to-image translation with conditional
 adversarial networks, in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017.
- ³⁰ D. Ulyanov, A. Vedaldi, and V. Lempitsky, Instance normalization: The missing ingredient
 ⁷⁰⁸ for fast stylization, arXiv preprint arXiv:1607.08022 (2016).
- ⁷⁰⁹ ³¹ D. Ulyanov, A. Vedaldi, and V. Lempitsky, Improved texture networks: Maximizing quality
 ⁷¹⁰ and diversity in feed-forward stylization and texture synthesis, in *Proceedings of the IEEE* ⁷¹¹ conference on computer vision and pattern recognition, pages 6924–6932, 2017.
- ³² D. P. Kingma and J. Ba, Adam: A method for stochastic optimization, arXiv preprint
 arXiv:1412.6980 (2014).

- T. Park, A. A. Efros, R. Zhang, and J.-Y. Zhu, Contrastive learning for unpaired image-to image translation, in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IX 16*, pages 319–345, Springer, 2020.
- ³⁴ D. Tomar, L. Zhang, T. Portenier, and O. Goksel, Content-preserving unpaired translation
 ⁷¹⁸ from simulated to realistic ultrasound images, in *International Conference on Medical Image* ⁷¹⁹ *Computing and Computer-Assisted Intervention*, pages 659–669, Springer, 2021.
- ³⁵ D. Torbunov, Y. Huang, H.-H. Tseng, H. Yu, J. Huang, S. Yoo, M. Lin, B. Viren, and Y. Ren,
 Rethinking CycleGAN: Improving Quality of GANs for Unpaired Image-to-Image Translation,
 arXiv preprint arXiv:2303.16280 (2023).
- ³⁶ B. Kim, G. Kwon, K. Kim, and J. C. Ye, Unpaired Image-to-Image Translation via Neural
 ⁷²⁴ Schr\" odinger Bridge, arXiv preprint arXiv:2305.15086 (2023).
- X. Ma, N. Anantrasirichai, S. Bolomytis, and A. Achim, PMT: Partial-Modality Translation
 Based on Diffusion Models for Prostate Magnetic Resonance and Ultrasound Image Registra tion, in *Annual Conference on Medical Image Understanding and Analysis*, pages 285–297,
 Springer, 2024.
- ⁷²⁹ ³⁸ H. Alqahtani, M. Kavakli-Thorne, G. Kumar, and F. SBSSTC, An analysis of evaluation
 ⁷³⁰ metrics of gans, in *International Conference on Information Technology and Applications* ⁷³¹ (*ICITA*), volume 7, 2019.
- ³⁹ A. Borji, Pros and cons of GAN evaluation measures: New developments, Computer Vision
 ³³ and Image Understanding 215, 103329 (2022).
- ⁴⁰ M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, GANs Trained by a
 Two Time-Scale Update Rule Converge to a Local Nash Equilibrium, in *Advances in Neural* ⁷³⁶ *Information Processing Systems*, volume 30, Curran Associates, Inc., 2017.
- ⁴¹ M. Bińkowski, D. J. Sutherland, M. Arbel, and A. Gretton, Demystifying mmd gans, arXiv
 preprint arXiv:1801.01401 (2018).
- ⁴² J. Liang et al., Weakly-supervised high-fidelity ultrasound video synthesis with feature de ⁷⁴⁰ coupling, in *International Conference on Medical Image Computing and Computer-Assisted* ⁷⁴¹ *Intervention*, pages 310–319, Springer, 2022.

43 742 743 744	C. Szegedy, V. Vanhoucke, S. loffe, J. Shlens, and Z. Wojna, Rethinking the inception architecture for computer vision, in <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 2818–2826, 2016.
44 745 746	A. Obukhov, M. Seitzer, P. Wu, S. Zhydenko, J. Kyl, and E. Lin, High-fidelity performance metrics for generative models in PyTorch, 2020.
747 ⁴⁵ 748	C. Bonferroni, Statistic theory of classes and calculation of probabilities, Volume in Honor of Riccardo della Volta. Florence: University of Florence , 1–62 (1937).
749 46	J. Cohen, Statistical power analysis for the behavioral sciences, routledge, 2013.
47 751 752	G. E. Mailloux, M. Bertrand, R. Stampfler, and S. Ethier, Local histogram information content of ultrasound B-mode echographic texture, Ultrasound in medicine & biology 11, 743–750 (1985).
 48 754 755 756 	D. China, F. Tom, S. Nandamuri, A. Kar, M. Srinivasan, P. Mitra, and D. Sheet, Ultra- compression: framework for high density compression of ultrasound volumes using physics modeling deep neural networks, in <i>2019 IEEE 16th International Symposium on Biomedical</i> <i>Imaging (ISBI 2019)</i> , pages 798–801, IEEE, 2019.
49 757 758	A. K. Tripathi, S. Mukhopadhyay, and A. K. Dhara, Performance metrics for image contrast, in <i>2011 International Conference on Image Information Processing</i> , pages 1–4, IEEE, 2011.
759 ⁵⁰ 760	J. R. De Leeuw, jsPsych: A JavaScript library for creating behavioral experiments in a Web browser, Behavior research methods 47, 1–12 (2015).
761 51 762	J. Wasserthal et al., TotalSegmentator: robust segmentation of 104 anatomic structures in CT images, Radiology: Artificial Intelligence 5 (2023).
 52 764 765 	A. Myronenko, D. Yang, Y. He, and D. Xu, Automated 3D Segmentation of Kidneys and Tumors in MICCAI KiTS 2023 Challenge, in <i>International Challenge on Kidney and Kidney</i> <i>Tumor Segmentation</i> , pages 1–7, Springer, 2023.
766 ⁵³ 767 768	L. Chen, H. Liao, W. Kong, D. Zhang, and F. Chen, Anatomy preserving GAN for realistic simulation of intraoperative liver ultrasound images, Computer Methods and Programs in Biomedicine 240 , 107642 (2023).