



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

Otras autorías: Vitale, Santiago; Orlando, José Ignacio; Díaz, Alejandro; Larrabide, Ignacio

Autoría ditelliana: Iarussi, Emmanuel

Año: 2025

Nota: Este document es una versión previa. La versión final está publicada en Medical Physics (eISSN: 2473-4209)

¿Cómo citar la versión final de este trabajo?

Vitale S, Orlando JI, Iarussi E, Díaz A, Larrabide I. Improving realism in abdominal ultrasound simulation combining a segmentation-guided loss and polar coordinates training. Med Phys. 2025; 52: 4540–4556. <https://doi.org/10.1002/mp.17801>

El presente documento se encuentra alojado en el Repositorio Digital de la **Universidad Torcuato Di Tella**, para su preservación, archivo y difusión

Dirección: <https://repositorio.utdt.edu/handle/20.500.13098/13822>

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 Universitario, Paraje Arroyo Seco. Email: santiago.vitale@pladema.exa.unicen.edu.ar

Abstract

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

¹National Scientific and Technical Research Council (CONICET), Buenos Aires, Argentina

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

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

42 I. Introduction

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

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

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

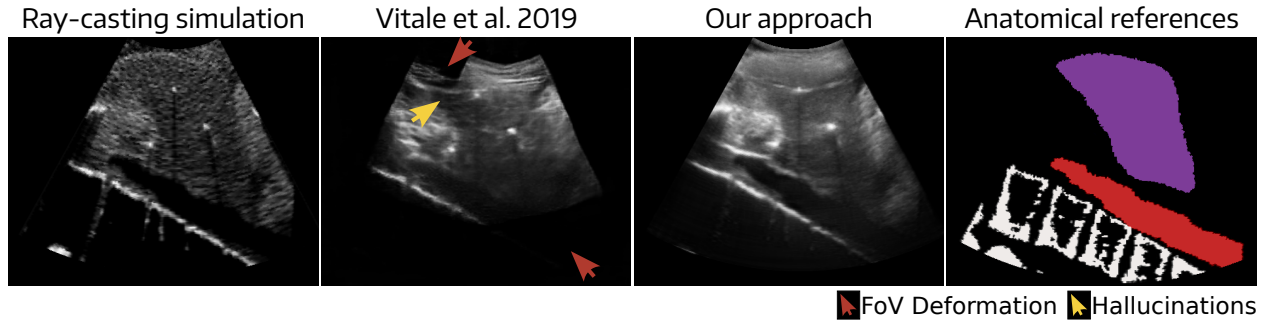


Figure 1: Examples of different artificial US scans obtained with a ray-casting model, our previous approach based on a standard CycleGAN model¹, and our improved method using a segmentation-guided loss and polar coordinates. Anatomical masks are provided as reference.

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

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

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

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

90 new individuals

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

93 II. Related work

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

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

$$\begin{aligned} \mathcal{L}_{\text{GAN}}(G_{\mathcal{A} \rightarrow \mathcal{R}}, D_{\mathcal{R}}, \mathcal{A}, \mathcal{R}) &= \mathbb{E}_{r \sim p_{\text{data}}(r)} [\log(D_{\mathcal{R}}(r) - 1)^2] \\ &\quad + \mathbb{E}_{a \sim p_{\text{data}}(a)} [\log(D_{\mathcal{R}}(G_{\mathcal{A} \rightarrow \mathcal{R}}(a))^2], \\ \mathcal{L}_{\text{GAN}}(G_{\mathcal{R} \rightarrow \mathcal{A}}, D_{\mathcal{A}}, \mathcal{A}, \mathcal{R}) &= \mathbb{E}_{a \sim p_{\text{data}}(a)} [\log(D_{\mathcal{A}}(a) - 1)^2] \\ &\quad + \mathbb{E}_{r \sim p_{\text{data}}(r)} [\log(D_{\mathcal{A}}(G_{\mathcal{R} \rightarrow \mathcal{A}}(r))^2], \end{aligned} \tag{1}$$

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

108 To allow unpaired image-to-image translation, the training scheme incorporates an additional
109 cycle-consistency loss \mathcal{L}_{cyc} . This term enforce that translations produced by one generator are
110 reversible and retain the original domain's characteristics (Step 1, Figure 2). Formally, a forward
111 cycle translates an image $a \in \mathcal{A}$ previously translated to domain \mathcal{R} back to \mathcal{A} (that is, $a \rightarrow$
112 $G_{\mathcal{A} \rightarrow \mathcal{R}}(a) \rightarrow G_{\mathcal{R} \rightarrow \mathcal{A}}(G_{\mathcal{A} \rightarrow \mathcal{R}}(a)) \approx a$). Similarly, a reverse cycle ensures an image $r \in \mathcal{R}$

113 translated to domain \mathcal{A} is brought back to \mathcal{R} (by doing $r \rightarrow G_{\mathcal{R} \rightarrow \mathcal{A}}(r) \rightarrow G_{\mathcal{A} \rightarrow \mathcal{R}}(G_{\mathcal{R} \rightarrow \mathcal{A}}(r)) \approx$
 114 r). \mathcal{L}_{cyc} can then be defined as the sum of two losses:

$$\begin{aligned} \mathcal{L}_{\text{cyc}}(G_{\mathcal{A} \rightarrow \mathcal{R}}, G_{\mathcal{R} \rightarrow \mathcal{A}}) = & \mathbb{E}_{r \sim p_{\text{data}}(r)} [\|G_{\mathcal{A} \rightarrow \mathcal{R}}(G_{\mathcal{R} \rightarrow \mathcal{A}}(r)) - r\|_1] \\ & + \mathbb{E}_{a \sim p_{\text{data}}(a)} [\|G_{\mathcal{R} \rightarrow \mathcal{A}}(G_{\mathcal{A} \rightarrow \mathcal{R}}(a)) - a\|_1]. \end{aligned} \quad (2)$$

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

$$\begin{aligned} \mathcal{L}_{\text{idt}}(G_{\mathcal{A} \rightarrow \mathcal{R}}, G_{\mathcal{R} \rightarrow \mathcal{A}}) = & \mathbb{E}_{a \sim p_{\text{data}}(a)} [\|G_{\mathcal{R} \rightarrow \mathcal{A}}(a) - a\|_1] \\ & + \mathbb{E}_{r \sim p_{\text{data}}(r)} [\|G_{\mathcal{A} \rightarrow \mathcal{R}}(r) - r\|_1]. \end{aligned} \quad (3)$$

117 III. Methods

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

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

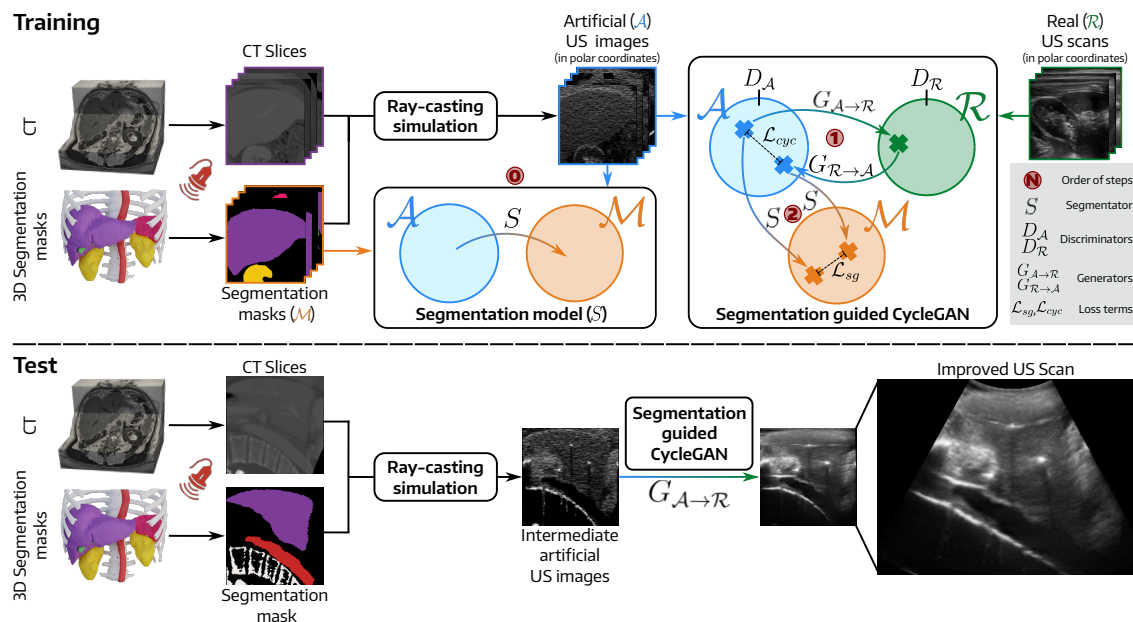


Figure 2: Schematic representation of the training (top) and testing (bottom) phases of our proposed approach for improving abdominal US simulation using a novel anatomically consistent image-to-image translation model.

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

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

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

$$\mathcal{L}_{\text{sg}}(G_{\mathcal{A} \rightarrow \mathcal{R}}, G_{\mathcal{R} \rightarrow \mathcal{A}}, S) = - \sum S(a) \log (S(G_{\mathcal{R} \rightarrow \mathcal{A}}(G_{\mathcal{A} \rightarrow \mathcal{R}}(a)))), \quad (4)$$

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

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

$$\begin{aligned} \mathcal{L}(G_{\mathcal{A} \rightarrow \mathcal{R}}, G_{\mathcal{R} \rightarrow \mathcal{A}}, D_{\mathcal{A}}, D_{\mathcal{R}}, S) &= \mathcal{L}_{\text{GAN}}(G_{\mathcal{A} \rightarrow \mathcal{R}}, D_{\mathcal{R}}, \mathcal{A}, \mathcal{R}) \\ &+ \mathcal{L}_{\text{GAN}}(G_{\mathcal{R} \rightarrow \mathcal{A}}, D_{\mathcal{A}}, \mathcal{R}, \mathcal{A}) \\ &+ \lambda_{\text{cyc}} \cdot \mathcal{L}_{\text{cyc}}(G_{\mathcal{A} \rightarrow \mathcal{R}}, G_{\mathcal{R} \rightarrow \mathcal{A}}) \\ &+ \lambda_{\text{idt}} \cdot \mathcal{L}_{\text{idt}}(G_{\mathcal{A} \rightarrow \mathcal{R}}, G_{\mathcal{R} \rightarrow \mathcal{A}}) \\ &+ \lambda_{\text{sg}} \cdot \mathcal{L}_{\text{sg}}(G_{\mathcal{A} \rightarrow \mathcal{R}}, G_{\mathcal{R} \rightarrow \mathcal{A}}, S), \end{aligned} \quad (5)$$

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

162 Notice that the identity loss and the segmentation-guided loss serve different purposes in the
 163 model. The identity term enforces that each generator maintains features from the target domain
 164 that are already present in the source domain. On the other hand, our segmentation-guided loss
 165 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

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

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

181 IV.A.2. Segmentation masks dataset

182 The anatomical masks used correspond to the cross-sectional slices extracted from the silver
183 corpus segmentations of the 13 CT volumes in Artificial US dataset. The original dataset in-
184 cluded segmentations of the spleen, liver, gallbladder, aorta, and kidneys. To provide additional
185 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

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

203 IV.B. Architectures

204 IV.B.1. Generator architecture

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

213 IV.B.2. Discriminator architecture

214 Following previous studies^{16,17,28}, we employed a 70×70 patchGAN²⁹ as the discriminator. The
215 network consists of four convolutional blocks, each with a 4×4 kernel and a stride of 2. Instance
216 normalization was used instead of batch normalization, as it has been shown to enhance diversity
217 and prevent mode collapse^{30,31}. Each block applies Leaky-ReLU activation, as commonly done
218 in patchGANs²⁹, progressively reducing spatial dimensions while increasing feature maps to 64,
219 128, 256, and 512, respectively. A final 1-filter convolution, followed by a sigmoid activation
220 function, produces the output probability for each patch.

221 IV.B.3. Segmentation model

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

233 IV.C. Model configuration

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

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

249 IV.D. Baselines for comparison

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

261 IV.E. Evaluation metrics & statistical analysis

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

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

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

$$\chi^2(h_A||h_B) = \frac{1}{2} \sum_{l=1..d} \frac{(h_A[l] - h_B[l])^2}{h_A[l] + h_B[l]}, \quad (6)$$

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

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

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

301 IV.F. User study-based evaluation

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

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

316 V. Results

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

322 V.A. Simulation performance

323 V.A.1. Qualitative evaluation

324 Figure 3 presents example simulations generated using the original CycleGAN in Cartesian coord-
325 inates¹, the same model in polar coordinates, and our SG-CycleGAN. The samples correspond
326 to different abdominal windows commonly used in clinical analyses.

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

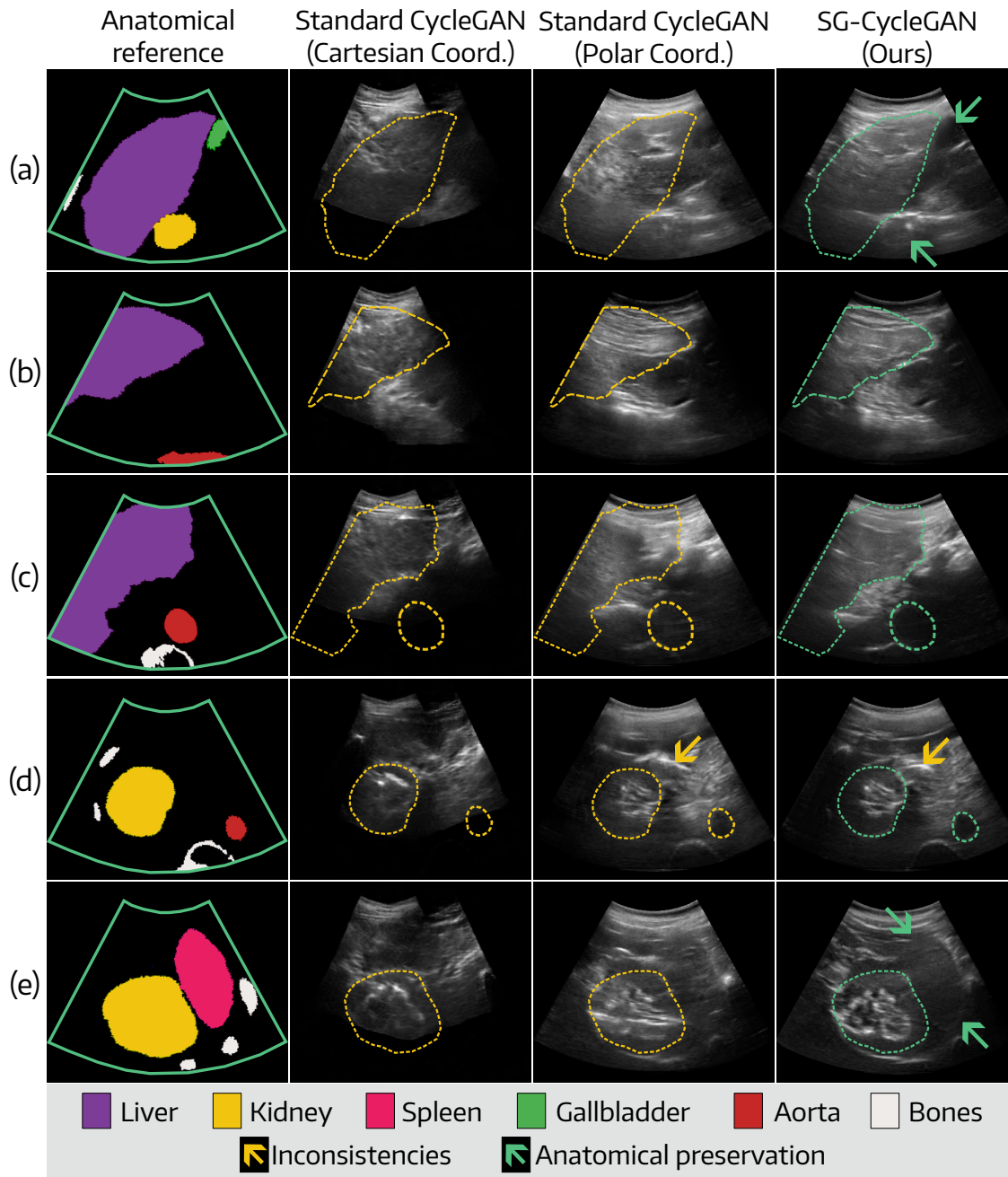


Figure 3: Qualitative results for abdominal US simulation obtained using a standard CycleGAN trained in Cartesian and polar coordinates and our proposed SG-CycleGAN approach. Dotted lines indicate inconsistent organs (yellow) and their improved counterparts (green). From top to bottom: (a) right subcostal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

332 all the organs that are present in the images.

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

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

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

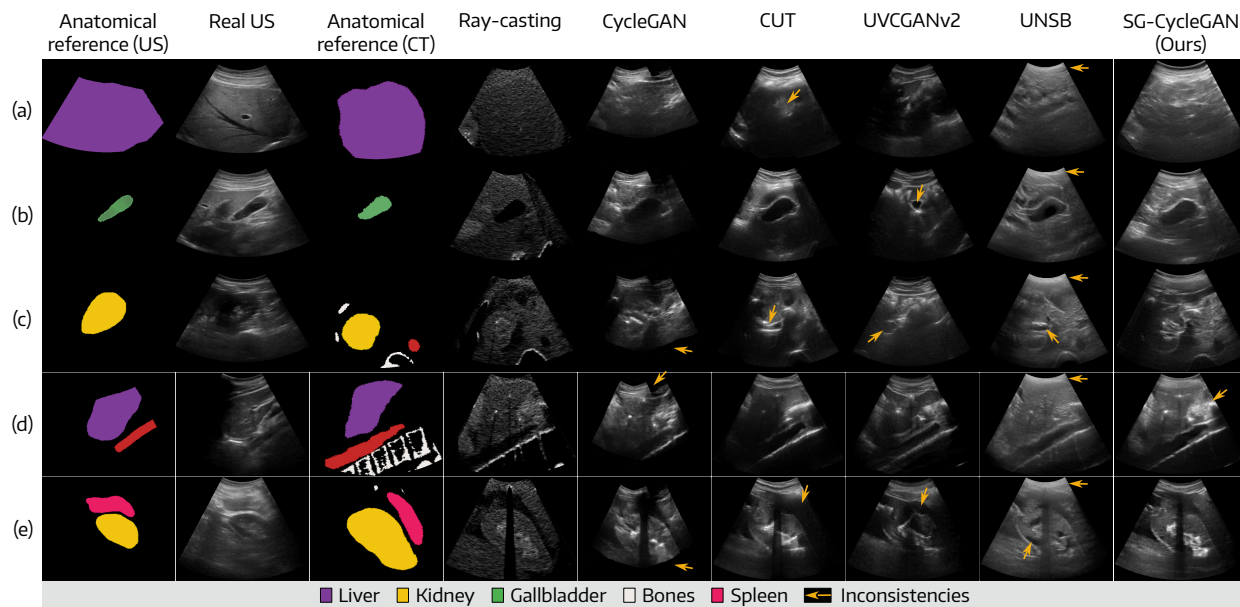


Figure 4: Qualitative examples for each model and their associated segmentations as reference. Yellow arrows indicate inconsistencies.

361 **V.A.2. Quantitative evaluation**

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

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

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

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

393 We also report training and inference time comparisons in the supplementary material.
394 SG-CycleGAN increased training time from 95 seconds (standard CycleGAN) to 127 seconds
395 per epoch, similar to CUT and notably faster than UVCGANv2 and UNSB. For inference, SG-
396 CycleGAN and CycleGAN were the fastest at 0.0813 seconds per scan, while other models required
397 2–3 times longer.

398 V.B. Ablation analysis

399 V.B.1. Quantitative evaluation

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

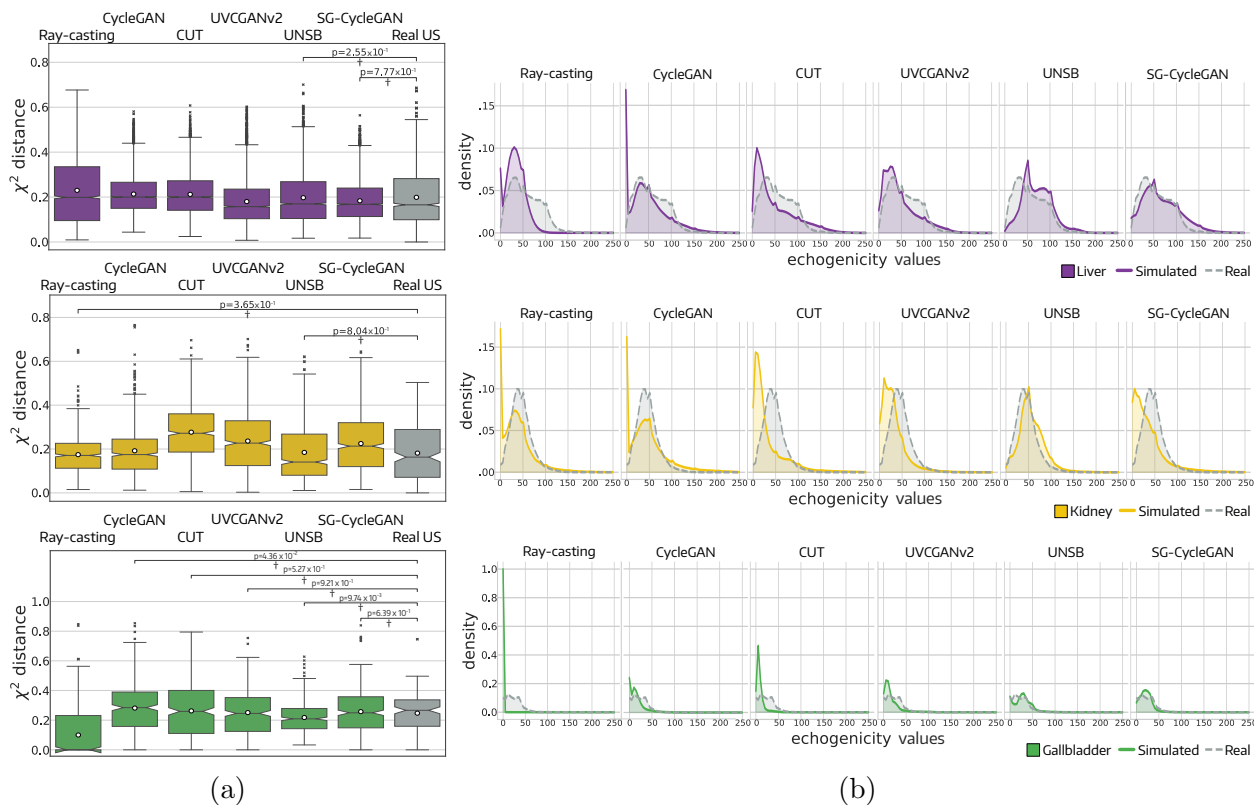


Figure 5: Organ-wise quantitative evaluation. (A) Box plots illustrating the distribution of pairwise χ^2 distances between pairs of simulated and real US images for each organ of interest (colored), and between pairs of real US images (gray). p-values are included for all comparison where no statistical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

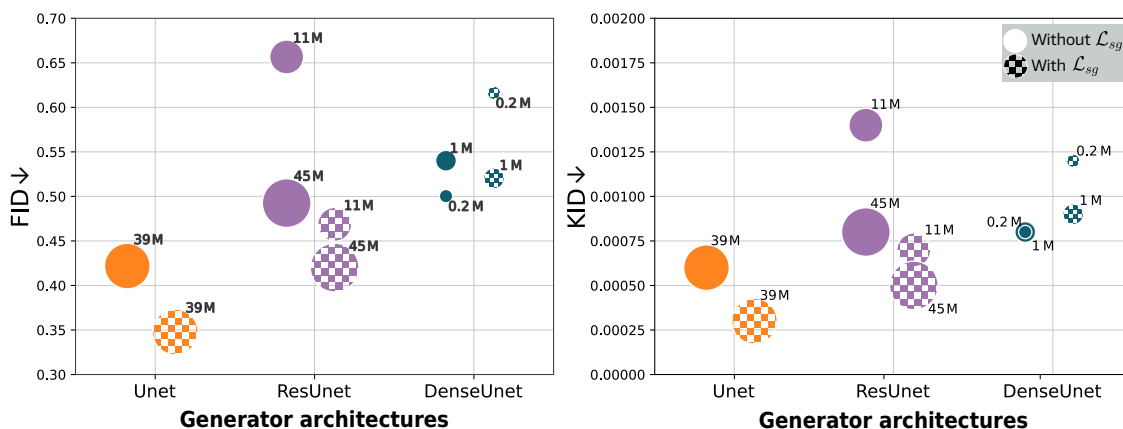


Figure 6: FID and KID results for different architectures of generator models. Each network was trained with (right) and without (left) our proposed loss term. The bubble size is proportional to the number of parameters of each model, indicated in millions (M) on top of each one.

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 (\dagger) 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	FID \downarrow	KID $\downarrow(\times 10^{-3})$	χ^2 [95% CI]		
	[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] _{†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] _{†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] _{†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] _{†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] _{†0.01}	0.18 [0.02 - 0.50] _{†0.03}	0.22 [0.05 - 0.45] _{†0.23}
SG-CycleGAN (ours)	0.33 [0.32 - 0.35]	0.28 [0.25 - 0.31]	0.18 [0.05 - 0.40] _{†0.13}	0.22 [0.03 - 0.48] _{0.33}	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 $\mathcal{R} \mathcal{A}$	213 213	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

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

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 (\dagger) 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 loss	Coordinate space	FID \downarrow	KID \downarrow ($\times 10^{-3}$)	χ^2		
			[95% CI]	[95% CI]	Liver	Kidney	Gallbladder
CG	Vanilla	C	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	P	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	P	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	P	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] _{0.34}	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 ($\mathcal{R} \mathcal{A}$)			213 293	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

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

428 V.B.3. Qualitative effect of the generator architecture

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

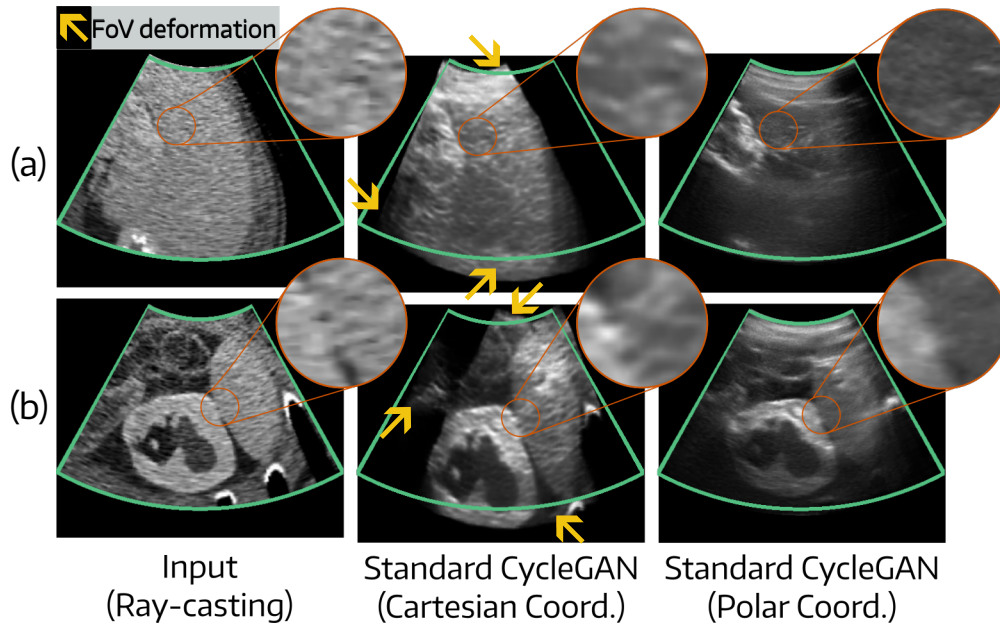


Figure 7: Comparison of simulated images with CycleGANs trained on different coordinate systems. Green boundaries indicate the original FoV.

439 V.B.4. User study

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

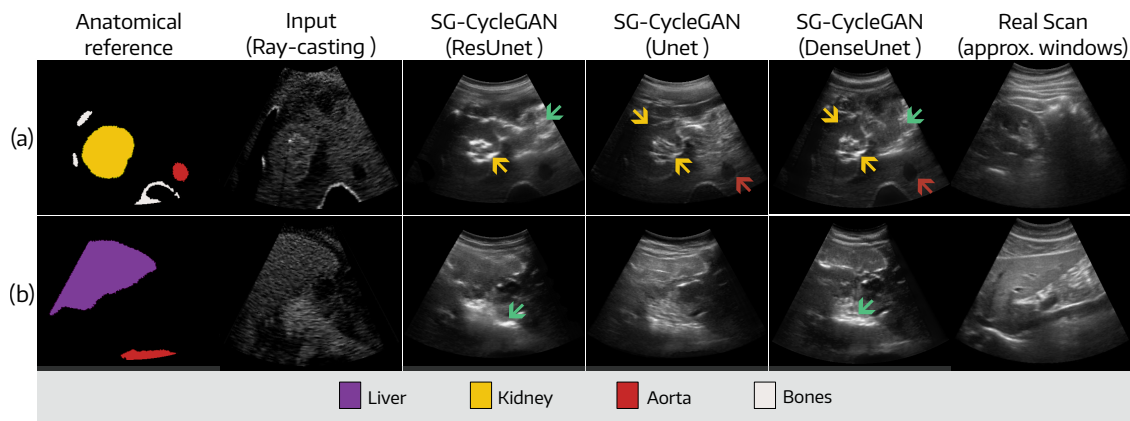


Figure 8: Comparison of simulation results obtained using an SG-CycleGAN with Unet, ResUnet, or DenseUnet based generator.

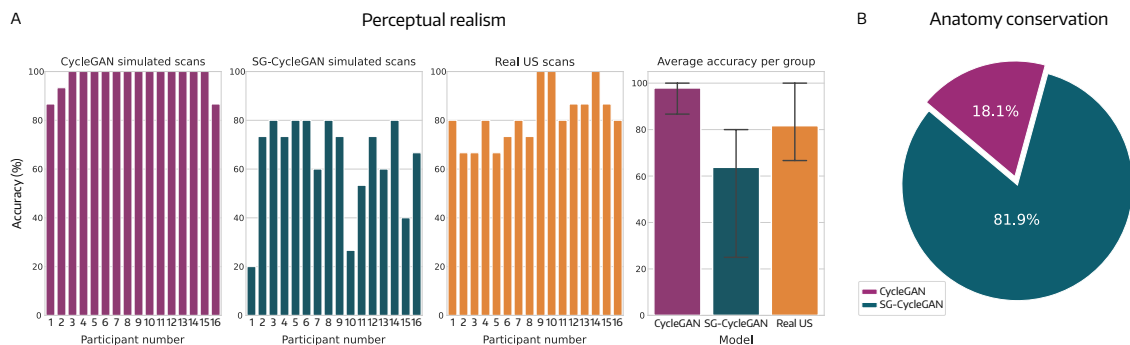


Figure 9: User study results. A) Classification accuracy for each simulation model and real scans as a bar per participant. Additionally, a bar plot with average accuracy per method. B) Pie chart comparing responses about which generative model performs better in terms of anatomy conservation.

451 VI. Discussion

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

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

458 To alleviate this issue, we proposed a segmentation-guided loss, penalizing segmentation
 459 mismatches before and after completing the cycle. This term prevents the generator $\mathcal{G}_{A \rightarrow R}$ to
 460 introduce artifacts that cannot be removed through the reversed cycle $\mathcal{G}_{R \rightarrow A}$, without any extra

461 annotation. The anatomical labels from ray-casting simulations suffice for training. As seen in
462 Figure 3 (green lines), our approach produces well-defined organ interfaces and homogeneous
463 speckle noise patterns. Compared to existing methods (Figure 4), our loss function preserves
464 anatomical structures while preventing hallucinated patterns within them.

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

475 VI.B. Impact of training in polar coordinates

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

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

490 VI.C. Influence of the generator architecture

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

500 VI.D. Advantages of SG-CycleGAN

501 Integrating all our proposed modifications into the standard CycleGAN framework resulted in a
502 robust generative model that outperforms several state-of-the-art approaches in realism. We
503 compared SG-CycleGAN against recent deep learning models, including Vision Transformers
504 (UVCANv2) and conditional diffusion models (UNSB). As shown in Table 1, these methods
505 reduced FID and KID scores relative to the ray-casting model, with Vision Transformers achiev-
506 ing the largest improvement. However, SG-CycleGAN achieved the lowest FID and KID values
507 ($p = 0.33 \times 10^{-3}$ and $p = 0.25 \times 10^{-3}$, respectively), with a very large effect size (Cohen’s
508 $d > 9.20$). Our model also closely matches real ultrasound (US) echogenicity distributions. As
509 shown in Table 2, χ^2 tests indicate no statistically significant differences in liver and gallbladder
510 echogenicities between SG-CycleGAN-generated images and real scans ($p > 0.008$). The effect
511 size is minimal (Cohen’s $d = 0.07$ for the gallbladder and $d = 0.13$ for the liver), suggesting
512 that our model generates tissue echogenicities within the natural variability of real US images.
513 While UNSB achieves a slightly better match for the liver ($d = 0.01$), our approach still per-
514 forms competitively, as showed in Figure 5 (B). From a qualitative perspective, SG-CycleGAN
515 produces more realistic scans. If the generated images were easily distinguishable from real ones,
516 expert classification accuracy would approach 100%. While this was true for standard CycleGAN,
517 experts misclassified 36% of SG-CycleGAN images as real (Figure 9). This suggests that our
518 model generates anatomically consistent and realistic US scans, making it a promising tool for
519 improving ultrasound training applications.

520 VI.E. Limitations

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

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

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

547 Notice that our image translation approach was trained and evaluated using images simulated
548 with a single ray-casting approach with a fixed configuration, and with real scans obtained from
549 a single US device. Consequently, it does not generalize to produce images from other probes or

550 devices. However, notice also that our proposed model is general enough to be retrained with
551 images from other sources. Hence, by changing \mathcal{A} and/or \mathcal{R} with sets of artificial and/or real
552 scans generated with a different simulator or US device, respectively, or under different imaging
553 setups, the model would adapt to produce new artificial images for other practical applications.

554 As with all generative models, another limitation of this study is the lack of a trustworthy
555 automated evaluation metric. The best approach for assessing the performance of US simulation
556 algorithms is to run user tests with US experts, where individual images are analyzed and ranked
557 based on their realism, without knowing their source. However, this becomes impractical for
558 ablation studies, which require a substantial number of comparisons across multiple models and
559 images. Furthermore, it is affected by subjective factors such as the level of experience of the
560 human graders and their fatigue while performing the assessment. Although we conducted a user
561 study with participants who are professionals specializing in abdominal ultrasound to add reliability
562 to our findings, we acknowledge that a larger sample size could provide additional insights into the
563 generalizability of the results. While the sample size is small, it enabled us to obtain meaningful
564 insights that allowed to complement the validation of our approach. Furthermore, it is important
565 to notice that most user studies in US simulation research use even smaller sample sizes (between
566 4 and 6^{17,18,42,53}) than the one presented in this work (16). To the best of our knowledge, only
567 one study used more experts for the validation than ours³⁴.

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

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

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

592 VII. Conclusions

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

609 Acknowledgments

610 This work is funded by ANPCyT PICTs 2020-0045 and PIP GI 2021-2023-11220200102472CO.
611 A Kaggle Open Data Research Grant also supported us with a financial grant to purchase the
612 GPU used for this research. We thank all the expert radiologists who participated in the user
613 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).
- 625 ³ J. Urbina, S. M. Monks, and S. B. Crawford, Simulation in Ultrasound Training for Obstetrics
626 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).
- 630 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).

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

- 633 ⁶ 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).
- 639 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.
- 643 ⁹ L. Bargsten and A. Schlaefler, SpeckleGAN: a generative adversarial network with an adaptive
644 speckle layer to augment limited training data for ultrasound image processing, *International*
645 *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.
- 649 ¹¹ 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.
- 655 ¹³ D. Tomar, L. Zhang, T. Portenier, and O. Goksel, Content-preserving unpaired translation
656 from simulated to realistic ultrasound images, in *Medical Image Computing and Computer*
657 *Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France,*
658 *September 27–October 1, 2021, Proceedings, Part VIII 24*, pages 659–669, Springer, 2021.
- 659 ¹⁴ L. Ruthotto and E. Haber, An introduction to deep generative modeling, *GAMM-Mitteilungen*
660 **44**, e202100008 (2021).
-

- 661 ¹⁵ F. Tom and D. Sheet, Simulating patho-realistic ultrasound images using deep generative
662 networks with adversarial learning, in *2018 IEEE 15th international symposium on biomedical
663 imaging (ISBI 2018)*, pages 1174–1177, IEEE, 2018.
- 664 ¹⁶ L. Zhang, T. Portenier, and O. Goksel, Learning ultrasound rendering from cross-sectional
665 model slices for simulated training, *International Journal of Computer Assisted Radiology
666 and Surgery* **16**, 721–730 (2021).
- 667 ¹⁷ J. Liang, X. Yang, Y. Huang, H. Li, S. He, X. Hu, Z. Chen, W. Xue, J. Cheng, and D. Ni,
668 Sketch guided and progressive growing GAN for realistic and editable ultrasound image syn-
669 thesis, *Medical Image Analysis* **79**, 102461 (2022).
- 670 ¹⁸ G. Pigeau, L. Elbatarny, V. Wu, A. Schonewille, G. Fichtinger, and T. Ungi, Ultrasound
671 image simulation with generative adversarial network, in *Medical Imaging 2020: Image-
672 Guided Procedures, Robotic Interventions, and Modeling*, volume 11315, pages 54–60, SPIE,
673 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).
- 677 ²⁰ 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.
- 680 ²¹ J. P. Cohen, M. Luck, and S. Honari, Distribution matching losses can hallucinate fea-
681 tures in medical image translation, in *Medical Image Computing and Computer Assisted
682 Intervention–MICCAI 2018: 21st International Conference, Granada, Spain, September 16-
683 20, 2018, Proceedings, Part I*, pages 529–536, Springer, 2018.
- 684 ²² P. Rubi, E. F. Vera, J. Larrabide, M. Calvo, J. D’Amato, and I. Larrabide, Comparison of
685 real-time ultrasound simulation models using abdominal CT images, in *12th international
686 symposium on medical information processing and analysis*, volume 10160, pages 55–63,
687 SPIE, 2017.

- 688 ²³ O. Ronneberger, P. Fischer, and T. Brox, U-net: Convolutional networks for biomedical
689 image segmentation, in *International Conference on Medical image computing and computer-*
690 *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.
- 694 ²⁵ O. Jimenez-del Toro et al., Cloud-based evaluation of anatomical structure segmentation
695 and landmark detection algorithms: VISCERAL anatomy benchmarks, *IEEE transactions on*
696 *medical imaging* **35**, 2459–2475 (2016).
- 697 ²⁶ 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.
- 700 ²⁷ S. Dangi and C. Linte, DenseUNet-K: A simplified Densely Connected Fully Convolutional
701 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).
- 704 ²⁹ P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, Image-to-image translation with conditional
705 adversarial networks, in *Proceedings of the IEEE conference on computer vision and pattern*
706 *recognition*, pages 1125–1134, 2017.
- 707 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.
- 712 ³² D. P. Kingma and J. Ba, Adam: A method for stochastic optimization, arXiv preprint
713 arXiv:1412.6980 (2014).
-

- 714 ³³ T. Park, A. A. Efros, R. Zhang, and J.-Y. Zhu, Contrastive learning for unpaired image-to-
715 image translation, in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow,
716 UK, August 23–28, 2020, Proceedings, Part IX 16*, pages 319–345, Springer, 2020.
- 717 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.
- 720 ³⁵ D. Torbunov, Y. Huang, H.-H. Tseng, H. Yu, J. Huang, S. Yoo, M. Lin, B. Viren, and Y. Ren,
721 Rethinking CycleGAN: Improving Quality of GANs for Unpaired Image-to-Image Translation,
722 arXiv preprint arXiv:2303.16280 (2023).
- 723 ³⁶ B. Kim, G. Kwon, K. Kim, and J. C. Ye, Unpaired Image-to-Image Translation via Neural
724 Schrödinger Bridge, arXiv preprint arXiv:2305.15086 (2023).
- 725 ³⁷ X. Ma, N. Anantrasirichai, S. Bolomytis, and A. Achim, PMT: Partial-Modality Translation
726 Based on Diffusion Models for Prostate Magnetic Resonance and Ultrasound Image Registra-
727 tion, in *Annual Conference on Medical Image Understanding and Analysis*, pages 285–297,
728 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.
- 732 ³⁹ A. Borji, Pros and cons of GAN evaluation measures: New developments, *Computer Vision
733 and Image Understanding* **215**, 103329 (2022).
- 734 ⁴⁰ M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, GANs Trained by a
735 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.
- 737 ⁴¹ M. Bińkowski, D. J. Sutherland, M. Arbel, and A. Gretton, Demystifying mmd gans, arXiv
738 preprint arXiv:1801.01401 (2018).
- 739 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.

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