

Self-Supervised Ultrasound to MRI Fetal Brain Image Synthesis

Jianbo Jiao, *Member, IEEE*, Ana I.L. Namburete, Aris T. Papageorgiou, and J. Alison Noble

Abstract—Fetal brain magnetic resonance imaging (MRI) offers exquisite images of the developing brain but is not suitable for second-trimester anomaly screening, for which ultrasound (US) is employed. Although expert sonographers are adept at reading US images, MR images which closely resemble anatomical images are much easier for non-experts to interpret. Thus in this paper we propose to generate MR-like images directly from clinical US images. In medical image analysis such a capability is potentially useful as well, for instance for automatic US-MRI registration and fusion. The proposed model is end-to-end trainable and self-supervised without any external annotations. Specifically, based on an assumption that the US and MRI data share a similar anatomical latent space, we first utilise a network to extract the shared latent features, which are then used for MRI synthesis. Since paired data is unavailable for our study (and rare in practice), pixel-level constraints are infeasible to apply. We instead propose to enforce the distributions to be statistically indistinguishable, by adversarial learning in both the image domain and feature space. To regularise the anatomical structures between US and MRI during synthesis, we further propose an adversarial structural constraint. A new cross-modal attention technique is proposed to utilise non-local spatial information, by encouraging multi-modal knowledge fusion and propagation. We extend the approach to consider the case where 3D auxiliary information (e.g., 3D neighbours and a 3D location index) from volumetric data is also available, and show that this improves image synthesis. The proposed approach is evaluated quantitatively and qualitatively with comparison to real fetal MR images and other approaches to synthesis, demonstrating its feasibility of synthesising realistic MR images.

Index Terms—Self-Supervised, Unpaired, Ultrasound, MRI.

I. INTRODUCTION

OBSTETRIC ultrasound (US) is the most commonly applied clinical imaging technique to monitor fetal development. Clinicians use fetal brain US imaging (fetal neurosonography) to detect abnormalities in the fetal brain and growth restriction. However, fetal neurosonography suffers from acoustic shadows and occlusions caused by the fetal skull. On the other hand, magnetic resonance imaging (MRI) is unaffected by the presence of bone and typically provides good and more complete spatial detail of the full anatomy [1].

Copyright (c) 2019 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

This work was supported in part by the EPSRC (EP/M013774/1 project Seebibyte, EP/R013853/1 project CALOPUS), the ERC (ERC-ADG-2015 694581, project PULSE), and the NIHR Biomedical Research Centre funding scheme.

Jianbo Jiao, Ana I.L. Namburete, and J. Alison Noble are with the Department of Engineering Science, University of Oxford, Oxford, UK (e-mail: jianbo.jiao@eng.ox.ac.uk).

Aris T. Papageorgiou is with the Nuffield Department of Women's & Reproductive Health, University of Oxford, Oxford, UK.

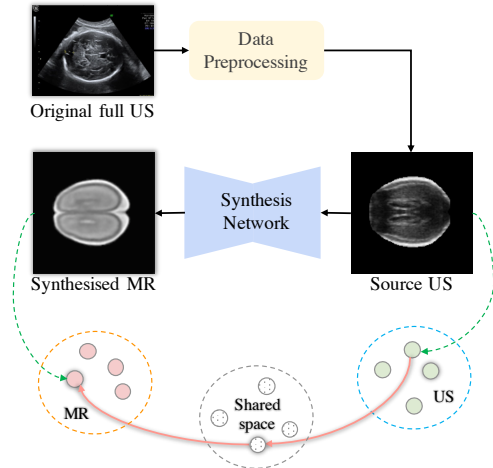


Fig. 1. **Top:** Overview of the proposed US-to-MR synthesis framework. **Bottom:** Assumption of the shared latent space.

Whereas MRI is costly and time-consuming, making it unsuitable for fetal anomaly screening, in the second and third trimesters it is often routinely used for assessment of the fetal brain [2].

Medical image synthesis has received growing interest in recent years. Most prior work to date has focused on the synthesis between MR and CT (computed tomography) images [3]–[5] or of retinal images [6], [7]. Simulation of US images has also been proposed to assist in automatic alignment of US and other modalities [8], [9]. Prior to the deep learning era, medical image synthesis was primarily based on segmentation and atlases. Taking MR-to-CT image synthesis as an example, in segmentation-based methods [10], [11], first an MR image is segmented into different tissue classes, and then the corresponding synthesised CT image is generated by intensity-filling for each class. On the other hand, atlas-based approaches [12], [13] first register the input MR image to an MR atlas by a transformation, followed by applying such transformation to a CT atlas to synthesise the corresponding CT image. Notwithstanding, the above approaches rely heavily on the segmentation and atlas quality, implying low-quality would directly lead to a poor synthesis. Methods based on convolutional neural networks (CNNs) have demonstrated promising performance for medical image synthesis in recent literature. For instance, given a large number of paired MR-CT data, some proposed methods [3], [4], [14] learn a mapping directly from MR to CT through a CNN architecture design. To alleviate the paired data restriction, other methods [5], [15] have converted the image synthesis

problem to image-to-image translation by a recently proposed CycleGAN architecture [16]. Even though the training data is not necessarily perfectly registered, weakly paired data (e.g., pairs from the same subject) or other types of supervision from additional tasks like dense segmentation are still required in these methods.

In this paper, we address the problem of US-to-MR image synthesis by a learning-based framework. Fig. 1 summarises the proposed framework. We design an anatomically constrained self-supervised model to learn the mapping from US to MR under an assumption that US and MR share a common representation in a latent space. The anatomical constraint considers both the latent space and the geometrical structures between the two modalities. To the best of our knowledge, this article presents the first attempt towards unpaired US-to-MR synthesis in a self-supervised manner. Qualitative and quantitative experiments demonstrate that the proposed approach generates realistic MR images, even with highly-imbalanced data.

Relationship to Preliminary Work [17]: An early version of this work was presented in [17]. In the current paper, we considerably expand the preliminary study by: 1) We further propose three new solutions to utilise 3D auxiliary information and boost the synthesis performance. Leveraging additional neighbouring inputs and predicting the position in 3D space as an auxiliary task are explored to achieve the goal. 2) In this version we provide a more detailed analysis of our framework and its new extensions. The detailed architecture of each network components are elaborated; more details with additional illustrations are included for the EdgeNet; new perspectives on the proposed cross-modal attention is included; a detailed analysis of the anatomical latent space is presented. 3) Additional experimental evaluations are included in this extension, with more training details; standard deviation for all the quantitative performance is reported for better understanding of the models; more qualitative results are presented with comparison to other solutions; an anatomy-preserving analysis is presented to evaluate the effectiveness of the proposed approach for both synthetic structures and real data; performance on the aforementioned 3D-based solutions is also reported with comparison to our preliminary results.

The main highlights of the paper are summarised as:

- We present an approach to synthesise MR-like images from unpaired US images;
- We propose an anatomy-aware deep neural network architecture with mono-directional consistency, to address the synthesis problem in a self-supervised manner;
- Based on the shared latent space assumption, we propose a latent space consistency constraint;
- We propose a cross-modal attention module that propagates information across modalities in the feature domain;
- We propose to leverage 3D auxiliary information to reduce ambiguity during image synthesis;
- Comprehensive experimental evaluation and analysis show that the proposed synthesis framework is able to generate high-quality MR-like images and performs favourably against other alternative methods.

The rest of this paper is organised as follows. Section II reviews related work and discusses the main differences to our approach. In Section III, we elaborate on the detail of our self-supervised US-to-MR image synthesis approach, with analysis of the network architecture design and the underlying representative features. Following that, we perform extensive experiments evaluating the effectiveness of the proposed framework both qualitatively and quantitatively in Section IV. In addition, potential applications derived from this work and the model generalisations are discussed in Section V. Finally, the paper is concluded in Section VI and possible future directions are discussed.

II. RELATED WORK

A. Medical Image Synthesis

Medical image synthesis or simulation aims to synthetically generate one imaging modality from another. Classical methods to achieve this have been based on segmentation and atlases. Segmentation-based approaches are straightforward, and in the case of CT synthesis from MRI, may, for instance, first segment the MR images into different parts (e.g., bony structure, soft tissue) and then assign the corresponding CT number to each part. In [18], the authors study the possibility of radiotherapy treatment planning using only MR by bone and water segmentation. Berker et al. propose to address the MRI-based attenuation correction problem by segmenting air, bone and tissues [10]. On the other hand, atlas-based methods generate a synthetic CT by deforming a CT atlas onto the patient space, where the required deformation is found by registering an MR atlas to the patient real MR image. Dowling et al. [19] generate pseudo-CT scans by nonrigid registration of an MRI atlas to an MRI scan. The authors of [12] use atlas-based regression to deform a collection of atlas CTs into a single pseudo-CT, based on the target MR and an atlas database. However, the segmentation-based approaches suffer from requiring a time-consuming segmentation, while the uncertainty in registration (e.g., missing tissues) is a critical inherent limitation of atlas-based methods. With the recent progress of deep learning in medical image analysis, CNN-based approaches have started to dominant the medical image synthesis. Zhao et al. [3] propose to directly optimise the mapping function for MR to CT synthesis, with reference to different 3D views. However, such a regression-based method may lead to blurred results if there are misalignments between CT and MR images. To handle this problem, some works [4], [6] augment the regression-based loss with another adversarial loss in a generative adversarial network framework [20]. Although blur caused by misalignment has been addressed, a training set of paired images (e.g., MRI and CT) is still necessary for the above models. Whereas such paired training data is very scarce for many clinical imaging applications.

B. Self-Supervised Learning

Self-supervised (also termed ‘unsupervised’ in some literature) learning is a learning technique that does not rely on the supervision from external label/annotation, i.e., the

learning is totally based on the available data itself. An auto-encoder (AE) [21] is one of the most basic self-supervised learning approaches, which optimises a self-reconstruction loss by recovering the input signal itself. Derivatives including the denoising auto-encoder (DAE) [22] and variational auto-encoder (VAE) [23] also focus on self-supervised learning but with different optimisation functions. Recently, Goodfellow et al. proposed the Generative Adversarial Network (GAN) [20] which learns to generate meaningful signals from random noise, by playing a minimax game. Based on the GAN framework, an architecture named CycleGAN [16] or DualGAN [24] is proposed to address the problem of image-to-image translation, without the dependency on paired training data. Consequently, such an architecture with cycle consistency has enabled a group of medical image synthesis methods in recent years, mainly focusing on MR-to-CT image synthesis [5], [15], [25], [26], and vice versa [15], [27]. Building on top of the original CycleGAN, in [5] and [26], additional loss terms have been introduced to further constrain the structure features. Zeng and Zheng [28] propose a hybrid GAN that combines a 3D generator and a 2D discriminator to synthesise CT from MR images, in a weakly-supervised manner. In [15], the authors propose a solution to MR-CT synthesis by a 3D CNN composed of mutually beneficial generators and segmentors with cycle- and shape-consistency. Paired data dependency has been alleviated to some extent by the above CycleGAN-based approaches. However, aligned or weakly-aligned data or auxiliary tasks are still necessary in these works. Besides, MR and CT are relatively, similar in anatomically appearance and relatively easier to align, when compared with ultrasound. To our knowledge, there is no prior work on cross-modal image synthesis from ultrasound (US) data, in a data-driven self-supervised manner.

C. Ultrasound Image Analysis

Different from the aforementioned medical imaging modalities MRI and CT, US imaging usually does not present as clear and sharp anatomical structures. However, its real-time and un-harmful properties make it a much more suitable choice for many medical screening scenarios, including fetal development monitoring. Prior fetal US image analysis work mainly focuses on fetal anatomy detection [29]–[33] and registration to other modalities [8], [9], [34], [35]. Maraci et al. [29] propose an approach to make the US diagnosis easier by combining simple US scanning protocols with machine learning solutions. Yaqub et al. [30] propose a random forest based classifier to categorise fetal US images. With the help of deep learning techniques, Chen et al. [31] present a CNN-based approach to locate the fetal abdominal standard plane in US videos. Some methods [33] utilise human eye-gaze data to assist standard plane detection. The fusion of tracked US with other modalities like CT and MRI has benefits for a variety of clinical applications. Wein et al. [34] develop methods to simulate US from CT in real-time, while in [35] the authors evaluate the performance of methods of MRI to US registration. Kuklisova et al. [8] propose a method for 3D fetal brain US and MRI registration by simulating a pseudo-US from an MR volume segmentation. While most existing

work focusing on the above US image analysis topics, there lacks a study on US to MRI synthesis in the literature.

III. METHOD

In this section, we elaborate the proposed approach for US-to-MR image synthesis in detail. Specifically, we first pre-process the US volumes by cropping and automatically aligning the US volumes as described in [36]. Following that, we manually align the MR volume to the same reference space. Then we propose a novel learning-based framework for unpaired US-to-MR synthesis, which is illustrated in Fig. 2 (detailed structure of the blue block in Fig. 1(a)). Given a source US image, the corresponding MR image is synthesised with reference to real MR data. In addition to pixel-level (*rec. loss*) constraints, a distribution similarity (*dis. loss*) is also incorporated to address the unpaired data property and ensure anatomical consistency. As our objective is to synthesis MR from US images, the overall design of the proposed framework is mono-directional, instead of bi-directional as in the CycleGAN architecture [16]. That is, we only have the forward cycle (i.e., US→MR→US) without the reverse cycle (i.e., MR→US→MR), and we experimentally find that such a design leads to less ambiguity for image synthesis in our case. Next, we elaborate on each of the proposed components in detail.

A. Anatomy-Aware Synthesis

Paired (e.g., same fetus at the same gestational age) fetal brain US and MR data is rare in clinical practice, and unavailable in our case. Even if it is available, US and MR are not simultaneously co-registered as has often been assumed in prior medical image synthesis methods [3], [4]. Hence it is infeasible to learn the mapping from US to MR directly by traditional CNN-based techniques for our task. Therefore, we propose to address the problem through a synthesis framework, by enforcing the synthesised MR images to lie in a similar distribution to real MR data. Throughout the synthesis, an important objective is to correctly map the clinically important anatomical structures between the two modalities. As a result, anatomy-aware constraints are specifically designed to implicitly preserve anatomy consistency.

1) *Anatomical Feature Extraction*: As paired data is unavailable, we assume that the US and MR images share an anatomical latent space (Fig. 1(b)). Building upon this assumption, instead of using the pixels in the image domain, we propose to extract the underlying anatomical features and synthesise images in the corresponding MR domain accordingly. Specifically, we leverage an autoencoder to extract the latent features, as shown in the bottom-left part of Fig. 2 (encoder-A→decoder-B). Assume the set of n source US images as $\{x_U^i\}_{i=1}^n$ where $x_U^i \in \mathcal{X}_U$ is the i^{th} image, the extracted anatomical feature is formally defined as $y^i = F(x_U^i)$ where $F(\cdot)$ is the encoder.

2) *Bi-directional Latent Space Consistency*: The above extracted latent features are fed into decoder-C to synthesise the corresponding MR image. As pixel-level supervision is unavailable for *Synth. MR*, we use a backward-inference path

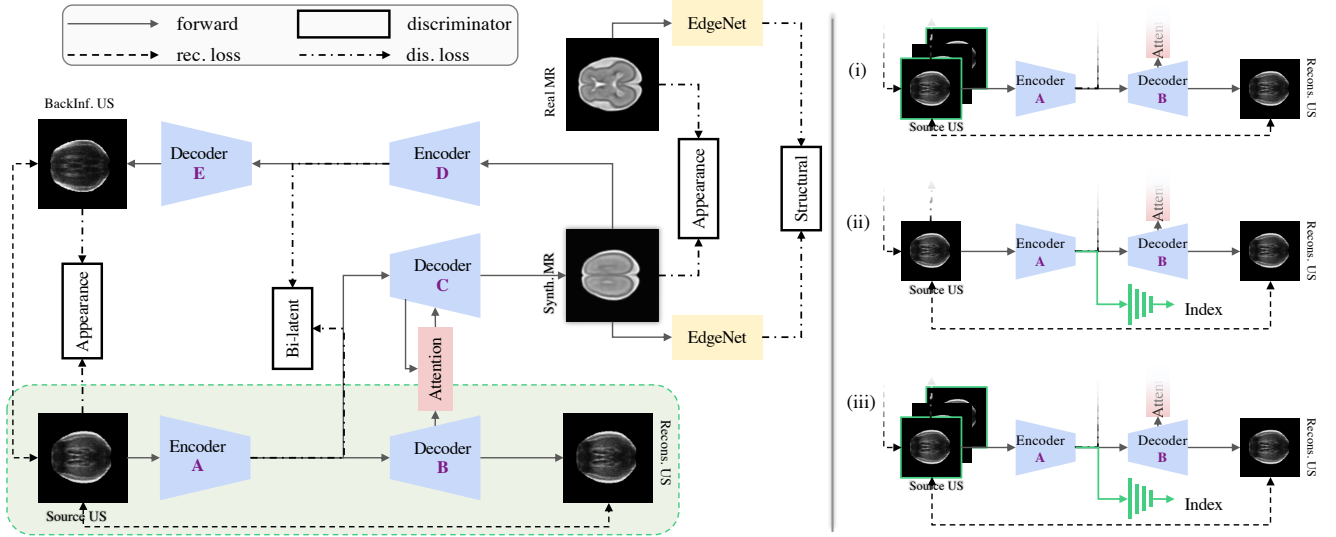


Fig. 2. **Left:** Architecture of the proposed US-to-MR synthesis framework. **Right:** Solutions to leverage 3D auxiliary information, each of which can be plugged into the green part in the left framework. (i) Augmented with neighbouring slices in the 3D volume; (ii) Predicting the index of the slice in the 3D volume; (iii) Both with augmented slices and index prediction.

(encoder-D→decoder-E) to recover the source US. Denoting the encoded latent feature (at the end of encoder-D) as y_b^i , we propose a bi-directional latent space consistency constraint, based on the assumption of shared latent space. As a result, y^i and y_b^i are forced to lie in a similar distribution by means of adversarial learning (*Bi-latent* block in Fig. 2).

3) *Structural Consistency:* Although the anatomical feature extraction module encodes the main structure of an US image, the image structure in the MR domain is quite different in appearance compared to that in the US domain. To synthesise realistic MR images, we further propose a constraint to enforce the structures of the *Synth. MR* and the *Real MR* to be similar. Noting the unpaired nature of our data, we choose to constrain the structural information to lie in a similar distribution. Specifically, the edge information of the *Synth. MR* and *Real MR* is extracted by an EdgeNet, following which a structural discriminator (*Structural* block in Fig. 2) is leveraged to compute the edge similarity.

B. Cross-modal Attention

Based on the aforementioned components, the MR image synthesis process is mainly guided by the latent features y^i extracted from encoder-A. To further leverage guidance across different modalities, we propose a cross-modal attention module between the US decoder-B and the MR decoder-C, as shown in Fig. 2 (the red *Attention* block). Specifically, the US features are reformulated as self-attention guidance for MR image synthesis, and such a guidance is applied to the MR features implicitly in an attentive manner (detailed in Fig. 3). The cross-modal attention module consists of several 1×1 convolutional layers and a skip connection, without any modification to the input feature dimension. This cross-modal attention module leverages guidance across different modalities (MR and US here) by cross-referencing the features from the two modalities. Specifically, the features from US are

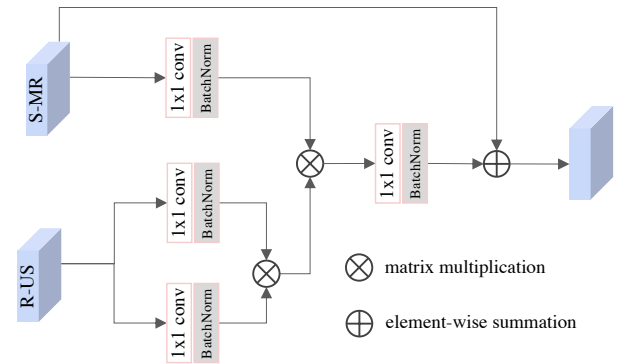


Fig. 3. Detailed architecture of the proposed cross-modal attention module.

combined with the features from MR by matrix multiplication, which can also be considered as an approximation of the mutual information acquisition across these two modalities. The 1×1 convolutions adapt the channel dimension for the subsequent multiplication. Since our target is to synthesise MR-like images, the original MR features are added back to the mutual information by a skip connection, which also keeps the feature dimension. A similar idea for single modality attention (also termed as self-attention [37] or non-local scheme [38]) has been shown to be effective to leverage neighbouring information in natural image analysis.

Denoting the features from the US reconstruction decoder-B (R-US in Fig. 3) as f_U and the features from the MR synthesis decoder-C (S-MR in Fig. 3) as f_M , the updated feature after the cross-modal attention module is defined as:

$$\tilde{f}_M = \eta(\delta(f_U)^T \otimes \phi(f_U) \otimes g(f_M)) + f_M, \quad (1)$$

where η, δ, ϕ, g are linear embedding functions which are implemented by 1×1 convolutions. The proposed cross-modal attention enforces the network to not only focus on local activations (favoured by CNNs) but also non-local context

information via both self- and cross-modal attention. Specifically, the local and non-local correlation is modelled by feature matrix multiplication (\otimes).

C. Joint Adversarial Objective Function

In this section, we now formally define the objective function that is used to train the proposed US-to-MR synthesis model. As aforementioned (and can be observed in Fig. 2), there are two forms of objective terms: one based on pixel-wise reconstruction (*rec. loss*) and the second a discriminator-based distribution similarity (*dis. loss*). The *rec. loss* is defined as an ℓ_1 -norm, while the *dis. loss* is achieved by a discriminator with adversarial learning [20].

1) *Generative Adversarial Learning*: The generative adversarial network (GAN) is a self-supervised learning framework proposed in [20], which consists of two modules, namely a generator and a discriminator. The main idea of a GAN is to play a minimax game with these two modules. The training of a GAN forces the distribution of the source data x to be similar to the target data (y) distribution. Suppose the mapping from x to a new data space is $G(x; \theta_g)$ where G is the generator with parameters θ_g . The discriminator $D(y; \theta_d)$ outputs a binary scalar indicating whether the data y is real or fake. Then the D net is trained to maximise the probability of assigning the correct label to both real data samples and samples from G . Meanwhile, the G net is trained to minimise $\log(1 - D(G(x)))$. The final objective function for a GAN is defined as: $\min_G \max_D \mathbb{E}[\log D(y)] + \mathbb{E}[\log(1 - D(G(x)))]$.

2) *Proposed Joint Objective Function*: Denoting the reconstructed US as $\hat{x}_U \in \hat{\mathcal{X}}_U$, latent feature as $y \in \mathcal{Y}$, synthesised MR as $\hat{x}_M \in \hat{\mathcal{X}}_M$, and real MR as $x_M \in \mathcal{X}_M$. The objective for the forward path (US \rightarrow MR) is defined as:

$$\{\min \mathcal{L}_F \mid \mathcal{L}_F = \lambda \mathcal{L}_{lat} + \gamma \mathcal{L}_{app} + \gamma \mathcal{L}_{stru}\}, \quad (2)$$

where

$$\mathcal{L}_{lat} = \mathbb{E}_{x_U \in \mathcal{X}_U} \|G_U(F(x_U)) - x_U\|_1, \quad (3)$$

$$\mathcal{L}_{app} = \mathbb{E}_{x_M \in \mathcal{X}_M} (D_{app}(x_M)) + \mathbb{E}_{y \in \mathcal{Y}} \log(1 - D_{app}(G_M(y))), \quad (4)$$

$$\mathcal{L}_{stru} = \mathbb{E}_{x_M \in \mathcal{X}_M} (D_{stru}(E(x_M))) + \mathbb{E}_{y \in \mathcal{Y}} \log(1 - D_{stru}(E(G_M(y)))). \quad (5)$$

Here the \mathcal{L}_{lat} , \mathcal{L}_{app} , \mathcal{L}_{stru} are loss terms for the reconstruction from latent space, appearance, and structural consistency, respectively. The first term \mathcal{L}_{lat} represents the generator loss while the following two terms indicate the discriminator loss. The decoder-B that is used to reconstruct *Recons.US* is represented by G_U , while the decoder-C for MR synthesis is represented by G_M . The \hat{x}_U and \hat{x}_M are defined as: $\hat{x}_U = G_U(F(x_U))$, $\hat{x}_M = G_M(y)$. The discriminators D_{app} and D_{stru} that each consists of four *conv* layers are used to measure the similarity for appearance and structure, respectively. The EdgeNet is represented by E , while parameters λ and γ are balancing weights for the objective terms, so that they lie on a similar scale.

Denoting the *BackInf.US* recovered from the *Synth.MR* as $\tilde{x}_U \in \tilde{\mathcal{X}}_U$ and the back-inferred feature at the end of encoder-D as $y^{back} \in \mathcal{Y}^{back}$, the objective for the backward path (*Synth.MR* \rightarrow *BackInf.US*) is defined as:

$$\{\min \mathcal{L}_B \mid \mathcal{L}_B = \lambda \mathcal{L}_{proj} + \gamma \mathcal{L}_{app}^{back} + \gamma \mathcal{L}_{bi}\}, \quad (6)$$

where

$$\mathcal{L}_{proj} = \mathbb{E}_{\tilde{x}_U \in \tilde{\mathcal{X}}_U, x_U \in \mathcal{X}_U} \|\tilde{x}_U - x_U\|_1, \quad (7)$$

$$\mathcal{L}_{app}^{back} = \mathbb{E}_{x_U \in \mathcal{X}_U} (D_{app}^{back}(x_U)) + \mathbb{E}_{y^{back} \in \mathcal{Y}^{back}} \log(1 - D_{app}^{back}(G_{BU}(y^{back}))), \quad (8)$$

$$\mathcal{L}_{bi} = \mathbb{E}_{y \in \mathcal{Y}} (D_{bi}(y)) + \mathbb{E}_{y^{back} \in \mathcal{Y}^{back}} \log(1 - D_{bi}(y^{back})). \quad (9)$$

Here the \mathcal{L}_{proj} , \mathcal{L}_{app}^{back} , \mathcal{L}_{bi} are loss terms for the back-inference reconstruction, backward appearance, and bi-directional latent space consistency, respectively. Similar to Eq. 2, parameters λ and γ are balancing weights. The decoder-E that is used to back recover *BackInf.US* is represented by G_{BU} and $\tilde{x}_U = G_{BU}(y^{back})$. The discriminators D_{app}^{back} and D_{bi} are used to compute the similarity for backward-inference and bi-directional latent space, respectively. Based on the above defined objective terms, the final joint loss function for our model training is defined as:

$$\mathcal{L} = \mathcal{L}_F + \mathcal{L}_B. \quad (10)$$

D. 3D Auxiliary Information

Here we investigate the possibility of leveraging 3D volumetric information to improve the synthesis. The proposed approaches to leverage 3D information are shown in Fig. 2-Right. Specifically, we propose three strategies to achieve the goal: 1) by adding neighbouring slices as augmented input; 2) by predicting the position/index of the current slice in the volume as an additional task; 3) with both the augmented input and the index prediction task. For simplicity, we only show the modified part (the green block in Fig. 2-Left) compared to the 2D-based model. These approaches are motivated by the constraints humans appear to have when viewing slice-wise volumetric data. We assume that if the model is able to utilise the 3D positional information (by either referring to neighbours or directly reasoning its position), it could have a more thorough understanding of the whole anatomical structure and alleviate synthesis ambiguity. The above modifications do not severely influence the original network architecture. The augmented input only leads to a channel number update (1 \rightarrow 3) for the first layer of Encoder-A, while all the decoders only output the middle slice. The index prediction branch is implemented by four convolutional layers with a fully-connected layer, in a regression manner. In the rest of this paper, we use the original 2D settings (by default) as aforementioned, unless otherwise specified.

E. Network Architecture

The detailed network architecture design and parameters are presented in Table I. The proposed network basically consists of convolutional (*conv*) layers, up-convolutional (*up-conv*) layers, and pooling (*maxpool*) layers. Specifically, each part is described as follows:

TABLE I

DETAILS OF THE NETWORK DESIGN. THE FIVE MAIN COMPONENTS (SUBNETWORKS) ARE PRESENTED WITH DETAILED PARAMETER SETTINGS. THE NUMBERS FOLLOW EACH *conv/up-conv* ARE THE KERNEL SIZE AND NUMBER OF CHANNELS, WHILE THE NUMBER FOLLOW *maxpool* IS THE SCALE. *FC* REPRESENTS A FULLY-CONNECTED LAYER AND THE NUMBERS FOLLOWING ARE THE INPUT AND OUTPUT DIMENSIONS. THE SYMBOL \odot REPRESENTS CONCATENATION, WHILE \oplus AND \otimes THE ELEMENT-WISE SUMMATION AND MATRIX MULTIPLICATION.

Encoder			Cross-modal Attention			Discriminator		
Layer	Input	Parameter	Layer	Input	Parameter	Layer	Input	Parameter
enc_1	US/MR img	conv, 3×3, 48 conv, 3×3, 48 maxpool, 2	convMR	S-MR	conv, 1×1, 72	disc_1	img/ edge/ feature	conv, 3×3, 64, stride=2
enc_2	enc_1	conv, 3×3, 48 maxpool, 2	convUS_1	R-US	conv, 1×1, 72	disc_2	disc_1	conv, 3×3, 128, stride=2
enc_3	enc_2	conv, 3×3, 48 maxpool, 2	convUS_2	R-US	conv, 1×1, 72	disc_3	disc_2	conv, 3×3, 256, stride=2
enc_4	enc_3	conv, 3×3, 48 maxpool, 2	fuse	convUS_1 \otimes convUS_2 \otimes convMR	conv, 1×1, 144	disc_4	disc_3	conv, 3×3, 512
enc_5	enc_4	conv, 3×3, 48 maxpool, 2	output	S-MR \oplus fuse	-	out	disc_4	conv, 3×3, 1
Decoder			EdgeNet			Volume Index Predictor		
Layer	Input	Parameter	Layer	Input	Parameter	Layer	Input	Parameter
dec_1	enc_5	conv, 3×3, 48 up-conv, 3×3, 48, stride=2 conv, 3×3, 96	g_hori	synthetic MR/ real MR	conv, 1×5, 1	Ind_1	enc_5	conv, 3×3, 48
dec_2	dec_1 \odot enc_4	conv, 3×3, 96 up-conv, 3×3, 96, stride=2 conv, 3×3, 96	g_vert	g_hori	conv, 5×1, 1	Ind_2	Ind_1	conv, 3×3, 32
dec_3	dec_2 \odot enc_3	conv, 3×3, 96 up-conv, 3×3, 96, stride=2 conv, 3×3, 96	s_hori	g_vert	conv, 3×3, 1	Ind_3	Ind_2	conv, 3×3, 16
dec_4	dec_3 \odot enc_2	conv, 3×3, 96 up-conv, 3×3, 96, stride=2 conv, 3×3, 96	s_vert	s_hori	conv, 3×3, 1	Ind_4	Ind_3	conv, 3×3, 4
dec_5	dec_4 \odot enc_1	conv, 3×3, 96 up-conv, 3×3, 96, stride=2				Out_Ind	Ind_4	FC, 16, 2
dec_6	dec_5 \odot US	conv, 3×3, 64 conv, 3×3, 32						
output	dec_6	conv, 3×3, out_channel						

1) *Encoder and Decoder*: All the encoders (and the decoders) share the same architecture as shown in the left part of Table I. The encoder takes either the US image or MR image as input and consists of five *enc* blocks. On the other hand, the decoder is composed of a mirrored architecture to the encoder to reconstruct/synthesise the target image. Each *enc* block consists of *conv* layers and *maxpool* layers, while the *dec* block in the decoder mainly consists of *conv* layers and *up-conv* layers. Skip connections are added between the encoder and decoder.

2) *Cross-modal Attention*: The cross-modal attention module is proposed to implicitly learn the feature fusion strategy between the US features and MR features. As illustrated in Fig. 3 and Table I, this module is implemented by several basic 1×1 *conv* layers, with matrix multiplication. The *conv* layers are utilised to adjust the feature dimension and combination weights for the following fusion.

3) *EdgeNet*: The EdgeNet extracts edges from the input image (either US or MRI). It consists of four *conv* layers, imitating the Canny edge detector (can also be other edge detectors). Specifically, the input images are first smoothed with Gaussian kernels and convolved with Sobel edge filters, followed by non-maximum suppression and thresholding. Since the edge maps extracted from the EdgeNet are further

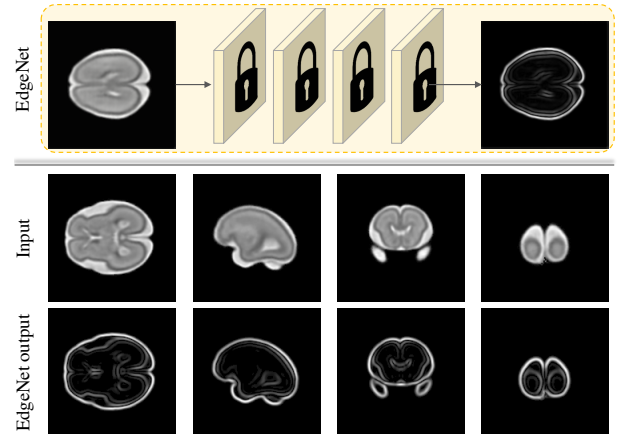


Fig. 4. **Top**: Illustration of the EdgeNet, where the lock symbol indicates a layer frozen. **Bottom**: Example results, where the first row shows the input MR images and the second row shows the detected edges from the EdgeNet.

fed into the discriminator for loss computation, all the parameters in this module are fixed, avoid updating. Illustration of the proposed EdgeNet, together with example generated edge maps are shown in Fig. 4 for reference.

4) *Discriminator*: There are four discriminators in our main architecture, measuring the similarity in terms of appearance, latent space, and structure. Each discriminator is composed of five *conv* layers and outputs a scalar value, indicating whether the input is real or fake. The discriminator essentially acts as a binary classifier.

5) *Volume Index Predictor*: The index prediction branch consists of four *conv* layers and an FC layer. All the *conv* layers are with a 3×3 kernel size and decrease in the channel dimension till the final fully-connected output layer. The index predictor takes the feature at the latent space as input and predicts the corresponding index of the original input US. Possible solutions of leveraging the index predictor are illustrated in Fig. 2-Right.

F. Anatomical Space Analysis

Our basic assumption is that the US and MR data share a similar anatomical latent space (as illustrated in Fig. 1). To better understand the anatomical property of the shared latent space, here we explore it by analysing the features between the encoder and decoder. Specifically, we visualise the feature maps at the end of Encoder-A and Encoder-D in Fig. 5. From the feature visualisation we can see that for the forward pass, the features focus more on the inner-part anatomical structures, while for the backward pass, the features primarily learn the overall structure of the brain. This is mainly due to the forward pass not only needing reconstruct the US itself but also it has to synthesise the corresponding MR image. As a result, the shared anatomical features are learned at this point. On the other hand, the backward pass aims to infer the original US image and is simultaneously constrained by the latent space from the forward pass, thus focusing more on the global structure. In addition to the features in the shared latent space, we also visualise the corresponding attention map (by the approach in [39]), which represents where the model pays most attention. Similarly, it can be observed that the forward pass focuses more on the internal structure while the backward pass attention depicts the boundary of the brain. The above analysis provides some evidence to validate our assumed anatomical latent space. Note that the assumption of bi-directional latent space consistency aims not to force the forward and backward features of the latent space to be *identical*, but to be similar in distribution, as determined by the discriminator.

IV. EXPERIMENTS

A. Data and Implementation Details

The training and evaluation of the proposed US-to-MR synthesis framework are based on a dataset consisting of healthy fetal brain US and MR volumes. We obtained the fetal US data from a multi-centre, ethnically diverse dataset [40] of 3D ultrasound scans collected from normal pregnancies. The MR data is obtained from the CRL fetal brain atlas [41] database and additional data scanned at Hammersmith Hospital. As proof of principle, we selected the gestational age of 23-week for US and MR data. In the aggregate, the whole dataset consists of 107 US and 2 MRI volumes. Around

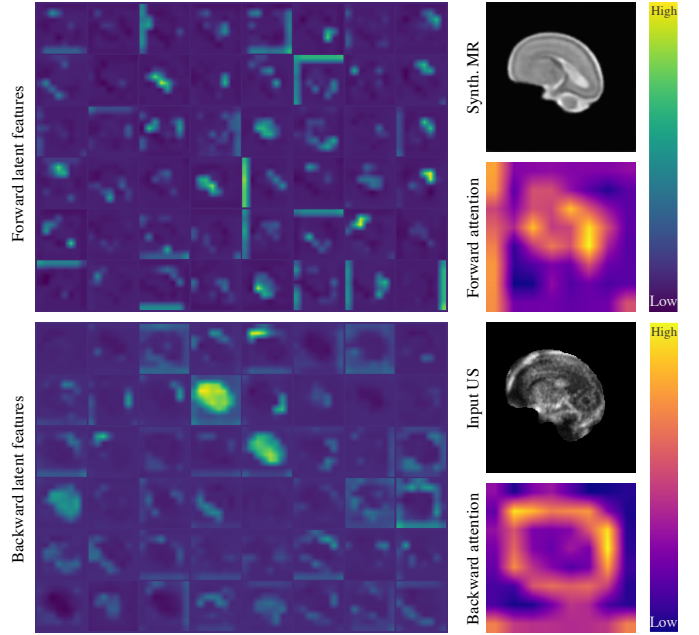


Fig. 5. Anatomical latent space visualisation. The top part shows the features (a grid of 48 feature maps) for the forward pass, i.e., at the end of Encoder-A, while the lower part shows the features for the backward pass, i.e., at the end of Encoder-D. The corresponding attention maps with the US and MR images are also shown on the side.

36,000 2D US slices and 600 MRI slices were extracted accordingly. 80% of the whole database accounted for the training and validation sets, while the remaining 20% acted as the testing set. As detailed in Table I, the proposed model was implemented by simple *conv*, *up-conv*, and *maxpool* layers. Skip connections were included for each encoder-decoder pair to enhance the structural details. The balancing weights λ and γ were empirically set to 10 and 1, respectively. All the images (US, MRI, and the corresponding edge maps) used in our framework are of size 160×160 pixels. The index range for the index prediction task is 160. The learning rate was initialised as 10^{-4} and decayed by half for every 20 training epochs. The whole model was trained for 100 epochs. The generator and discriminator are optimised iteratively, i.e. updating the generator for every update of the discriminator. Our whole model was implemented using the PyTorch framework and trained on an Nvidia Titan V GPU in an end-to-end manner. Taking an US image as input, only the A, B, C and Attention blocks in Fig. 2 are remained during the model inference, without the whole backward path and all the discriminators. Implementation of the proposed approach is available online¹.

B. Evaluation Metrics

The commonly used evaluation metrics like PSNR (Peak Signal-to-Noise Ratio) or SSIM (Structural Similarity) are not applicable in our study, as US-MR data is not paired. As a result, we propose to leverage two other alternative metrics for the quality evaluation of our synthesised MR images: 1) the MOS (Mean Opinion Score) and 2) the Deformation score

¹<https://bitbucket.org/JianboJiao/ssus2mri>

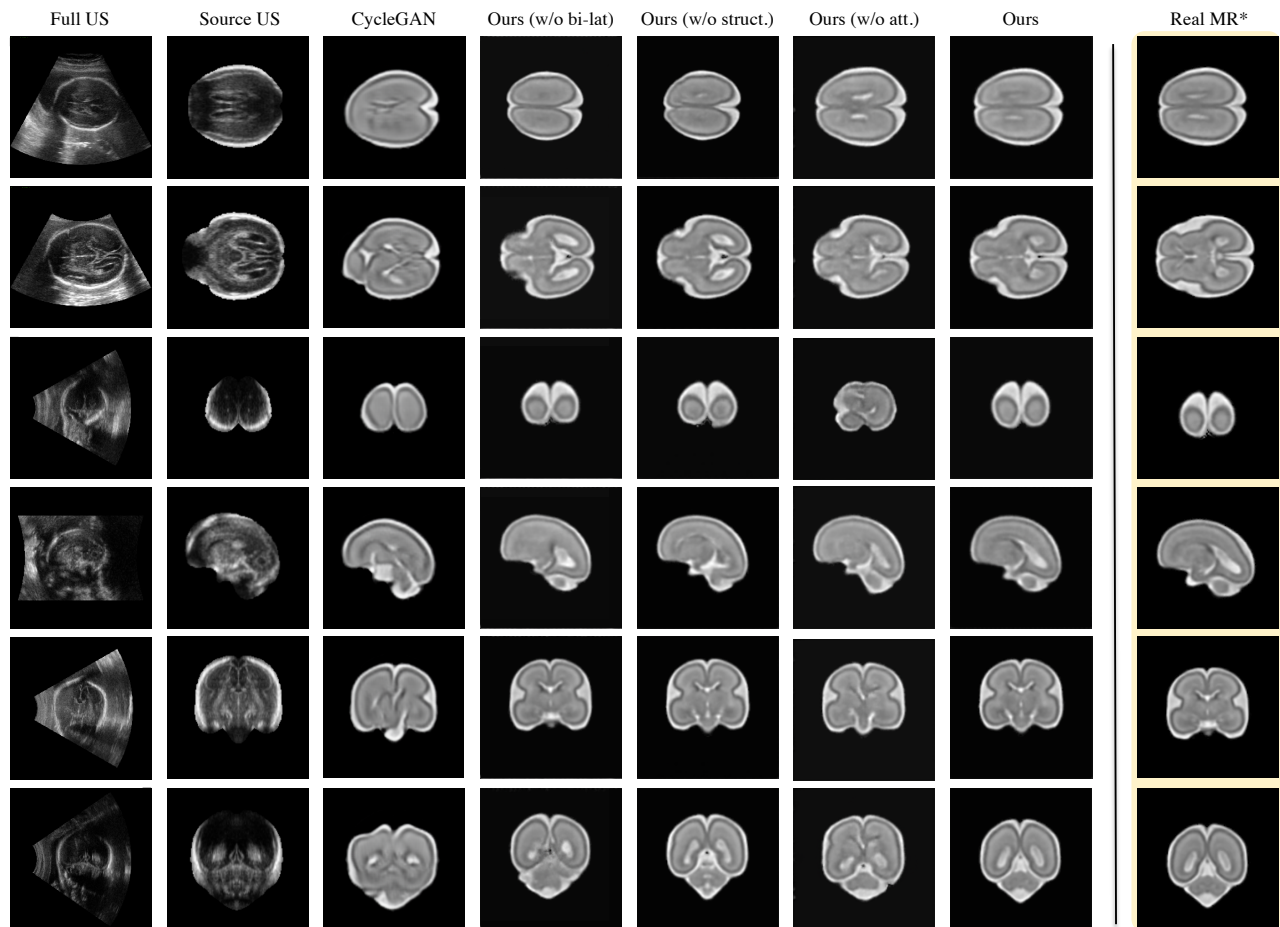


Fig. 6. Qualitative performance on the US-to-MR image synthesis. Each row shows an example sample and from left to right: the original full US, pre-processed source US, synthesised MR by CycleGAN [16] and our approach (with its counterparts). Some real MR samples are also shown for reference in the last column. *Note that the highlighted *Real MR* is **NOT** the exact corresponding MR to the input US, instead illustrative examples for visual comparison.

(registration metric based on Jacobian). The MOS measures the quality of a given image by a rating score between 1 and 5: 1 indicates *inferior* while 5 indicates *superior*. A user-study was performed to achieve the MOS performance, given participants from two groups (2 medical experts and 11 beginners), in which each observer was shown with 80 samples. For the Deformation score, an FFD-based [42] deformable registration was applied to the synthesised MR to register to a real MR at a similar imaging plane. The average Jacobian (normalised to $[0,1]$) of the required deformation to complete such registration was computed as the score consequently. The underlying assumption is that a synthesised MRI with high-quality tends to have a lower Jacobian when registering to the real MRI.

C. Qualitative and Quantitative Performance

In this section, we firstly present qualitative results of the synthesised MR and secondly quantitatively evaluate synthesis results. In the testing phase, given a test 2D US image as input, the corresponding MR image is synthesised accordingly. Several US example inputs with corresponding synthesised MR images are shown in Fig. 6. Note that the last column (*Real MR*) is **not** in direct correspondence to the input US,

instead is only presented for reference. These reference MR images are selected from a similar 3D position to the input US images. As a result, although these *Real MR* images are not perfectly aligned to the input US, we assume they are valid references for readers to compare the visual performance. It can be observed from the results in Fig. 6 that the visual appearance of our synthesised MR images is very similar to the real ones. Further the results generated using our approach are visually superior to the results from an alternative approach, CycleGAN [16] that is widely used for image synthesis tasks. Additionally, the anatomical structures between the source US and the synthetic MR are well preserved.

The quantitative performance is reported in Table II, where the evaluation metrics of MOS and deformation are presented. Furthermore, the proposed approach is compared with several alternative methods including an autoencoder (AE), GAN [20], and CycleGAN [16]. The presented results suggest that the performance of the proposed US-to-MR synthesis framework surpasses the other CNN-based architectures.

D. Ablation Study

To better understand the effectiveness of each proposed components, we performed an ablation study by removing

TABLE II

QUANTITATIVE PERFORMANCE FOR MRI SYNTHESIS WITH COMPARISON TO SEVERAL ALTERNATIVE APPROACHES ON MOS SCORE AND DEFORMATION SCORE. THE STANDARD DEVIATION (\pm STD.) IS ALSO SHOWN. MOS THE HIGHER THE BETTER, WHILE DEFORMATION THE LOWER THE BETTER.

Settings		AE	GAN	CycleGAN	Ours (w/o bi-lat)	Ours (w/o struct.)	Ours (w/o att.)	Ours	Real
MOS \uparrow	Expert	1.00 ± 0.00	2.05 ± 1.12	2.50 ± 0.53	3.05 ± 0.80	3.45 ± 1.30	3.30 ± 1.14	3.90 ± 0.81	4.35 ± 0.97
	Beginner	1.01 ± 0.03	2.75 ± 1.28	3.42 ± 0.36	3.69 ± 0.71	3.87 ± 0.74	3.65 ± 0.73	4.08 ± 0.42	4.23 ± 0.67
Deformation \downarrow		0.97 ± 0.09	0.78 ± 0.46	0.66 ± 0.46	0.55 ± 0.22	0.65 ± 0.39	0.47 ± 0.31	0.46 ± 0.24	0.00 ± 0.00

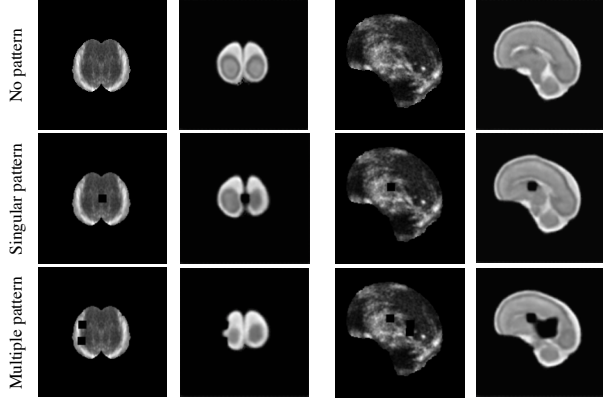


Fig. 7. Anatomy-preserving performance for US-to-MR image synthesis on synthetic patterns. For each example the first column shows the input US and the second column the synthesised MRI by our approach.

each component at a time: the bi-directional latent consistency module (*w/o bi-lat*), the structural consistency module (*w/o struct.*), and the cross-modal attention module (*w/o att.*). The corresponding qualitative and quantitative results are shown in Fig. 6 and Table II, respectively. It can be observed from the results that the model performs worse when removing any of the above components. Specifically, the bi-directional latent space and the cross-modal attention contribute the most to the model performance, which validates our initial assumption of the importance of the shared anatomical space. The structural consistency contributes more to the detailed structures, which is revealed by the deformation metric. The above qualitative and quantitative results support the inclusion of each proposed component in our model.

E. Anatomy-Preserving Analysis

In order to evaluate the performance of our approach with respect to the anatomy-preserving property, we perform an analysis based on synthetic abnormal data and real pseudo-paired data:

a) Synthetic Abnormal Data: We first randomly apply some synthetic patterns to the input US images and evaluate whether these patterns are preserved during the synthesis of the corresponding MR images. While various patterns can be applied, here for simplicity, we use a square pattern, which we apply in two ways: singular pattern and multiple patterns. Note that our trained model is directly applied to these data without any fine-tuning. Some examples are illustrated in Fig. 7. We can see from the results that the applied pattern regions are well preserved in the synthesised MR images. Note that in

TABLE III
QUANTITATIVE EVALUATION ON OUR SYNTHESISED MR IMAGES FOR SYNTHETIC PATTERN PRESERVING ANALYSIS.

Settings	AE	GAN	CycleGAN	Ours
PSNR (dB) \uparrow	31.56 ± 3.91	34.63 ± 13.05	43.07 ± 15.07	99.37 ± 1.55

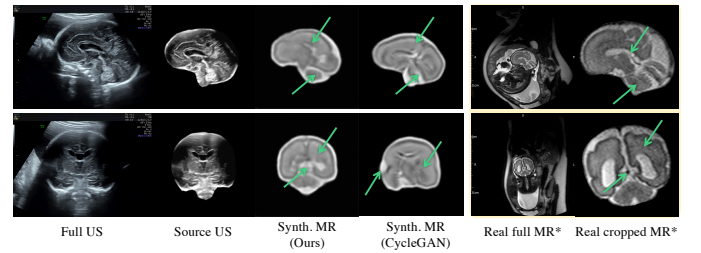


Fig. 8. Anatomy-preserving performance for US-to-MR image synthesis on real pseudo-paired data. Key anatomical structures (marked by arrows) are preserved when compared to the real MR images. *The presented MR examples on the right side are not exactly aligned to the left US images.

the second example of the multiple-pattern, a larger pattern is generated in the MRI compared to that in the input US. We speculate that this is caused by the small distance between the multiple patterns that makes the synthesis ambiguous for the network. To quantify the anatomy-preserving property of the proposed model, we further calculate the similarity between the original pattern (multiple version) and the patterns preserved in the synthesised MRI, using the same testing set as aforementioned. The corresponding quantitative result is reported in Table III, where we perform a comparison to the alternative architectures of AE, GAN [20], and CycleGAN [16]. We can see from the result that our approach performs much better than the others, which reveals the anatomy-preserving property of the proposed method.

b) Real Pseudo-Paired Data: We obtained some anonymised US data with corresponding MRI of the same subject from the John Radcliffe Hospital. Such data can be considered as “pseudo-paired”, as these US-MR pairs are not exactly corresponding (infeasible to be captured at the exactly same time). We first pre-process the data by the same scheme as aforementioned and then feed the US data into our network. Example synthetic MR results are shown in Fig. 8, in which we can see that the anatomical structures are well preserved by the synthesis process, with comparison to CycleGAN and the real MR data from the same subject (right side). Note that our trained model is directly applied to these data without any fine-tuning. In addition, we also present a quantitative comparison to other methods in Table IV. Specifically, a registration [43]

TABLE IV
QUANTITATIVE EVALUATION FOR US-TO-MR IMAGE SYNTHESIS ON REAL PSEUDO-PAIRED DATA.

Settings	AE	GAN	CycleGAN	Ours
SSIM \uparrow	0.1595 \pm 0.0127	0.2271 \pm 0.0702	0.6041 \pm 0.0483	0.6250\pm0.0586

TABLE V
QUANTITATIVE EVALUATION ON OUR SYNTHESISED MR IMAGES WITH COMPARISON TO THE 3D AUXILIARY APPROACHES. 3D-I, II, III ARE THE CORRESPONDING APPROACHES SHOWN IN FIG. 2 (RIGHT).

Settings	Ours (base)	Ours (3D-i)	Ours (3D-ii)	Ours (3D-iii)
Deformation \downarrow	0.46 \pm 0.24	0.59 \pm 0.33	0.44 \pm 0.19	0.60 \pm 0.37

is performed between the synthesised MR images and the real cropped MR images and the SSIM (Structural Similarity Index Measure) metric is used to measure the performance. We can see that the proposed method outperforms the other alternative solutions, which again validates the effectiveness of our approach.

F. 3D Auxiliary Analysis

When 3D volumetric US data is available, our synthesis framework can be easily adapted to leverage the additional information, as described in Section III-D. Here we analyse the effectiveness of including such auxiliary information. Specifically, we present the deformation score of the three approaches (see Fig. 2-right) with comparison to our 2D-based approach, in Table V. Qualitative results with one failure case (the third row) of our base model are also shown in Fig. 9. From the results, we can observe that the 3D-ii solution (i.e., predicting the slice index as an auxiliary task) performs the best among all the solutions. Although both augmenting with neighbouring slices (3D-i) and slice index prediction (3D-ii) provide additional 3D guidance for the target task, by feeding additional slices, ambiguity is also introduced for the discriminative prediction, which leads to slightly worse performance. On the other hand, the task of directly reasoning the slice index is based on the features from the latent space, which shares anatomical information with the two data modalities and is a less ambiguous task. Note that all these 3D-auxiliary solutions perform better than the other alternative architectures shown in Table II. The slice index reasoning task also performs quite well, with an accuracy of 86%. Through our experiment, we found that lower performance in this task leads to a less accurate image synthesis task, which on the other side validates the effectiveness of the slice index reasoning.

V. DISCUSSIONS

A. Potential Applications

a) *Annotation Transfer between US and MRI:* Data annotation from human experts is a long-standing challenge for data-driven medical image models, especially for those data (e.g., US) where anatomical structures are difficult to recognise. Thus, well-trained experts are necessary to annotate such data, which is labour-intensive work. Even though, the

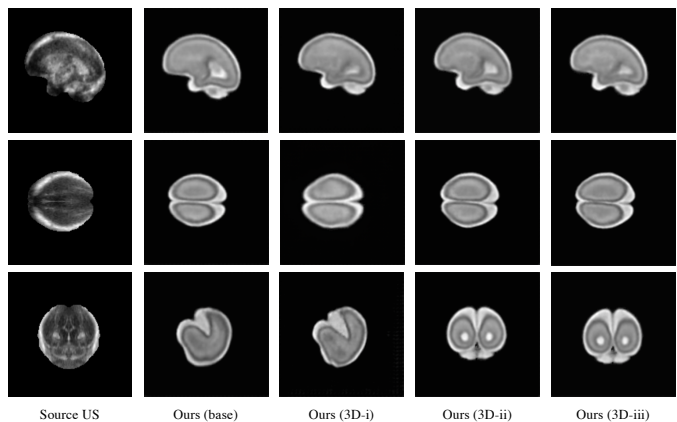


Fig. 9. Qualitative performance for US-to-MR image synthesis by 3D auxiliaries. Our base model without 3D auxiliaries is included for comparison.

accuracy of the annotation cannot be always guaranteed. On the other hand, if the non-anatomical data has an anatomical correspondence in which annotation is much easier, the labelling effort can be largely mitigated. By using the model proposed in this work, corresponding MRI data can be generated for each US image. Since the data annotation for MRI is much more efficient compared to that for US, the annotation work can be first done on MRI and then transferred to the US.

b) *Data Augmentation for Learning-Based Approaches:* A large amount of data is essential for deep learning models, whereas the acquisition of some data modalities is rather challenging. Fetal MRI is a kind of such data, where routinely MRI scan is usually not provided in most cases. As a result, data-driven deep learning models for fetal MRI related problems are infeasible to scale. In this case, if the correlation between fetal MRI and sufficient routine US scans can be bridged, training data for MRI can be generated. The proposed approach in this work may be well-suited for this. By feeding US images into our model, abundant corresponding MRI-like data can be synthesised to train deep learning models for fetal MRI.

B. Generalisation

The input to our model is pre-processed (cropped and aligned) US data instead of full US scans. Although promising synthesis results have been achieved, one limitation of our model lies in the input data assumptions. The generalisation to full US scans will require further research. Since the proposed method is data-driven and our model was trained with 23-week data, it is not an optimal model for other gestational ages. Similarly, as our model is trained on healthy data, it may not be the optimal model for new US images with real pathologies. This is a challenging problem for modern deep learning approaches and could be possibly addressed by transfer learning strategies. Without real paired data for evaluation, we cannot determine if the synthesised images are realistic enough to transfer all useful diagnostic information. Looking into this would be a natural next step towards assessing potential clinical utility. Finally, the proposed approach is a general framework that can be readily extended to other body parts (e.g., heart), medical imaging modalities (e.g., CT) and clinical application domains.

VI. CONCLUSION

In this paper, we have presented an original method to synthesise MR-like fetal brain images from unpaired US images, via a novel anatomy-aware self-supervised framework. Specifically, shared latent features between the two modalities (US and MR) are first extracted, from which the target MRI is synthesised under a group of anatomy-aware constraints. A cross-modal attention module is introduced to incorporate non-local guidance across the two modalities. An investigation to leverage 3D volumetric auxiliaries is also presented. Experimental results demonstrate the effectiveness of the proposed framework both qualitatively and quantitatively, with comparison to alternative CNN architectures.

We believe the proposed framework to be useful within analysis tasks such as the alignment between US and MRI and for communicating US findings to obstetricians and patients. The generalisation to full US images is another interesting direction worth investigation. Given more paired examples in the future, the synthesis quality would also be improved.

ACKNOWLEDGMENT

The authors would like to thank Andrew Zisserman for many helpful discussions, the volunteers for assessing images, and NVIDIA Corporation for the Titan V GPU donation. Ana Namburete is grateful for support from the UK Royal Academy of Engineering under its Engineering for Development Research Fellowships scheme.

REFERENCES

- [1] D. Pugash, P. C. Brugger, D. Bettelheim, and D. Prayer, "Prenatal ultrasound and fetal mri: the comparative value of each modality in prenatal diagnosis," *European journal of radiology*, vol. 68, no. 2, pp. 214–226, 2008.
- [2] D. Bulas and A. Egloff, "Benefits and risks of mri in pregnancy," in *Seminars in perinatology*, vol. 37, no. 5. Elsevier, 2013, pp. 301–304.
- [3] Y. Zhao *et al.*, "Towards mr-only radiotherapy treatment planning: Synthetic ct generation using multi-view deep convolutional neural networks," in *MICCAI*, 2018, pp. 286–294.
- [4] D. Nie *et al.*, "Medical image synthesis with context-aware generative adversarial networks," in *MICCAI*, 2017, pp. 417–425.
- [5] H. Yang *et al.*, "Unpaired brain mr-to-ct synthesis using a structure-constrained cyclegan," in *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*, 2018, pp. 174–182.
- [6] P. Costa *et al.*, "End-to-end adversarial retinal image synthesis," *IEEE TMI*, vol. 37, no. 3, pp. 781–791, 2018.
- [7] —, "Towards adversarial retinal image synthesis," *arXiv preprint arXiv:1701.08974*, 2017.
- [8] M. Kuklisova-Murgasova *et al.*, "Registration of 3d fetal neurosonography and mri," *Medical image analysis*, vol. 17, no. 8, pp. 1137–1150, 2013.
- [9] A. P. King *et al.*, "Registering preprocedure volumetric images with intraprocedure 3-d ultrasound using an ultrasound imaging model," *IEEE TMI*, vol. 29, no. 3, pp. 924–937, 2010.
- [10] Y. Berker *et al.*, "Mri-based attenuation correction for hybrid pet/mri systems: a 4-class tissue segmentation technique using a combined ultrashort-echo-time/dixon mri sequence," *Journal of nuclear medicine*, vol. 53, no. 5, p. 796, 2012.
- [11] G. Delpon *et al.*, "Comparison of automated atlas-based segmentation software for postoperative prostate cancer radiotherapy," *Frontiers in oncology*, vol. 6, p. 178, 2016.
- [12] J. Sjölund, D. Forsberg, M. Andersson, and H. Knutsson, "Generating patient specific pseudo-ct of the head from mr using atlas-based regression," *Physics in Medicine & Biology*, vol. 60, no. 2, p. 825, 2015.
- [13] C. Catana *et al.*, "Towards implementing an mr-based pet attenuation correction method for neurological studies on the mr-pet brain prototype," *Journal of nuclear medicine*, vol. 51, no. 9, p. 1431, 2010.
- [14] S. Roy, J. A. Butman, and D. L. Pham, "Synthesizing ct from ultrashort echo-time mr images via convolutional neural networks," in *International Workshop on Simulation and Synthesis in Medical Imaging*, 2017, pp. 24–32.
- [15] Z. Zhang, L. Yang, and Y. Zheng, "Translating and segmenting multimodal medical volumes with cycle-and shape-consistency generative adversarial network," in *CVPR*, 2018, pp. 9242–9251.
- [16] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *ICCV*, 2017.
- [17] J. Jiao, A. I. Namburete, A. T. Papageorghiou, and J. A. Noble, "Anatomy-aware self-supervised fetal mri synthesis from unpaired ultrasound images," in *International Workshop on Machine Learning in Medical Imaging*. Springer, 2019.
- [18] Y. K. Lee *et al.*, "Radiotherapy treatment planning of prostate cancer using magnetic resonance imaging alone," *Radiotherapy and oncology*, vol. 66, no. 2, pp. 203–216, 2003.
- [19] J. A. Dowling *et al.*, "An atlas-based electron density mapping method for magnetic resonance imaging (mri)-alone treatment planning and adaptive mri-based prostate radiation therapy," *International Journal of Radiation Oncology* Biology* Physics*, vol. 83, no. 1, pp. e5–e11, 2012.
- [20] I. Goodfellow *et al.*, "Generative adversarial nets," in *NeurIPS*, 2014, pp. 2672–2680.
- [21] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle, "Greedy layer-wise training of deep networks," in *NeurIPS*, 2007, pp. 153–160.
- [22] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P.-A. Manzagol, "Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion," *Journal of machine learning research*, vol. 11, no. Dec, pp. 3371–3408, 2010.
- [23] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [24] Z. Yi, H. Zhang, P. Tan, and M. Gong, "Dualgan: Unsupervised dual learning for image-to-image translation," in *ICCV*, 2017, pp. 2849–2857.
- [25] J. M. Wolterink, A. M. Dinkla, M. H. Savenije, P. R. Seevinck, C. A. van den Berg, and I. Išgum, "Deep mr to ct synthesis using unpaired data," in *International Workshop on Simulation and Synthesis in Medical Imaging*. Springer, 2017, pp. 14–23.
- [26] Y. Hiasa, Y. Otake, M. Takao, T. Matsuoka, K. Takashima, A. Carass, J. L. Prince, N. Sugano, and Y. Sato, "Cross-modality image synthesis from unpaired data using cyclegan," in *International workshop on simulation and synthesis in medical imaging*. Springer, 2018, pp. 31–41.
- [27] A. Chartsias, T. Joyce, R. Dharmakumar, and S. A. Tsafaris, "Adversarial image synthesis for unpaired multi-modal cardiac data," in *International workshop on simulation and synthesis in medical imaging*. Springer, 2017, pp. 3–13.
- [28] G. Zeng and G. Zheng, "Hybrid generative adversarial networks for deep mr to ct synthesis using unpaired data," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2019, pp. 759–767.
- [29] M. A. Maraci, R. Napolitano, A. Papageorghiou, and J. A. Noble, "Searching for structures of interest in an ultrasound video sequence," in *International Workshop on Machine Learning in Medical Imaging*. Springer, 2014, pp. 133–140.
- [30] M. Yaqub, B. Kelly, A. T. Papageorghiou, and J. A. Noble, "Guided random forests for identification of key fetal anatomy and image categorization in ultrasound scans," in *MICCAI*. Springer, 2015, pp. 687–694.
- [31] H. Chen *et al.*, "Standard plane localization in fetal ultrasound via domain transferred deep neural networks," *IEEE journal of biomedical and health informatics*, vol. 19, no. 5, pp. 1627–1636, 2015.
- [32] C. F. Baumgartner *et al.*, "Sononet: real-time detection and localisation of fetal standard scan planes in freehand ultrasound," *IEEE TMI*, vol. 36, no. 11, pp. 2204–2215, 2017.
- [33] Y. Cai, H. Sharma, P. Chatelain, and J. A. Noble, "Multi-task sononet: detection of fetal standardized planes assisted by generated sonographer attention maps," in *MICCAI*. Springer, 2018, pp. 871–879.
- [34] W. Wein, S. Brunke, A. Khamene, M. R. Callstrom, and N. Navab, "Automatic ct-ultrasound registration for diagnostic imaging and image-guided intervention," *Medical image analysis*, vol. 12, no. 5, pp. 577–585, 2008.
- [35] Y. Xiao *et al.*, "Evaluation of mri to ultrasound registration methods for brain shift correction: The curious2018 challenge," *arXiv preprint arXiv:1904.10535*, 2019.
- [36] A. I. Namburete, W. Xie, M. Yaqub, A. Zisserman, and J. A. Noble, "Fully-automated alignment of 3d fetal brain ultrasound to a canonical reference space using multi-task learning," *Medical image analysis*, vol. 46, pp. 1–14, 2018.

- [37] H. Zhang, I. Goodfellow, D. Metaxas, and A. Odena, "Self-attention generative adversarial networks," *arXiv preprint arXiv:1805.08318*, 2018.
- [38] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *CVPR*, 2018, pp. 7794–7803.
- [39] S. Zagoruyko and N. Komodakis, "Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer," *arXiv preprint arXiv:1612.03928*, 2016.
- [40] A. T. Papageorghiou *et al.*, "International standards for fetal growth based on serial ultrasound measurements: the fetal growth longitudinal study of the intergrowth-21st project," *The Lancet*, vol. 384, no. 9946, pp. 869–879, 2014.
- [41] A. Gholipour *et al.*, "Construction of a deformable spatiotemporal mri atlas of the fetal brain: evaluation of similarity metrics and deformation models," in *MICCAI*. Springer, 2014, pp. 292–299.
- [42] D. Rueckert, L. I. Sonoda, C. Hayes, D. L. Hill, M. O. Leach, and D. J. Hawkes, "Nonrigid registration using free-form deformations: application to breast mr images," *IEEE TMI*, vol. 18, no. 8, pp. 712–721, 1999.
- [43] I. Aganj, J. E. Iglesias, M. Reuter, M. R. Sabuncu, and B. Fischl, "Mid-space-independent deformable image registration," *NeuroImage*, vol. 152, pp. 158–170, 2017.