

MR Image Synthesis
ResViT

TMI

2022

ResViT: Residual vision transformers for
multi-modal medical image synthesisOnat Dalmaz, Mahmut Yurt, and Tolga Çukur*, *Senior Member*

Abstract—Generative adversarial models with convolutional neural network (CNN) backbones have recently been established as state-of-the-art in numerous medical image synthesis tasks. However, CNNs are designed to perform local processing with compact filters, and this inductive bias compromises learning of contextual features. Here, we propose a novel generative adversarial approach for medical image synthesis, ResViT, that leverages the contextual sensitivity of vision transformers along with the precision of convolution operators and realism of adversarial learning. ResViT's generator employs a central bottleneck comprising novel aggregated residual transformer (ART) blocks that synergistically combine residual convolutional and transformer modules. Residual connections in ART blocks promote diversity in captured representations, while a channel compression module distills task-relevant information. A weight sharing strategy is introduced among ART blocks to mitigate computational burden. A unified implementation is introduced to avoid the need to rebuild separate synthesis models for varying source-target modality configurations. Comprehensive demonstrations are performed for synthesizing missing sequences in multi-contrast MRI, and CT images from MRI. Our results indicate superiority of ResViT against competing CNN- and transformer-based methods in terms of qualitative observations and quantitative metrics.

Index Terms—medical image synthesis, transformer, residual, vision, adversarial, generative, unified

I. INTRODUCTION

Medical imaging plays a pivotal role in modern healthcare

nonlinear differences in tissue contrast across modalities [8]–[13]. Unsurprisingly, recent adoption of deep learning methods for solving this difficult problem has enabled major performance leaps [14]–[21]. In learning-based synthesis, network models effectively capture a prior on the joint distribution of source-target images [22]–[24]. Earlier studies using CNNs for this purpose reported significant improvements over traditional approaches [22], [23], [25]–[28]. Generative adversarial networks (GANs) were later introduced that leverage an adversarial loss to increase capture of detailed tissue structure [24], [29]–[35]. Further improvements were attained by leveraging enhanced architectural designs [36]–[39], and learning strategies [40]–[42]. Despite their prowess, prior learning-based synthesis models are fundamentally based on convolutional architectures that use compact filters to extract local image features [43], [44]. Exploiting correlations among small neighborhoods of image pixels, this inductive bias reduces the number of model parameters to facilitate learning. However, it also limits expressiveness for contextual features that reflect long-range spatial dependencies [45], [46].

Medical images contain contextual relationships across both healthy and pathological tissues. For instance, bone in the skull or CSF in the ventricles broadly distribute over spatially contiguous or segregated brain regions, resulting in dependencies among distant voxels. While pathological tissues have less regular anatomical priors, their spatial distribution (e.g., location, quantity, shape) can still show disease-specific

MR Image Synthesis

ResViT

Background & Goal

▶ Background

- CNN based models: Only learns local information and fails to utilize long-range contextual information
- Transformer based models: Leverages long-range contextual information, but has high computational cost

▶ Goal

- ResViT that translates between multi-modal imaging data combines the sensitivity of Vision Transformers to global context, the localization power of CNNs, and the realism of Adversarial learning

MR Image Synthesis ResViT

Network Architecture

▶ Generator subnetworks

- Encoder: To capture a hierarchy of localized features of source images using Convolutional layers
- Information Bottleneck: To distill task-relevant information in the encoded features using CNN + Transformer
- Decoder: To distill multi-modal images in separate channels using Transposed convolutional layers
- Parameter Sharing Transformers: To avoid inevitably elevate memory demand and risk of overfitting, weight sharing strategy is adopted where the model weights for the transformer encoder are tied across separate ART blocks.

▶ Discriminator subnetworks

- Discriminator is based on a conditional PatchGAN architecture

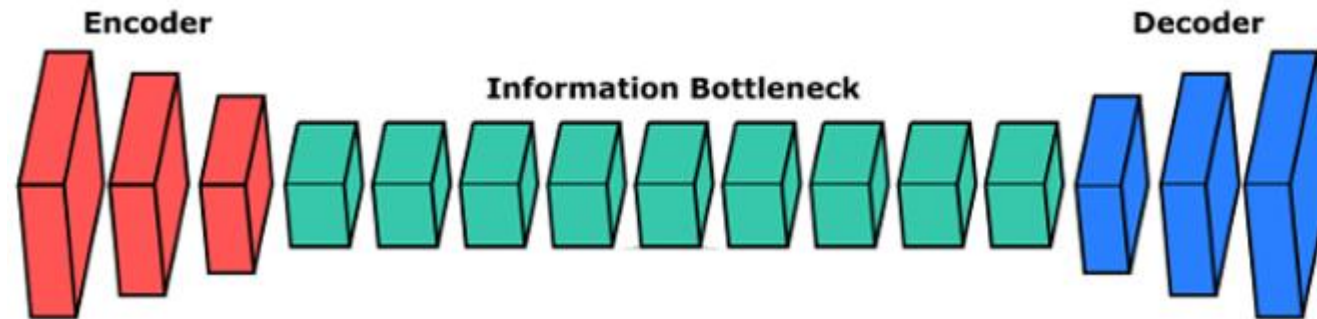


Fig 1. Overview of ResViT Generator

MR Image Synthesis ResViT

Network Architecture

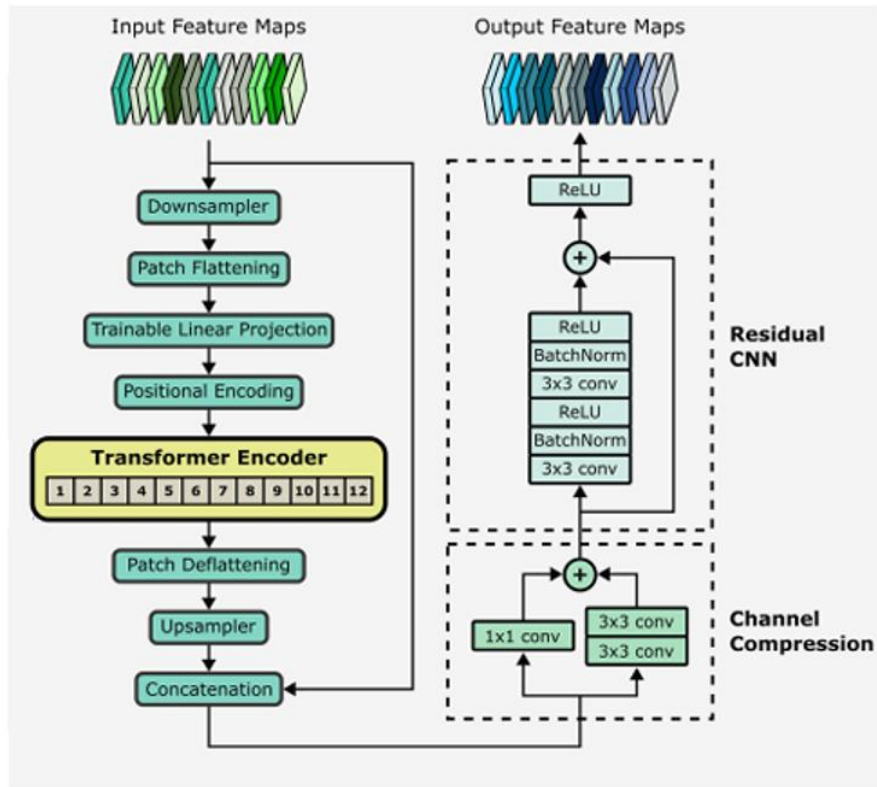


Fig 2. Overview of ART Block

Information Bottleneck

- To maintain both localization power and contextual sensitivity, ART blocks that aggregate information from residual convolutional and transformer branches

Residual Blocks

- The reason for using Residual Transformer and Residual CNN module is to learn a unified representation that considers local + contextual information together

Channel Compression

- The process of optimally combining the outputs of CNN and Transformer

MR Image Synthesis ResViT

Network Architecture

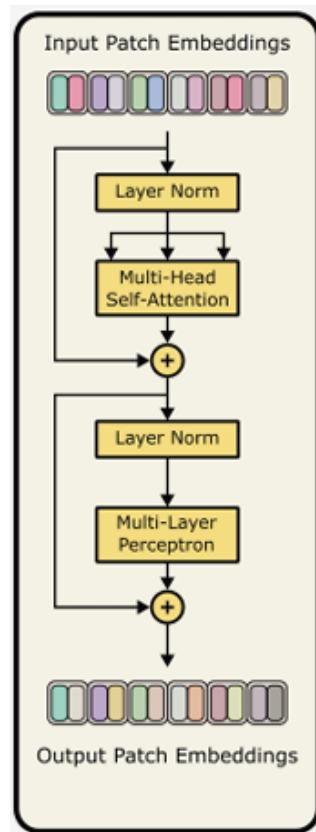


Fig 3. Overview of Transformer

▶ Transformer

- The role of learning long range context information

▶ Multi-Head Self Attention

- Learning the relationship between input features with multiple attention heads

▶ Multi-layer Perceptron

- Learn more useful representations by nonlinearly transforming the output of the transformer

MR Image Synthesis ResViT

Network Architecture

► Decoder

- Synthesize all contrasts within the multi-modal protocol regardless of the specific source-target configuration

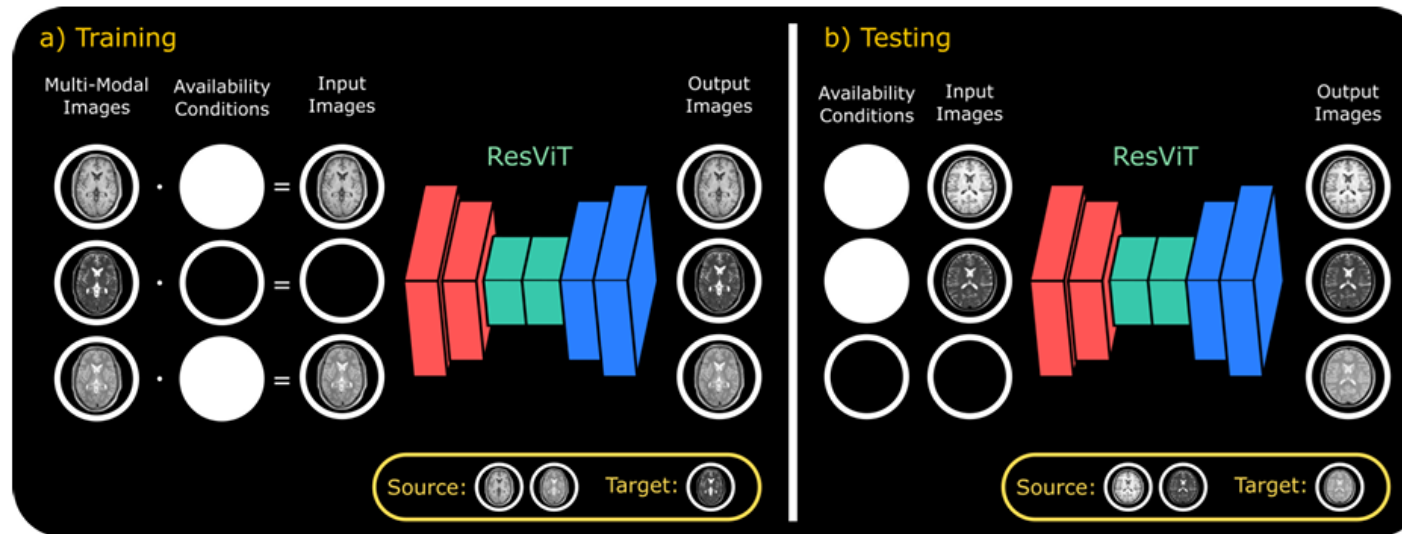


Fig 2. a) Training Process: ResViT takes as input the entire set of images within the multi-modal protocol, Multiple configurations of source-target modalities are expressed as availability conditions in ResViT. b) Inference Process: For each test subject, a specific source-target configuration is determined through availability conditions. This allows the model to perform appropriate modality transformation.

MR Image Synthesis

ResViT

Loss Functions

▶ Pixel-wise L1 loss

- Minimize the difference between the real target image and the synthetic image

▶ Pixel-wise consistency loss

- Minimize the difference between the real source image and the reconstructed source image

▶ Adversarial loss

- Adversarial loss is based on Least-Squares GAN loss

$$L_{pix} = \sum_{i=1}^I (1 - a_i) \mathbb{E}[\| (X^G)_i - m_i \|_1] \quad L_{rec} = \sum_{i=1}^I a_i \mathbb{E}[\| G(X^G)_i - m_i \|_1]$$

Eq 1. pixel-wise L1 loss

$$L_{rec} = \sum_{i=1}^I a_i \mathbb{E}[\| G(X^G)_i - m_i \|_1]$$

Eq 2. pixel-wise consistency loss

$$L_{adv} = -\mathbb{E}[D(X^D(acquired))^2] - \mathbb{E}[(D(X^D(synthetic)) - 1)^2]$$

Eq 3. adversarial loss

MR Image Synthesis

ResViT

Results

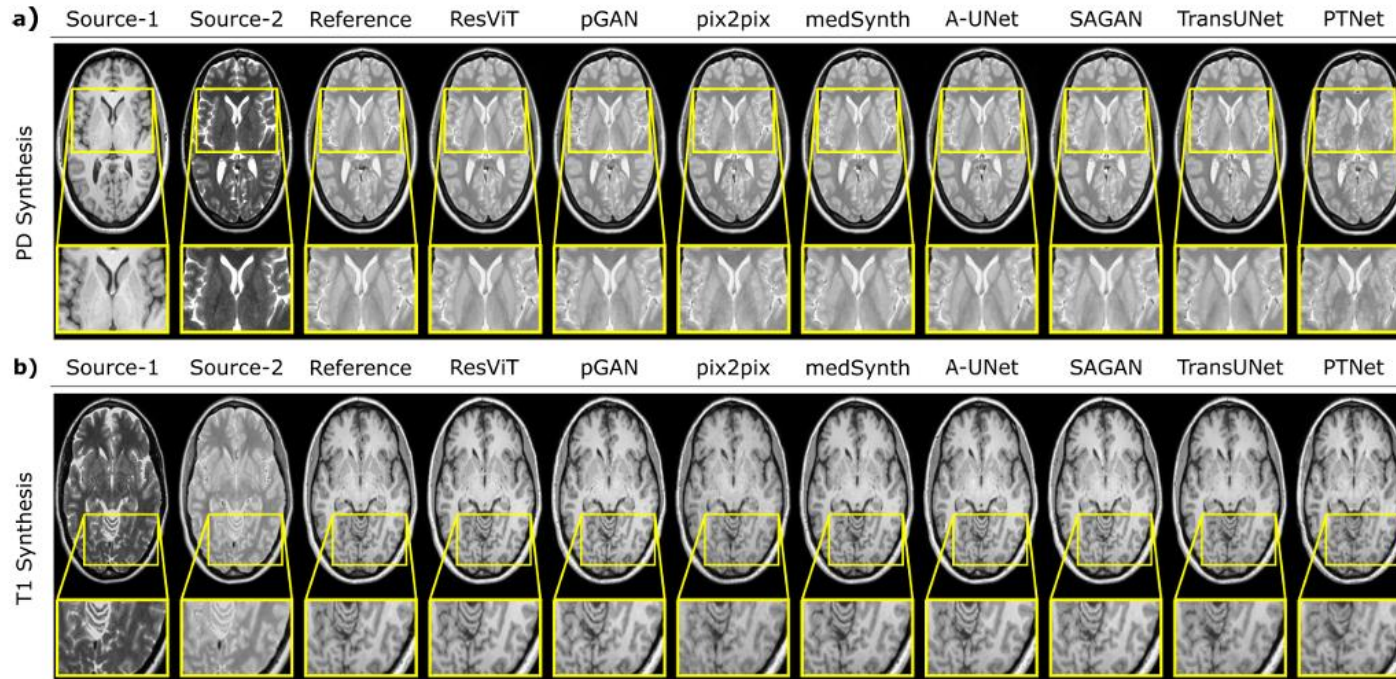


Fig 3. Demonstrated on the IXI dataset for two representative many-to-one synthesis tasks

| | T ₁ , T ₂ → PD | | T ₁ , PD → T ₂ | | T ₂ , PD → T ₁ | | T ₂ → PD | | PD → T ₂ | |
|-----------|--------------------------------------|--------------|--------------------------------------|--------------|--------------------------------------|--------------|---------------------|--------------|---------------------|--------------|
| | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| ResViT | 33.92 | 0.977 | 35.71 | 0.977 | 29.58 | 0.952 | 32.90 | 0.972 | 34.24 | 0.972 |
| | ±1.44 | ±0.004 | ±1.20 | ±0.005 | ±1.37 | ±0.011 | ±1.20 | ±0.005 | ±1.09 | ±0.005 |
| pGAN | 32.91 | 0.966 | 33.95 | 0.965 | 28.71 | 0.941 | 32.20 | 0.963 | 33.05 | 0.963 |
| | ±0.94 | ±0.005 | ±1.06 | ±0.006 | ±1.08 | ±0.013 | ±1.00 | ±0.005 | ±0.95 | ±0.007 |
| pix2pix | 32.25 | 0.974 | 33.62 | 0.973 | 28.35 | 0.949 | 30.72 | 0.956 | 30.74 | 0.950 |
| | ±1.24 | ±0.006 | ±1.31 | ±0.009 | ±1.24 | ±0.016 | ±1.28 | ±0.007 | ±1.63 | ±0.012 |
| medSynth | 33.23 | 0.967 | 32.66 | 0.963 | 28.43 | 0.938 | 32.20 | 0.964 | 30.41 | 0.956 |
| | ±1.09 | ±0.005 | ±1.30 | ±0.007 | ±1.01 | ±0.013 | ±1.10 | ±0.006 | ±3.98 | ±0.025 |
| A-UNet | 32.24 | 0.963 | 32.43 | 0.959 | 28.95 | 0.916 | 32.05 | 0.960 | 33.32 | 0.961 |
| | ±0.92 | ±0.014 | ±1.36 | ±0.007 | ±1.21 | ±0.013 | ±1.04 | ±0.009 | ±1.08 | ±0.007 |
| SAGAN | 32.50 | 0.964 | 33.71 | 0.965 | 28.62 | 0.942 | 32.07 | 0.963 | 32.96 | 0.962 |
| | ±0.93 | ±0.005 | ±1.00 | ±0.006 | ±1.10 | ±0.013 | ±0.98 | ±0.006 | ±1.01 | ±0.007 |
| TransUNet | 32.53 | 0.968 | 32.49 | 0.960 | 28.21 | 0.941 | 30.90 | 0.960 | 31.73 | 0.958 |
| | ±0.97 | ±0.005 | ±1.18 | ±0.008 | ±1.30 | ±0.013 | ±1.35 | ±0.006 | ±1.44 | ±0.008 |
| PTNet | 30.92 | 0.952 | 32.62 | 0.954 | 27.59 | 0.923 | 31.58 | 0.958 | 30.84 | 0.947 |
| | ±0.99 | ±0.006 | ±1.96 | ±0.019 | ±1.36 | ±0.021 | ±1.30 | ±0.007 | ±2.54 | ±0.033 |

Table 1. Performance of task-specific synthesis models in the IXI dataset

MR Image Synthesis ResViT

Results

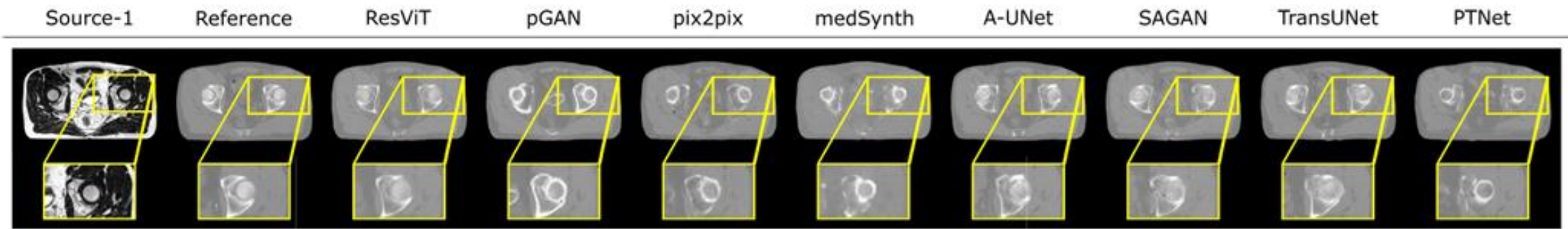


Fig 4. Demonstrated on the pelvic MRI-CT dataset for the T2-weighted MRI→CT task

| | | ResViT | pGAN | pix2pix | medSynth | A-UNet | SAGAN | TransUNet | PTNet |
|----------------|------|---------------|--------|---------|----------|--------|--------|-----------|--------|
| MRI ↑ CT | PSNR | 28.45 | 26.80 | 26.53 | 26.36 | 27.80 | 27.61 | 27.76 | 26.11 |
| | | ±1.35 | ±0.90 | ±0.45 | ±0.63 | ±0.63 | ±1.02 | ±1.03 | ±0.93 |
| | SSIM | 0.931 | 0.905 | 0.898 | 0.894 | 0.913 | 0.910 | 0.914 | 0.900 |
| | | ±0.009 | ±0.008 | ±0.004 | ±0.009 | ±0.004 | ±0.006 | ±0.009 | ±0.015 |

Table 2. Performance for the across-modality synthesis task(T2-w MRI → CT)

MR Image Synthesis ResViT

Results

| | T ₁ , T ₂ → PD | | | T ₁ , T ₂ → FLAIR | | | MRI → CT | | |
|------------|--------------------------------------|---------------|--------------|---|---------------|--------------|--------------|---------------|--------------|
| | PSNR | SSIM | FID | PSNR | SSIM | FID | PSNR | SSIM | FID |
| ResViT | 33.92 | 0.977 | 14.47 | 25.84 | 0.886 | 18.58 | 28.45 | 0.931 | 60.28 |
| | ±1.44 | ±0.004 | | ±1.13 | ±0.014 | | ±1.35 | ±0.009 | |
| w/o trans. | 32.91 | 0.966 | 14.56 | 24.96 | 0.868 | 19.21 | 26.73 | 0.899 | 95.38 |
| modules | ±0.96 | ±0.005 | | ±1.10 | ±0.005 | | ±0.91 | ±0.008 | |
| w/o conv. | 33.49 | 0.971 | 14.84 | 25.11 | 0.874 | 20.30 | 28.19 | 0.922 | 60.16 |
| modules | ±1.34 | ±0.005 | | ±1.02 | ±0.014 | | ±1.15 | ±0.009 | |
| w/o adv. | 33.75 | 0.977 | 15.80 | 22.95 | 0.891 | 40.68 | 28.58 | 0.932 | 65.49 |
| loss | ±1.45 | ±0.005 | | ±1.93 | ±0.015 | | ±1.13 | ±0.007 | |

Table 3. Performance of ResViT and Ablated of Transformer Modules, Convolutional Modules, or Adversarial loss

| | T ₁ , T ₂ → PD | | T ₁ , T ₂ → FLAIR | | MRI → CT | |
|------------------------------------|--------------------------------------|---------------|---|---------------|--------------|---------------|
| | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| ResViT | 33.92 | 0.977 | 25.84 | 0.886 | 28.45 | 0.931 |
| | ±1.44 | ±0.004 | ±1.13 | ±0.014 | ±1.35 | ±0.009 |
| w/o pre-training | 33.55 | 0.971 | 24.86 | 0.881 | 27.94 | 0.912 |
| | ±1.25 | ±0.005 | ±1.28 | ±0.016 | ±1.25 | ±0.009 |
| w/o del. insertion | 33.35 | 0.977 | 24.89 | 0.873 | 28.01 | 0.924 |
| | ±1.13 | ±0.004 | ±1.18 | ±0.015 | ±1.27 | ±0.008 |
| w/o pre-training or del. insertion | 33.58 | 0.971 | 24.74 | 0.869 | 27.66 | 0.913 |
| | ±1.16 | ±0.005 | ±1.30 | ±0.016 | ±0.78 | ±0.006 |

Table 5. Performance of ResViT and Ablated of pre-training and delayed insertion procedures for Transformers

| | T ₁ , T ₂ → PD | | T ₁ , T ₂ → FLAIR | | MRI → CT | |
|--|--------------------------------------|---------------|---|---------------|--------------|---------------|
| | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| A ₁ - A ₆ | 33.92 | 0.977 | 25.84 | 0.886 | 28.45 | 0.931 |
| | ±1.44 | ±0.004 | ±1.13 | ±0.014 | ±1.35 | ±0.009 |
| A ₁ - A ₆ (untied weights) | 33.72 | 0.973 | 25.19 | 0.879 | 28.16 | 0.923 |
| | ±1.23 | ±0.005 | ±1.18 | ±0.014 | ±1.11 | ±0.007 |
| A ₁ | 33.51 | 0.971 | 24.98 | 0.883 | 28.06 | 0.921 |
| | ±1.15 | ±0.005 | ±1.60 | ±0.015 | ±1.31 | ±0.008 |
| A ₆ | 33.78 | 0.977 | 25.25 | 0.880 | 27.95 | 0.921 |
| | ±1.34 | ±0.004 | ±1.20 | ±0.014 | ±1.22 | ±0.008 |

Table 4. Performance of ResViT and Ablated of weight tying and individual transformer modules

| | T ₁ , T ₂ → PD | | T ₁ , T ₂ → FLAIR | | MRI → CT | |
|------------------------------------|--------------------------------------|---------------|---|---------------|--------------|---------------|
| | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| ResViT | 33.92 | 0.977 | 25.84 | 0.886 | 28.45 | 0.931 |
| | ±1.44 | ±0.004 | ±1.13 | ±0.014 | ±1.35 | ±0.009 |
| w/o skip around conv. modules | 28.24 | 0.942 | 25.02 | 0.864 | 26.94 | 0.906 |
| | ±1.27 | ±0.009 | ±0.98 | ±0.016 | ±0.73 | ±0.007 |
| w/o skip around trans. modules | 31.53 | 0.962 | 24.06 | 0.868 | 27.08 | 0.908 |
| | ±1.26 | ±0.006 | ±1.28 | ±0.014 | ±0.80 | ±0.006 |
| ART with unlearned down/upsampling | 33.73 | 0.969 | 25.33 | 0.884 | 28.16 | 0.931 |
| | ±1.19 | ±0.005 | ±1.11 | ±0.014 | ±1.04 | ±0.007 |
| ART w/o down/upsampling | 31.51 | 0.961 | 23.61 | 0.867 | 26.79 | 0.915 |
| | ±1.27 | ±0.006 | ±1.53 | ±0.015 | ±0.62 | ±0.006 |

Table 6. Performance of ResViT and Variants ablated

MR Image Synthesis

ResViT

Implications & Limitations

► Implications

- Uniquely introduce many-to-one synthesis models and a unified model that generalizes across multiple source-target configurations
- Propose a hybrid architecture that combines localization capabilities of CNNs with contextual sensitivity of transformers

► Limitations

- Learn the distribution of the target modality implicitly, without explicitly evaluating the likelihood → Mode collapse
- Lack of reliability in network mapping because images are generated using the "One-Shot Sampling" method