

PADM: A Physics-Aware Diffusion Model for Attenuation Correction (Supplementary Material)

Anonymous WACV Algorithms Track submission

Paper ID 3311

001 1. Experiments Details

002 1.1. System & Hyper-parameters

003 The experiments are conducted on a system featuring an
004 AMD EPYC 7402P CPU and an NVIDIA A6000 GPU. In
005 this section, we also provide further implementation details
006 for PADM and other methods.

007 1.1.1. Our proposed model - PADM

008 We implement PADM as described in Section 4. In con-
009 trast to the original BBDM setup, we do not embed the
010 preprocessed NAC and AC images into the latent space of
011 VQGAN, since their relatively low resolution diminishes
012 the benefit of such compression. Instead, diffusion is di-
013 rectly performed in the image space. The Brownian Bridge
014 timesteps are set to 500 for training and reduced to 10 for
015 inference, striking a balance between reconstruction quality
016 and computational efficiency.

017 1.1.2. Pix2Pix

018 Pix2Pix [2] is a conditional GAN framework designed for
019 paired image-to-image translation. In our setup, the gen-
020 erator adopts a ResNet backbone with 9 residual blocks,
021 while the discriminator follows the 70×70 PatchGAN de-
022 sign [2]. The model is optimized with Adam ($\beta_1 = 0.5$)
023 using the least-squares adversarial loss (LSGAN). Instance
024 normalization is consistently applied to both the generator
025 and discriminator. Training runs for 100 epochs at a fixed
026 learning rate of 2×10^{-4} , followed by a linear decay across
027 an additional 100 epochs. To improve stability, a buffer of
028 50 previously generated images is maintained. Our imple-
029 mentation is adapted from the public repository at <https://github.com/junyanz/pytorch-CycleGAN-and-pix2pix>.
030
031

032 1.1.3. Palette

033 Palette [6] is an image-to-image diffusion framework that
034 adopts an exponential moving average (EMA) strategy from
035 the very first training step, with a decay coefficient of
036 0.9999 to ensure stability. The model is optimized using

Table 1. Training and optimization hyperparameters for PADM.

Component	Configuration
Denoising Network	
Channels	128
Channel multipliers	1, 4, 8
Residual blocks per downsample	2
Attention resolutions	32, 16, 8
# of Trainable Parameters	237.11M
EMA	
Start step	30,000
Decay	0.995
Update interval	16
Batch size	8
Learning Rate Scheduler (ReduceLROnPlateau)	
Max learning rate	1.0×10^{-4}
Min learning rate	5.0×10^{-7}
Factor	0.5
Patience	3,000
Cooldown	3,000
Threshold	1.0×10^{-4}

Adam with a fixed learning rate of 5×10^{-5} and no weight
decay. A linear beta schedule is applied over 1000 training
timesteps, starting from 1×10^{-6} and ramping up to 0.01.
During inference, sampling is carried out with 100 steps.
Our implementation follows the official open-source re-
lease available at <https://github.com/Janspiry/Palette-Image-to-Image-Diffusion-Models>

044 1.1.4. BBDM

045 BBDM [4] formulates image-to-image translation as a dif-
046 fusion process defined by a stochastic Brownian bridge in
047 the VQ-GAN latent domain. For experimental consistency,
048 no pretrained VQ-GAN is used for initialization. The net-
049 work is optimized using Adam with a fixed learning rate of

050 1×10^{-4} , $\beta_1 = 0.9$, and $\beta_2 = 0.999$. A learning rate
051 schedule is applied, consisting of a 3000-step cooldown,
052 multiplicative decay of 0.5, and a minimum threshold of
053 5×10^{-7} . To stabilize training, an exponential moving ave-
054 rage (EMA) with decay factor 0.995 is introduced after
055 30k iterations. The diffusion process employs 500 steps
056 during training, while inference leverages DDIM [7] with
057 only 10 steps for efficiency. All experiments are imple-
058 mented in accordance with the official codebase at <https://github.com/xuekt98/BBDM.git>.
059

060 1.1.5. ResViT

061 ResViT [1] is a GAN-based framework for medical image
062 synthesis, featuring a generator with Aggregated Residual
063 Transformer (ART) blocks that unify convolutional and
064 transformer modules to jointly exploit local detail and
065 global context. The training of ResViT was performed in
066 two stages. In the first stage, we pre-trained the CNN
067 backbone (ART blocks without transformers) to learn local
068 structural representations. This was implemented with
069 learning rate of 0.0002, and 100 epochs (50 epochs with
070 fixed learning rate followed by 50 epochs of linear de-
071 cay). The resulting CNN weights were saved as initial-
072 ization for the second stage. In the fine-tuning stage, the
073 full ResViT model, including both convolutional and trans-
074 former modules. The network was initialized with the pre-
075 trained CNN backbone, while both transformer and resid-
076 ual modules were jointly optimized. Fine-tuning was per-
077 formed for 50 epochs (25 epochs with fixed learning rate
078 followed by 25 epochs of linear decay) with a higher learn-
079 ing rate of 0.001 to accelerate convergence. Our implemen-
080 tation is adapted from the official repository at <https://github.com/icon-lab/ResViT>.
081

082 1.1.6. Reg-GAN

083 RegGAN [3] is a GAN-based framework designed for med-
084 ical image-to-image translation. In our setup, all input im-
085 ages are normalized to the range $[-1, 1]$ and resampled to
086 256×256 . The generator follows the CycleGAN backbone
087 with two downsampling convolution blocks, nine residual
088 blocks, and two upsampling deconvolution blocks, while
089 the discriminator consists of four convolutional layers. A
090 U-Net-based registration network is incorporated to esti-
091 mate deformation fields for noise correction. The model
092 is trained using the Adam optimizer with learning rate
093 1×10^{-4} , $\beta_1 = 0.5$, $\beta_2 = 0.999$, and weight decay
094 1×10^{-4} , with batch size set to 1. Training is conducted
095 for 80 epochs (approximately 640K iterations). The overall
096 objective combines multiple loss terms with fixed weights:
097 adversarial loss ($\lambda_{adv} = 1$), L1 loss ($\lambda_{L1} = 100$), cycle-
098 consistency loss ($\lambda_{cyc} = 10$), correction loss ($\lambda_{corr} = 20$),
099 and smoothness loss ($\lambda_{smooth} = 10$). Our implemen-
100 tation is adapted from the official repository at <https://github.com/Kid-Liet/Reg-GAN>.
101

1.1.7. UNIT

102 UNIT [5] is an unsupervised image-to-image translation
103 framework that combines variational autoencoders (VAEs)
104 and generative adversarial networks (GANs) under the as-
105 sumption of a shared latent space between domains. The
106 model employs the Adam optimizer with a learning rate
107 of 1×10^{-4} and momentum parameters $\beta_1 = 0.5$ and
108 $\beta_2 = 0.999$, with no weight decay. Training alternates
109 between updating the generator and discriminator networks
110 to maintain stability. The architecture consists of domain-
111 specific encoders and decoders linked by a shared latent
112 representation, enabling bidirectional translation. Training
113 is performed with a batch size of 16 over 150 epochs. Infer-
114 ence is achieved by encoding an input image into the shared
115 latent space and decoding it into the target domain. The im-
116 plementation is based on the official UNIT code release at
117 <https://github.com/mingyuliutw/UNIT>
118

1.2. Detail of the Evaluation Metrics

119 We utilize Root Mean Squared Error (RMSE), Structural
120 Similarity Index Measure (SSIM), and Peak Signal-to-
121 Noise Ratio (PSNR) metrics to assess the quality of gen-
122 erated AC images. RMSE measures the average error, indi-
123 cating accuracy. SSIM reflects visual quality by assessing
124 changes in the image structure, whereas PSNR gauges im-
125 age fidelity.
126

1.2.1. Root Mean Squared Error (RMSE)

127 The Root Mean Squared Error (RMSE) measures the square
128 root of the average of the squared differences between the
129 predicted values and the actual values. RMSE penalizes
130 larger errors more strongly than MAE and is widely used
131 to evaluate regression and reconstruction performance in
132 machine learning and computer vision. A lower RMSE in-
133 dicates better model accuracy. The formula for RMSE is
134 given by:
135

$$136 \text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

137 where y_i is the ground-truth value, \hat{y}_i is the predicted value,
138 and n is the number of samples.

1.2.2. Structural Similarity Index Measurement (SSIM)

139 SSIM is a perceptual metric that quantifies how similar two
140 images are in terms of visual quality. It evaluates image
141 fidelity using three components: (i) luminance, (ii) contrast,
142 and (iii) structural information. The score ranges from -1 to
143 1 , where a value of 1 indicates perfect similarity. Models
144 achieving higher SSIM values produce outputs that more
145 closely resemble the reference images. The formulation of
146 SSIM is expressed as:
147

$$148 \text{SSIM}(x, y) = \frac{(2\mu_x\mu_y + c_1)(2\sigma_{xy} + c_2)}{(\mu_x^2 + \mu_y^2 + c_1)(\sigma_x^2 + \sigma_y^2 + c_2)}$$

149 where μ_x and μ_y are the mean values of x and y , σ_x^2 and σ_y^2
150 are the variances of x and y , σ_{xy} is the covariance of x and
151 y , and c_1 and c_2 are constants to stabilize the division.

152 1.2.3. Peak Signal-to-Noise Ratio (PSNR)

153 This metric quantifies the proportion between the peak possible
154 signal power and the noise power that distorts the signal, and it is usually reported in decibels (dB). In image
155 quality assessment, a larger PSNR value reflects superior fidelity. When two images are exactly the same, the PSNR
156 tends toward infinity. As a rule of thumb, a PSNR of 40 dB
157 or above indicates that visual differences between the original and reconstructed images are barely perceptible to the
158 human eye. The PSNR is formally defined as:
159
160
161

$$162 \text{PSNR} = 10 \cdot \log_{10} \left(\frac{\text{MAX}^2}{\text{MSE}} \right)$$

163 where MAX is the maximum possible pixel value of the image and MSE is the mean squared error between the original
164 and generated images.
165

166 2. Additional Qualitative Results

167 Supplementary results on cardiac SPECT slices are illustrated in Figures 1. Alongside each visualization, error
168 maps are displayed to reveal deviations from the ground truth and facilitate clearer method comparisons.
169
170

171 References

- 172 [1] Onat Dalmaz, Mahmut Yurt, and Tolga Cukur. Resvit: Residual vision transformers for multimodal medical image synthesis. *IEEE Transactions on Medical Imaging*, 41(10): 2598–2614, 2022. 2
- 173
174
175
176 [2] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1125–1134, 2017. 1
- 177
178
179
180
181 [3] Lingke Kong, Chenyu Lian, Detian Huang, Yanle Hu, Qichao Zhou, et al. Breaking the dilemma of medical image-to-image translation. *Advances in Neural Information Processing Systems*, 34:1964–1978, 2021. 2
- 182
183
184
185 [4] Bo Li, Kaitao Xue, Bin Liu, and Yu-Kun Lai. BBDM: Image-to-image translation with brownian bridge diffusion models. In *Proceedings of the 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1952–1961, 2023. 1
- 186
187
188
189
190 [5] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. *Advances in neural information processing systems*, 30, 2017. 2
- 191
192
193 [6] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *Proceedings of the 2022 ACM Special Interest Group on Computer Graphics and Interactive Techniques Conference*, pages 1–10, 2022. 1
- 194
195
196
197
198

- [7] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *Computing Research Repository arXiv Preprints arXiv:2010.02502*, 2020. 2

199
200
201

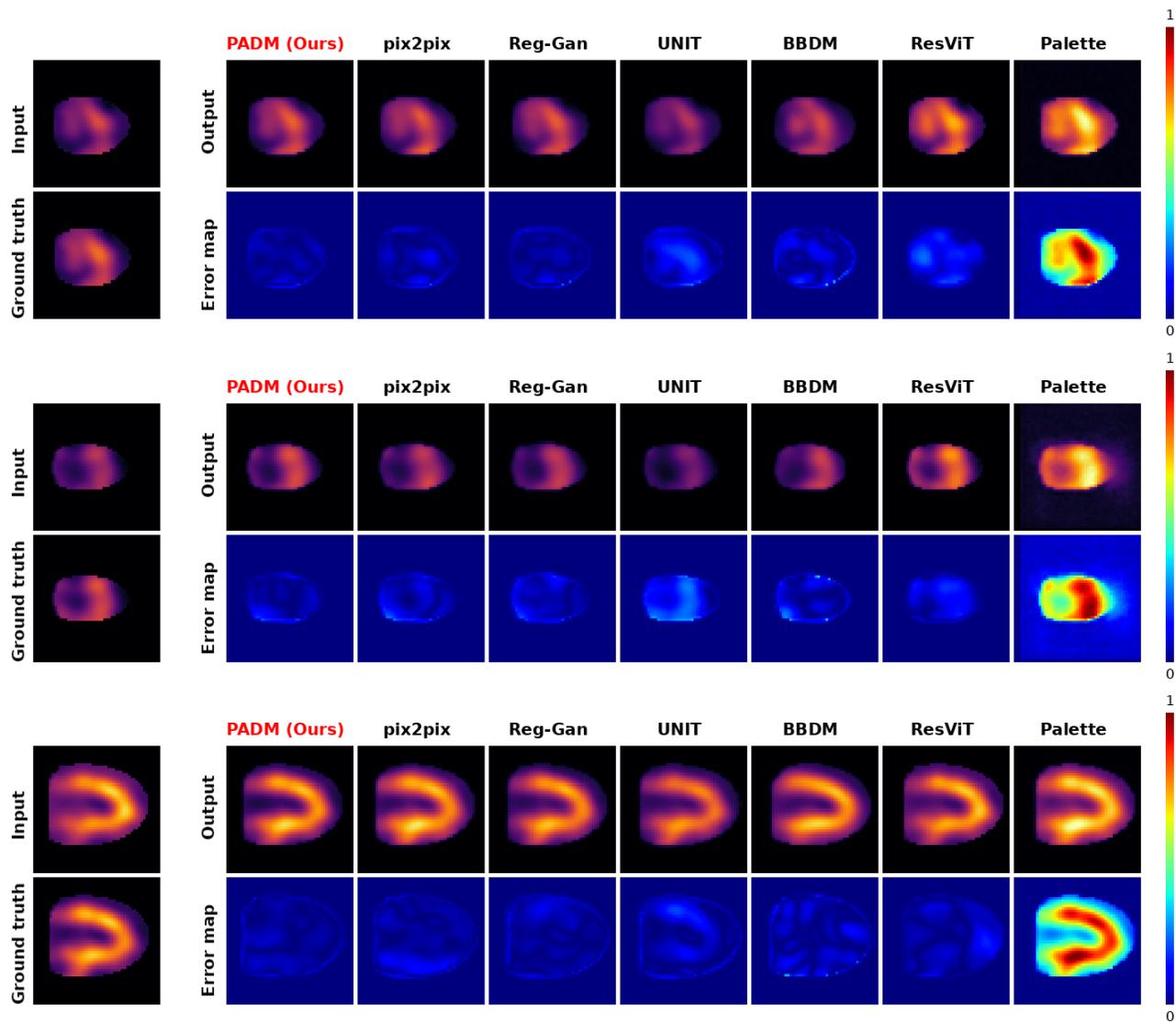


Figure 1. Qualitative comparison of reconstructed images across three standard views: horizontal long axis (top), short axis (middle), and vertical long axis (bottom).