# An Attention-based Multi-Scale Feature Learning Network for Multimodal Medical Image Fusion

Meng Zhou, Xiaolan Xu, and Yuxuan Zhang

**Abstract**—Medical images play an important role in clinical applications. Multimodal medical images could provide rich information about patients for physicians to diagnose. The image fusion technique is able to synthesize complementary information from multimodal images into a single image. This technique will prevent radiologists switch back and forth between different images and save lots of time in the diagnostic process. In this paper, we introduce a novel Dilated Residual Attention Network for the medical image fusion task. Our network is capable to extract multi-scale deep semantic features. Furthermore, we propose a novel fixed fusion strategy termed Softmax-based weighted strategy based on the Softmax weights and matrix nuclear norm. Extensive experiments show our proposed network and fusion strategy exceed the state-of-the-art performance compared with reference image fusion methods on four commonly used fusion metrics. Our code is available at https://github.com/simonZhou86/dilran.

**Index Terms**—Medical Image Fusion, Computer Vision, Attention Mechanism, Residual Mechanism, Dilated Convolution, Multi-scale

✦

## 1 INTRODUCTION

MEDICAL imaging plays an increasingly prominent role in clinical diagnosis. Multimodal medical images can provide information from various aspects, thus helping physicians confirm the diagnosis and make decisions regarding future treatment. For example, Magnetic Resonance Imaging (MRI) images provide the soft-tissue structures of the body; computerized tomography (CT) scans provide bone structure and high-density tissue information; Positron Emission Tomography (PET), and Single-Photon Emission Computed Tomography (SPECT) images can show the metabolic activity of the cells of tissues. To acquire adequate information, physicians must analyze many different medical images, which is time-consuming and laborious. Multimodal medical image fusion (MMIF) can merge the complementary information of original images and present the required information in one fused image. This is clinically significant because physicians can now access more detailed information about disease-related changes, thus providing patients with more comprehensive and precise medical treatment and support.

In recent years, many deep learning-based approaches have been proven successful in image fusion. These methods can extract features from input images and construct a fused image with the information needed. For instance, Wang et al. [1] introduced a CNN-based medical image fusion algorithm. This method employed the trained Siamese convolutional network to fuse the pixel activity information and implemented a contrast pyramid to decompose the source images. For more generalized image fusion, Xu

et al. [2] proposed an unsupervised and unified densely connected network, FusionDN, which is capable of various kinds of image fusion tasks.

In this work, we propose a novel end-to-end feature learning framework for multimodal medical image fusion with clear edge information and detailed textures. To be more precise, we focus on anatomical (CT) and functional (MRI) image fusion. The framework contains a feature extractor, a fixed fusion strategy that does not involve any adjustable parameters, and image reconstruction. Our approach combines input images and generates a fused image with more detailed textures and less information loss. Our fusion results exceed the state-of-the-art performance both qualitatively in subjective vision and quantitatively in multiple fusion metrics.

The main contributions of this paper are summarized as follows:

1) We propose a novel Dilated Residual Attention Network (DILRAN) for the feature extraction module. DILRAN coalesces the advantages of the residual attention network, the pyramid attention network, and the dilated convolutions. The proposed network has a faster convergence speed and can extract multi-scale deep semantic features.
2) We introduce a novel fusion strategy, Softmax Feature Weighted Strategy, which achieves a good result and outperforms other fusion strategies.
3) Extensive experiments show our proposed framework and fusion strategy exceed the state-of-the-art performance based on objective fusion metrics and subjective image quality.

---

- Meng Zhou is with the Department of Computer Science, University of Toronto, and The Hospital for Sick Children, Toronto, Canada.
  E-mail: simonzhou@cs.toronto.edu
- Xiaolan Xu and Yuxuan Zhang are with the Department of Computer Science, University of Toronto, Canada.
  E-mail: {landyxu,yuxuan}@cs.toronto.edu

## 2 RELATED WORK

There are several existing approaches for imaging fusion. Most traditional image fusions are based on the transform

domain and are at pixel, feature and decision-level [3]. The usage of fuzzy logic and neural networks with multi-scale decomposition can handle uncertainties and improve efficiency for fusion images. Other iconic works for utilizing traditional image processing algorithms for fusion tasks include Possion Image Editing [4], cross bilateral filter [5], and non-subsampled contourlet transform [6]. Tian et al. [7] proposed an improved version of pulse coupled neural network (PCNN) to manage NSCT sub-images using a shallow learning approach. This new PCNN calculation determines the linking strength parameter using the local area singular value decomposition of the structural information factor.

Recent advances in deep learning (DL) have led to the successful use of convolutional neural networks (CNN) in various imaging tasks such as classification [8], [9] and image super-resolution [10]. CNN has the ability to capture and extract features from images, which can then be used for image reconstruction. For example, Liu et al. [11] applied a CNN to multi-focus image fusion, using the network to generate a weight map of pixel activity during fusion. Hermessi et al. [12] introduced a CNN-based method in the shearlet domain for extracting feature maps of high-frequency fusion. Li et al. [13] proposed DenseFuse, a deep learning architecture for infrared and visible image fusion that is trained using dense blocks. Due to its limitation for only working at a single scale, Song et al. [14] proposed a multi-scale DenseNet (MSDNet) to overcome it. They encode the multi-scale mechanism with three filters of varying sizes for effectively capturing features at different scales. Increasing the width of the encoder network can also improve the amount of detail in the fused image. Despite this, DL has great potential in the field of medical artificial intelligence due to its ability to fit complex data, its ability to learn automatically, and its multitask adaptability.

A multi-generator multi-discriminator conditional generative adversarial network is presented by Huang et al. [15] to fuse functional information and structural information including texture details and dense structure information. This network presents better visual effects and also preserves the approximate maximum amount of information in several MMIF datasets. Fu et al. [16] introduced the residual pyramid attention structure: MSPRAN, which combines the advantages of residual attention and pyramid attention mechanisms in the fusion task. The framework extracts and keeps more information than a single residual attention or pyramid attention mechanism as the number of layers increases and maintains better deep features and expression capabilities. Another network structure (MSDRA) on double residual attention [17] combines a residual network and attention to acquire important detailed features while avoiding network gradient vanish or explosion.

## 3 PROPOSED METHOD

In this section, we provide a thorough discussion of the fusion framework we proposed and the corresponding loss function.

### 3.1 Overview

Figure 1 summarizes the overall pipeline. We introduce a novel end-to-end framework for medical image fusion. Our proposed fusion algorithm consists of a feature extractor, a fusion module, and a reconstruction module. The feature extractor aims to formulate deep semantic features of input images ($I_1$, $I_2$), and then use them as inputs to the fusion module. We introduce Dilated Residual Attention Network (DILRAN) for the feature extraction module. Next, the fusion module aims to fuse the two extracted feature maps into one map containing features from both original feature maps. The module determines for a certain pixel in the fused feature map, whether it comes from $I_1$, $I_2$, or both. The fusion rule can have many possible solutions, i.e., $max(I_1^i, I_2^i)$ for the $i$th pixel or the weighted average of both. Finally, the reconstruction module aims to reconstruct the fused image from the output of the fusion module. A sequence of convolutional layers is used when reconstructing the final fused image. Section 3.3, 3.4, and 3.5 provide an in-depth explanation of each module we mentioned above.

### 3.2 Dilated Residual Attention Network

The design principle of DILRAN is inspired by two state-of-the-art mechanisms: residual attention [9] and pyramid attention [18], as well as a recent paper [16] that relies on the residual pyramid attention mechanism. The classical residual attention network for image classification tasks [9], which is realized by imposing the attention mechanism in the residual network architecture. Figure 2 shows an example of the residual attention network. The top network path is called the trunk branch, which consists of a sequence of convolutional layers for feature processing. The bottom path is called the soft mask branch, which is formed by a downsample-upsample block with an activation function. The combination of the trunk and the soft mask branch build up the residual attention network following Equation (1).
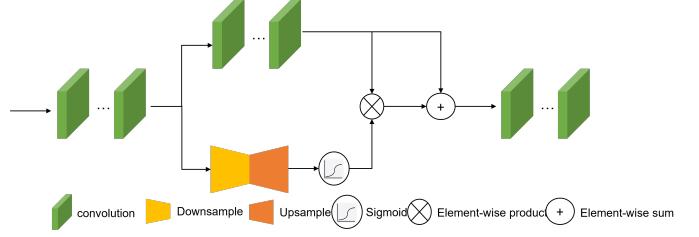


Fig. 2: Residual attention network architecture.

$$F(x) = (1 + S(x)) * T(x) \qquad (1)$$

Where $x$ is the feature map from the previous layer, $F(\cdot)$ is the output feature map, $S(\cdot)$ is the function for the soft mask branch, and $T(\cdot)$ is the convolution operation.

The advantage of the residual attention network is that it has a faster convergence speed without gradient vanishing or explosion, however, it may not be able to extract and learn deep features. Hence, the other mechanism we adopt is the pyramid attention network [18]. To be more specific, we focus on the feature pyramid attention that could provide deeper semantic features and attention advantages to the output feature map. A sample pyramid attention network architecture is shown in Figure 3. The convolution block usually contains a sequence of convolutional layers that
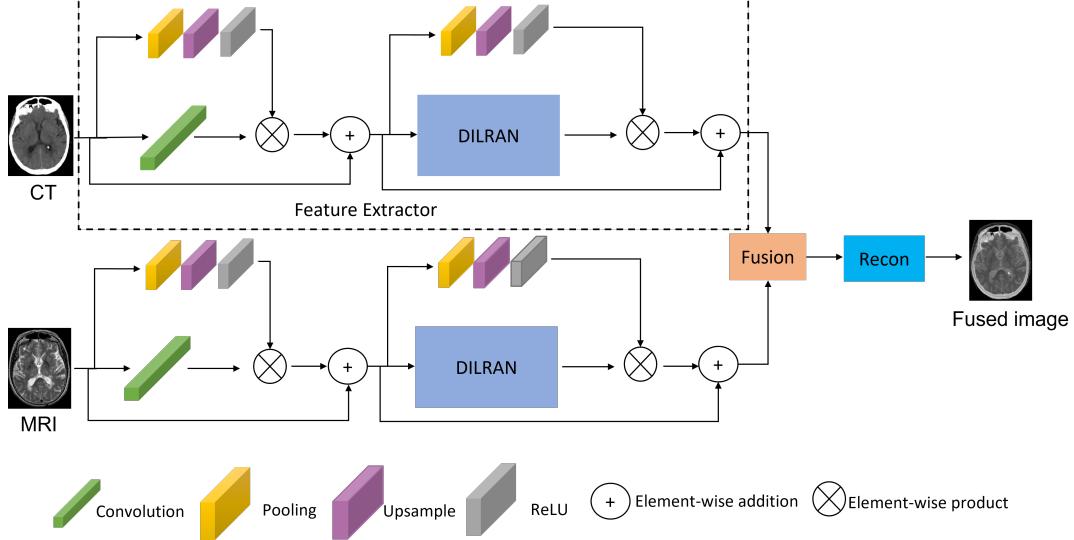
Fig. 1: The overall pipeline of the proposed method, which consists of a feature extractor, a fuser, and an image reconstructor

capture features at different scales and receptive fields. We adapt the idea to replace convolutions with larger kernel filters with a sequence of convolutions with smaller kernel filters [19]. Hence, the $CB1$ in our method is a single $3 \times 3$ convolutions; the $CB2$ consists of *two* $3 \times 3$ convolutions (represents a single $5 \times 5$ convolutions), and the $CB3$ consists of *three* $3 \times 3$ convolutions (represents a single $7 \times 7$ convolutions). The output features can be computed by the following Equation (2).
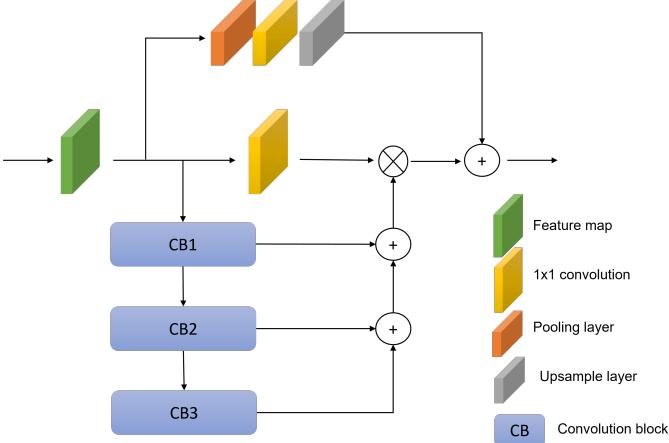


Fig. 3: Pyramid attention network architecture.

$$F(x) = (1 + P_1(P_2(P_3(x)))) * C(x) \qquad (2)$$

Where $x$ is the feature map from the previous layer, $F(\cdot)$ is the output feature map, $P_i(\cdot), i \in [1, 2, 3]$ is the parameters of the corresponding $CB_i, i \in [1, 2, 3]$ in the feature pyramid attention network, $C(\cdot)$ is the convolution operation.

In [16], they sequentially downsample the input feature maps to formulate the pyramid structure, and upsample again when performing the element-wise summation as

shown in Figure 3. However, this may lose the local information and fine details in the image. To solve this problem, we leverage the $\{1, 3, 5\}$-dilated convolution [20] on shallow features of the original input image to extract the multi-scale information instead of downsampling the feature map. The receptive field is expanded using three different dilated convolutions to improve the discriminative multi-scale feature extraction ability of the model. Once the multi-scale features are extracted, we concatenate those features channel-wise, and then the residual-pyramid attention paradigm is used to further extract deep features. The deep features are the output of the feature extraction module and are used in both the fusion and reconstruction modules. Our designed DILRAN architecture is shown in Figure 4.
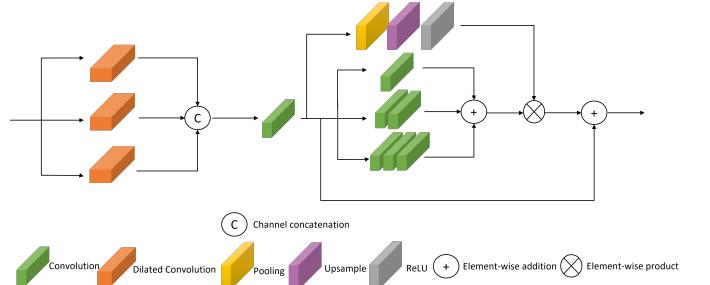


Fig. 4: Proposed Dilated Residual Attention Network architecture. Three pyramid convolution blocks contain one, two, and three consecutive $3 \times 3$ convolutional layers, respectively.

### 3.3 Feature Extraction Module

The detailed architecture for the feature extractor is shown in Figure 1. The input image is first passed to a residual attention network with a $1 \times 1$ convolutional layer to obtain a 64-dimensional shallow feature map. Then, the shallow feature is passed to the DILRAN module to obtain deep semantic features. To ensure the final output feature map

contains both the local and global information of the original input image, and to stabilize the training process, another residual attention network is utilized between shallow and deep semantic features. Finally, the output feature map is obtained and used in the following modules.

### 3.4 Fusion Strategy

The fusion strategy in the fusion module is used to fuse the extracted features of input images into a single feature map. In this section, We introduce a novel fusion strategy termed "Softmax Feature Weighted Strategy" that exceeds the state-of-the-art performance. The advantage of the proposed fusion strategy is validated in Section 4.

#### 3.4.1 *Softmax-based weighted strategy*

We obtain two output feature maps $f_1, f_2$ from the extraction module for input images $I_1, I_2$, respectively. The output feature map from the extraction module can be used to generate the corresponding weight map that indicates the amount of contribution of each pixel to the fused feature map [21]. First, to get the weight map, we take the Softmax [22] operation to the feature map which can be realized by Equation (3), where $x_i$ is the $i$th channel of the output feature map $x$.

$$S(x_i) = \frac{\exp(x_i)}{\sum_i \exp(x_i)} \quad (3)$$

After we obtained the Softmax output, we compute the matrix nuclear norm ($\|\cdot\|_*$), which is the summation of its singular values. Finally, we obtain the weights for the output feature map by taking the weighted average of the maximum value of the nuclear norm. The formula is given in Equation (4).

$$W_1 = \frac{\phi(\|S(x_i)^1\|_*)}{\phi(\|S(x_i)^1\|_*) + \phi(\|S(x_i)^2\|_*)} \quad (4)$$

$$W_2 = \frac{\phi(\|S(x_i)^2\|_*)}{\phi(\|S(x_i)^1\|_*) + \phi(\|S(x_i)^2\|_*)} \quad (5)$$

Where $S(x_i)^j, j \in [1,2]$ is the weight map after Softmax operation for the feature map $f_j, j \in [1,2]$, $\phi(\cdot)$ can be any arbitrary functions. The final fused feature map is then given by $f = W_1 * f_1 + W_2 * f_2$.

### 3.5 Reconstruction Module

The input of the reconstruction module is the fused feature map from the fusion module. The reconstruction module is used to generate the fused human-visible image from the fused feature map while retaining as many details as possible. Our reconstruction module utilizes three consecutive $3\times3$ convolutional layers with output channels 64, 32, and 1 to reduce the channel from 64 to 1, indicating the output image is gray-scale. An overview of the reconstruction module is shown in Figure 5.
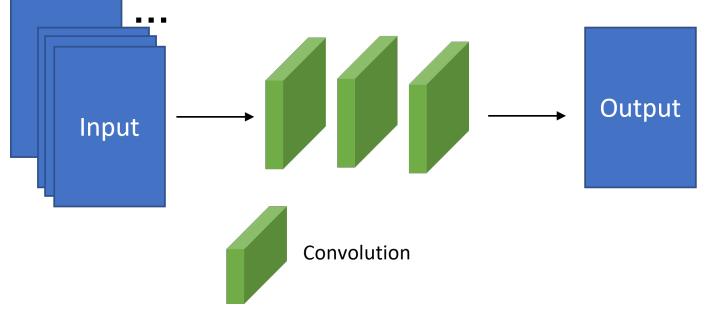


Fig. 5: Image reconstruction module architecture. There are three consecutive $3 \times 3$ convolutional layers, each with the output channel 64, 32, and 1, respectively.

### 3.6 Loss function

We hypothesize that in order to make the fused image as close as two input images, we need to minimize the distance between the image reconstructed from deep semantic features and the original input image. Then, after we apply the fusion strategy and obtain the fused deep features, the reconstructor will produce a more realistic and reasonable fused image. Our goal is to coalesce input images to form an output fused image that retains as many details as the input. In other words, the fused image should have a strong correlation with both input images. We argue that the gap between fused images and input images can be mainly measured by the common regression loss function "mean squared error" (MSE). We also add gradient loss [23], [24] to model the fine details of textures in the reconstructed image and perceptual loss [25] to model the high-level semantic similarity between reconstructed and input images. In detail, our loss function is defined as follows:

$$MSE = \sum_{j=1}^{M} \|I_o - I_j\|_F^2 \quad (6)$$

$$\nabla_I = \sum_{j=1}^{M} \|\nabla I_o - \nabla I_j\|_2^2 \quad (7)$$

$$Percep = \frac{1}{C * W_i * H_i} \sum_{j=1}^{M} \sum_{k=1}^{C} \|f_i^k(I_o) - f_i^k(I_j)\|_2^2 \quad (8)$$

$$\mathcal{L}(\theta) = \frac{1}{W * H}(MSE + \lambda_1 * \nabla_I) + \lambda_2 * Percep \quad (9)$$

Equation (6) is a variant of the MSE loss, where $M$ is the number of input images, $I_o$ is the output image, $\|\cdot\|_F$ is the matrix Frobenius norm. Image gradient loss is given in Equation (7), which is realized by the $\mathcal{L}_2$ norm of the image gradient in $x$ and $y$-direction. Equation (8) is the perceptual loss [25], where $f_i^k(x)$ is the the $k$th channel in $i$th layer (with size $W_i \times H_i$) from the pre-trained VGG16 network [26] with input image $x$, and $C$ is the number of channels. We prefer $i$ to be large, i.e., the deeper layer of the VGG network. Finally, the total loss function is given in Equation (9), $\theta$ is the set of network weights to be optimized, $\lambda_1, \lambda_2$ are weight balancing factors of the gradient and perceptual

loss, respectively. The total loss has been normalized by the total pixels in the image with width $W$ and height $H$.

## 4 EXPERIMENTAL RESULTS

In this section, we will discuss the dataset we used, the experimental setup, evaluation metrics, and the experimental results we obtained.

### 4.1 Data

We adapt the commonly used medical image fusion dataset "The Harvard Whole Brain Atlas" [1] [27] for this work. The dataset we obtained contains three sets of co-registered multimodal medical image pairs: MRI-CT (184 pairs), MRI-PET (269 pairs), and MRI-SPECT (357 pairs). All images are with size $256 \times 256$, MRI and CT are single-channel images with pixel intensities in the range of $[0, 255]$, and PET and SPECT are three-channel images with pixel intensities in the range of $[0, 255]$. In terms of the computational complexity of our model and the time given for this work, we only focus on the MRI-CT fusion task.

### 4.2 Experimental Setup

For the MRI-CT fusion task, we have 184 image pairs. We randomly select 20 image pairs for testing, so that these data are not included in the training process. For the rest of the data, we further split to training and validation data with the ratio of $0.8 : 0.2$, i.e., 80% of the data for training, 20% for validation. All images are normalized in the range of $[0, 1]$ before training. During the training phase, we take out the fusion module, and only train the feature extractor and reconstructor as discussed in Section 3.6. The stochastic gradient descent with Adam update rule [28] is utilized as the optimizer, the learning rate is set to 0.0001 and the model is trained for 100 epochs. The weight balancing factors $\lambda_1$ and $\lambda_2$ for image gradient loss and perceptual loss are all set to 0.2, and we use the third VGG block [26] output to calculate the perceptual loss. To prove the effectiveness of the proposed method, we compare it with the following methods: zero-shot learning for medical image fusion [21], MSRPAN [16], and MSDRA [17]. The parameters of these methods are the default setting suggested by the authors. The experimental environment is a Windows 10 (64 bit) system computer with the Intel Core I7 CPU, NVIDIA GeForce RTX 3050 Ti GPU, and 16 GB of RAM.

### 4.3 Evaluation Metrics

There are many fusion metrics have been proposed in the past few years. However, different metrics reflect the different perspectives of fusion performance. Hence, we select six different, but commonly used fusion metrics in this work: Peak signal-to-noise ratio (PSNR), Structural Similarity (SSIM) [29], Feature SSIM (FSIM) [30], Mutual information (MI) [31], pixel-wise feature MI ($FMI_{pixel}$) [32], and Information Entropy (EN). PSNR is the ratio between the maximum signal power and the signal noise power. It is used to measure the quality of the reconstructed image. The larger the PSNR value is, the better the output image

1. Available at https://www.med.harvard.edu/aanlib/

quality. Since PSNR is a reference-based metric, for input image $I_1, I_2$ and fused image $I_o$, we first compute the $PSNR(I_1, I_o)$ and $PSNR(I_2, I_o)$, and then the average of these two numbers to obtain the final PSNR value. SSIM is a metric to measure the structural similarity of two images from brightness, contrast, and structure. The SSIM is calculated in the same way as the PSNR. FSIM is a full reference fusion metric that measures how close the phase congruency and magnitude of the gradient are between the fused image and source image(s). MI measures the similarity of the image intensity distributions between two or more images, which is calculated in the same way like the PSNR as well. Pixel-wise feature MI is a non-reference image fusion metric that measures the mutual information of the image features between fused image and source image(s). Finally, information entropy measures the amount of information in the fused image. For all metrics, the larger the value, the better the fusion quality.

### 4.4 Selection of fusion strategy

As we discussed in Section 3.4 and Equation (4), our candidate functions for $\phi(\cdot)$ are $max()$, $mean()$, and $sum()$ functions to determine the weights from the matrix nuclear norm. We also investigate the squared version of the three functions above, e.g., $\phi(\cdot) = max()^2$. Next, we provide detailed quantitative results when different fusion strategies are used in Table 1. FER [16] and FL1N [17] are two fusion strategies proposed previously for the same task. The SFNN is what we proposed in this work, and selection methods are described above. Our proposed fusion strategy performed well on four metrics (PSNR, FMI, FSIM, EN). Figure 6 shows the qualitative results of different fusion strategies. For SFNN, we select two strategies that produce the best image quality visually, and also notice that different choices of $\phi$ in our proposed strategy do not affect the metrics very much. Although our SFNN-$max^2$ achieves the best result in terms of the objective metrics, it has a low visual fidelity, hence we select the second optimal SFNN-$max$ as our ultimate strategy in this work. Compared with the FER strategy [16] in Figure 6c, our results have better fidelity, and the inner tissue boundary is more clear. Similarly, our results produce a more brightness edge than FL1N strategy [17] in Figure 6d, which is better when separating between the edge and tissues.

TABLE 1: Comparison between different fusion strategies, bold values represent the best results, and values in blue represent the second-best results.

|  | PSNR | SSIM | MI | FMI-Pixel | FSIM | Entropy |
| --- | --- | --- | --- | --- | --- | --- |
| FER [16] | 13.944 | 0.736 | 4.551 | 0.876 | 0.806 | 8.720 |
| FL1N [17] | 15.979 | 0.739 | 4.569 | 0.878 | 0.813 | 9.782 |
| SFNN ($mean$) | 15.876 | 0.740 | **4.578** | 0.876 | 0.812 | 9.772 |
| SFNN ($mean^2$) | 16.311 | **0.741** | 4.567 | 0.883 | 0.819 | 9.812 |
| SFNN ($max$) | 16.413 | 0.740 | 4.558 | **0.891** | 0.820 | 9.816 |
| SFNN ($max^2$) | **18.258** | 0.740 | 4.451 | **0.891** | **0.829** | **9.887** |
| SFNN ($sum$) | 15.876 | 0.740 | **4.578** | 0.876 | 0.812 | 9.772 |
| SFNN ($sum^2$) | 16.311 | 0.740 | **4.578** | 0.876 | 0.819 | 9.812 |

(a) CT      (b) MRI

(c) FER [16]      (d) FL1N [17]
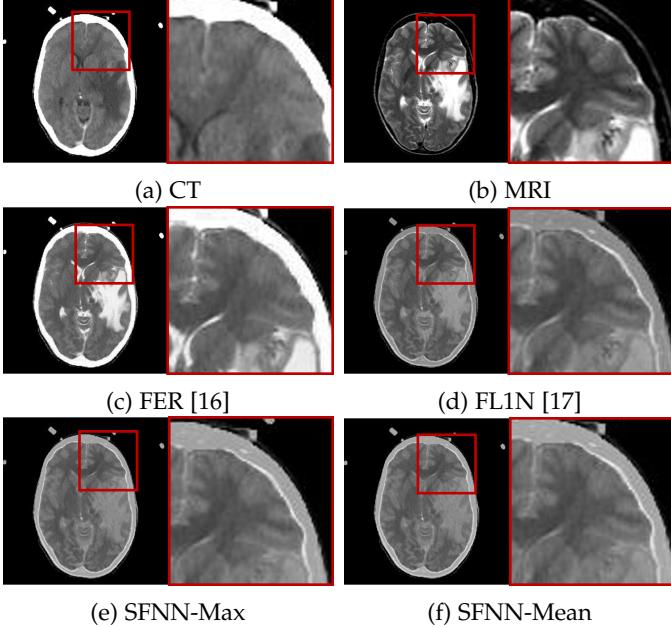
(e) SFNN-Max      (f) SFNN-Mean

Fig. 6: Qualitative comparison of different fusion strategies between CT and MRI. (a) and (b) are source images, (c) is the FER fusion strategy [16], (d) is the FL1N fusion strategy [17], and (e), (f) are the strategies we proposed in this work.

(a) CT      (b) MRI

(c) Zero-learning [21]      (d) MSPRAN [16]
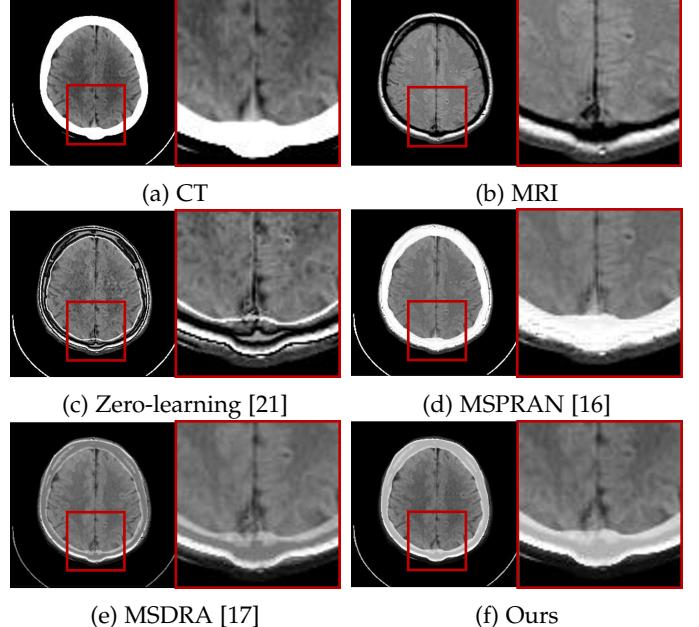
(e) MSDRA [17]      (f) Ours

Fig. 7: Qualitative comparison of different methods between CT and MRI. (a) and (b) are source images, (c) is the zero-learning method [21], (d) is the MSPRAN method [16], (e) is the MSDRA method [17], and (f) is the method we proposed in this work.

## 4.5 Comparison with other methods

Table 2 presents the quantitative results of the MRI-CT fusion task between our proposed method and other baseline methods. Our method performs well on multiple fusion metrics: PSNR, pixel-wise FMI, FSIM, and Entropy. The largest PSNR value indicates our fused image contains less noise and achieves the best image quality. The largest FMI value demonstrates our fused image contains as many source image features as possible. The largest FSIM value suggests our fused image has less information loss at the feature level. Finally, the largest entropy shows our fused image contains more information and details.

TABLE 2: Comparison between different methods, bold numbers represent optimal values

| | PSNR | SSIM | MI | FMI-Pixel | FSIM | Entropy |
|---|---|---|---|---|---|---|
| Zero-shot [21] | 13.525 | 0.681 | 4.633 | 0.836 | 0.738 | 4.279 |
| MSDRA [17] | 15.693 | 0.697 | 4.586 | 0.867 | 0.797 | 8.167 |
| MSRPAN [16] | 14.528 | **0.741** | **4.652** | 0.874 | 0.808 | 8.969 |
| Ours | **16.413** | 0.740 | 4.558 | **0.891** | **0.820** | **9.816** |

From Figure 7, it is also obvious from the enlarged red box that the proposed method preserves the edge and detail information well. Compared to other methods, Figure 7c does not retain the edge information from CT, and details of inner tissues are also distorted; Figure 7d has a favorable visual appearance, but the fidelity is low and lost the tissue boundary information from the MRI source image. Figure 7e has an uncleared boundary so it is hard to differentiate between boundaries and tissues. Finally, our proposed method has a better intensity contrast between edges and tissues, retains important information from both source images, and results in better fidelity.

## 4.6 Ablation study

To test the effectiveness of our proposed loss function, we ablate the image gradient loss and the perceptual loss as described in Equation (7) and (8). Hence, the loss function becomes only the MSE loss in Equation (6). Table 3 shows the quantitative results for the two loss functions on six different fusion metrics.

TABLE 3: Ablation study on the loss function, bold numbers represent optimal values.

| | PSNR | SSIM | MI | FMI-Pixel | FSIM | Entropy |
|---|---|---|---|---|---|---|
| Only MSE | **16.519** | 0.739 | 4.428 | 0.890 | **0.820** | 8.72 |
| All loss | 16.413 | **0.740** | **4.558** | **0.891** | 0.820 | **9.816** |

We observe that using the ablated loss function (MSE loss) only achieves a higher PSNR value by a small margin. The rest of the metrics except FSIM, are all worse than using all three losses (Equation (9)). FSIM values are equal across both loss functions. It is evident that gradient and perceptual loss will force the network to learn the high-level semantic features and focus more details on the features that would potentially contain important information, rather than purely minimizing pixel intensity differences using the traditional MSE loss.

## 5 DISCUSSION AND LIMITATIONS

Through the comparison of different fusion results visually, we can see that the fused image using our proposed framework have a clearer representation of boundaries and details of inner tissues. The image produced by our framework is more natural than other baseline methods, where other images are either a little bit sharp, contain artificial noise, or the

boundary information is not clear enough. The quantitative results validate our findings by performing well on several fusion metrics.

We do not test our proposed framework on the other two fusion tasks (MRI-SPECT and MRI-PET). The SPECT and PET images are 3-channel and MRI images are 1-channel gray-scale. To achieve the fusion task, we could transfer the 3-channel images into YCbCr color space and only focus on the Y channel for structural details and brightness information. The Y channel and the single channel MRI can be used as inputs to our framework. Then, the output can be interpreted as the fused Y-channel. We combine it with Cb and Cr channels and convert it back to RGB to get the final fused 3-channel image. We expect our proposed method also performs well on the 3-channel image fusion task.

Our proposed fusion strategy does not perform well in some cases (see Figure 8) and may introduce fusion noise to the result. In future work, we will focus on reducing the noise while maintaining the overall fusion quality.



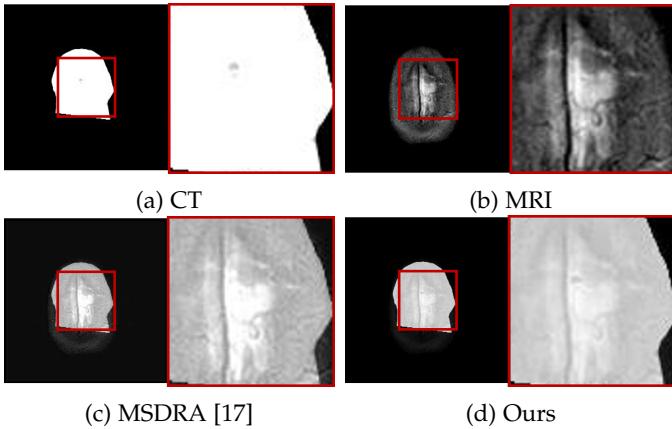(a) CT  (b) MRI

(c) MSDRA [17]  (d) Ours

Fig. 8: Failure case of our proposed method, the features from CT do not handle properly and the details from MRI can barely be distinguished. MSDRA [17] produces a slightly better visual appearance but it suffers from the artificial noise in the image.

## 6 CONCLUSION

In this paper, we propose a novel network architecture based on multi-scale feature extraction and fusion strategy for multimodal medical image fusion. The original multimodal images are passed into the feature extractor to multi-scale deep semantic features. The feature maps obtained from the feature extractor are fused based on the fixed Softmax-based weighted feature fusion strategy we proposed in this work. The fused feature map is used as the input to the reconstruction module to obtain the final fused image. Our fusion strategy is fixed, which means that there is no parameter that needs to be updated in both the training and inference phase. Thus, it can achieve real-time image fusion. Extensive experiments show our proposed method is superior to several baseline methods in both subjective visual appearance and objective fusion metrics.

Our network architecture is based on the MSRPAN [16], which improves on their pyramid structure. Experiments also validate the effectiveness of our method compared to other reference methods. However, our method performs not well on some corner cases and might introduce artificial noise in the fused image as we discussed previously. Despite the limitations, our method could provide a reliable reference for disease diagnosis in the real-life clinical routine.

## REFERENCES

[1] K. Wang, M. Zheng, H. Wei, G. Qi, and Y. Li, "Multi-modality medical image fusion using convolutional neural network and contrast pyramid," *Sensors*, vol. 20, no. 8, p. 2169, 2020.

[2] H. Xu, J. Ma, Z. Le, J. Jiang, and X. Guo, "Fusiondn: A unified densely connected network for image fusion," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 12 484–12 491.

[3] H. Hermessi, O. Mourali, and E. Zagrouba, "Multimodal medical image fusion review: Theoretical background and recent advances," *Signal Processing*, vol. 183, p. 108036, 2021. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S016516842100075X

[4] P. Pérez, M. Gangnet, and A. Blake, "Poisson image editing," in *ACM SIGGRAPH 2003 Papers*, 2003, pp. 313–318.

[5] B. Shreyamsha Kumar, "Image fusion based on pixel significance using cross bilateral filter," *Signal, image and video processing*, vol. 9, no. 5, pp. 1193–1204, 2015.

[6] A. L. Da Cunha, J. Zhou, and M. N. Do, "The nonsubsampled contourlet transform: theory, design, and applications," *IEEE transactions on image processing*, vol. 15, no. 10, pp. 3089–3101, 2006.

[7] Y. Tian, Y. Li, and F. Ye, "Multimodal medical image fusion based on nonsubsampled contourlet transform using improved pcnn," in *2016 IEEE 13th International Conference on Signal Processing (ICSP)*. IEEE, 2016, pp. 799–804.

[8] I. Ozer, Z. Ozer, and O. Findik, "Noise robust sound event classification with convolutional neural network," *Neurocomputing*, vol. 272, pp. 505–512, 2018. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0925231217312675

[9] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, X. Wang, and X. Tang, "Residual attention network for image classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3156–3164.

[10] D. Zhang, J. Shao, G. Hu, and L. Gao, "Sharp and real image super-resolution using generative adversarial network," pp. 217–226, 2017.

[11] Y. Liu, X. Chen, H. Peng, and Z. Wang, "Multi-focus image fusion with a deep convolutional neural network," *Information Fusion*, vol. 36, pp. 191–207, 2017. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S1566253516302081

[12] H. Hermessi, O. Mourali, and E. Zagrouba, "Convolutional neural network-based multimodal image fusion via similarity learning in the shearlet domain," *Neural Comput. Appl.*, vol. 30, no. 7, p. 2029–2045, oct 2018. [Online]. Available: https://doi.org/10.1007/s00521-018-3441-1

[13] H. Li and X.-J. Wu, "Densefuse: A fusion approach to infrared and visible images," *IEEE Transactions on Image Processing*, vol. 28, no. 5, pp. 2614–2623, 2019.

[14] X. Song, X.-J. Wu, and H. Li, "Msdnet for medical image fusion," in *International conference on image and graphics*. Springer, 2019, pp. 278–288.

[15] J. Huang, Z. Le, Y. Ma, F. Fan, H. Zhang, and L. Yang, "Mgmdcgan: Medical image fusion using multi-generator multi-discriminator conditional generative adversarial network," *IEEE Access*, vol. 8, pp. 55 145–55 157, 2020.

[16] J. Fu, W. Li, J. Du, and Y. Huang, "A multiscale residual pyramid attention network for medical image fusion," *Biomedical Signal Processing and Control*, vol. 66, p. 102488, 2021.

[17] W. Li, X. Peng, J. Fu, G. Wang, Y. Huang, and F. Chao, "A multi-scale double-branch residual attention network for anatomical–functional medical image fusion," *Computers in Biology and Medicine*, vol. 141, p. 105005, 2022.

[18] H. Li, P. Xiong, J. An, and L. Wang, "Pyramid attention network for semantic segmentation," *arXiv preprint arXiv:1805.10180*, 2018.

[19] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2818–2826.

[20] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *arXiv preprint arXiv:1511.07122*, 2015.

[21] F. Lahoud and S. Süsstrunk, "Zero-learning fast medical image fusion," in *2019 22th International Conference on Information Fusion (FUSION)*. IEEE, 2019, pp. 1–8.

[22] J. Bridle, "Training stochastic model recognition algorithms as networks can lead to maximum mutual information estimation of parameters," *Advances in neural information processing systems*, vol. 2, 1989.

[23] C. Ma, Y. Rao, Y. Cheng, C. Chen, J. Lu, and J. Zhou, "Structure-preserving super resolution with gradient guidance," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 7769–7778.

[24] C. Fabbri, M. J. Islam, and J. Sattar, "Enhancing underwater imagery using generative adversarial networks," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 7159–7165.

[25] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *European conference on computer vision*. Springer, 2016, pp. 694–711.

[26] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[27] D. Summers, "Harvard whole brain atlas: www. med. harvard. edu/aanlib/home. html," *Journal of Neurology, Neurosurgery & Psychiatry*, vol. 74, no. 3, pp. 288–288, 2003.

[28] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[29] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.

[30] L. Zhang, L. Zhang, X. Mou, and D. Zhang, "Fsim: A feature similarity index for image quality assessment," *IEEE transactions on Image Processing*, vol. 20, no. 8, pp. 2378–2386, 2011.

[31] G. Qu, D. Zhang, and P. Yan, "Information measure for performance of image fusion," *Electronics letters*, vol. 38, no. 7, p. 1, 2002.

[32] M. B. A. Haghighat, A. Aghagolzadeh, and H. Seyedarabi, "A non-reference image fusion metric based on mutual information of image features," *Computers & Electrical Engineering*, vol. 37, no. 5, pp. 744–756, 2011.