

RFN-Nest: An end-to-end residual fusion network for infrared and visible images

Hui Li^a, Xiao-Jun Wu^{a,*}, Josef Kittler^b

^a*Jiangsu Provincial Engineering Laboratory of Pattern Recognition and Computational Intelligence,
School of Artificial Intelligence and Computer Science, Jiangnan University,
214122, Wuxi, China*

^b*The Center for Vision, Speech and Signal Processing, University of Surrey, Guildford, GU2 7XH, UK.*

Abstract

In the image fusion field, the design of deep learning-based fusion methods is far from routine. It is invariably fusion-task specific and requires a careful consideration. The most difficult part of the design is to choose an appropriate strategy to generate the fused image for a specific task in hand. Thus, devising learnable fusion strategy is a very challenging problem in the community of image fusion. To address this problem, a novel end-to-end fusion network architecture (RFN-Nest) is developed for infrared and visible image fusion. We propose a residual fusion network (RFN) which is based on a residual architecture to replace the traditional fusion approach. A novel detail-preserving loss function, and a feature enhancing loss function are proposed to train RFN. The fusion model learning is accomplished by a novel two-stage training strategy. In the first stage, we train an auto-encoder based on an innovative nest connection (Nest) concept. Next, the RFN is trained using the proposed loss functions. The experimental results on public domain data sets show that, compared with the existing methods, our end-to-end fusion network delivers a better performance than the state-of-the-art methods in both subjective and objective evaluation. The code of our fusion method is available at <https://github.com/hli1221/imagefusion-rfn-nest>.

Keywords:

image fusion, end-to-end network, nest connection, residual network, infrared image, visible image

1. Introduction

Due to the physical limitations of imaging sensors, it is very difficult to capture an image of a scene that is of uniformly good quality. Image fusion plays an important role in this context. Its aim is to reconstruct a perfect image of the scene from multiple samples that provide complementary information about the visual content. It has many applications, such as object tracking [1] [2] [3], self-driving and video surveillance [4]. The fusion task requires algorithms to generate a single image which amalgamates the complementary information conveyed by different source images [5][6][7].

Image fusion involves three key processes: feature extraction, fusion strategy and reconstruction. Most of the existing fusion research focuses on one or more of these elements to improve the fusion performance. The existing fusion methods can be classified into two categories: traditional algorithms and deep learning-based methods. In the traditional algorithm category, multi-scale transform methods [8][9][10][11] are widely applied to extract multi-scale features from the source images. The feature channels are combined by an appropriate fusion strategy. Finally, the fused image is reconstructed by an inverse multiscale transform. Obviously, the fusion performance of these algorithms is highly dependent on the feature extraction method used.

Following this direction, sparse representation (SR) [12] and low-rank representation (LRR) [13][14] have been applied to extract salient features from the source images. In SR and LRR based fusion methods [15][16][17][18], the sliding window technique is used to decompose source images into image patches. A matrix is constructed using these image patches, in which each column is a reshaped image patch. This matrix is fed into SR (or LRR) to calculate SR (or LRR) coefficients which are considered as image features. By virtue of this operation, the image fusion problem is transformed to the one of coefficient fusion. The fused coefficients are generated by an appropriate fusion strategy and used to reconstruct the fused image in the SR (or LRR) framework. Beside the above, other SR based approaches and other signal processing methods [19][20] have been suggested in the literature.

Although the traditional fusion methods have achieved good fusion performance, they have drawbacks: (1) The fusion performance highly depends on handcrafted features [15][18][21], as it is difficult to find a universal feature extraction method for different fusion tasks; (2) Different fusion strategies may be required to work with different features; (3) For SR and LRR based methods, the dictionary learning is very time-consuming; (4) Complex source images pose a challenge for SR (or LRR) based fusion methods.

To overcome these drawbacks, deep learning based fusion methods have been developed, which can be grouped into three categories according to the three key elements of the fusion pro-

*Corresponding author email: wu_xiaojun@jiangnan.edu.cn

cess: deep feature extraction, fusion strategy and end-to-end training. In the feature extraction direction, deep learning methods are utilized to extract deep representation of the information conveyed by the source images [22][23][24][25][26]. Different fusion strategies have been suggested to reconstruct the fused image. In other fusion methods [27][28], deep learning is also used to design the fusion strategy. In [27][28], convolutional sparse representation and convolutional neural network are utilized to generate a decision map for the source images. Using the learned decision map, the fused images are obtained by appropriate post-processing. Although these fusion methods achieve good fusion performance, the fusion strategy and the post-processing are tricky to design. To avoid the limitations of handcrafted solutions, some end-to-end fusion frameworks were presented (FusionGAN [29], FusionGANv2 [30], DDcGAN [31]). These frameworks are based on adversarial learning which avoids the shortcoming of handcrafted features and fusion strategies. However, even the state of the art methods, FusionGANv2 [30] and DDcGAN [31], face challenges to preserve image detail adequately. To preserve more detail background information from visible images, a nest connection based autoencoder fusion network (NestFuse [32]) was proposed. Although NestFuse obtains good performance in detail information preservation, the fusion strategy is still not learnable.

To address these problems, in this paper, we propose a novel end-to-end fusion framework (RFN-Nest). Our network contains three parts: an encoder network, residual fusion network (RFN) which is designed to extract fused multi-scale deep features, and a decoder network based on nest connection [33]. Although the encoder and decoder architecture of the proposed network is similar to the NestFuse [32], the fusion strategy, the training strategy and the loss function are totally different.

Firstly, instead of fusing handcrafted features in NestFuse [32], several simple yet efficient learnable fusion networks (RFN) have been designed and inserted into the autoencoder architecture. With the RFN, the autoencoder-based fusion network is upgraded to an end-to-end fusion network. Secondly, as RFN is a learnable structure, it is important that the encoder and decoder exhibit powerful feature extraction and feature reconstruction abilities, respectively. Thus, we develop a two-stage training strategy to train our fusion network (encoder, decoder and RFN networks). Thirdly, to train the proposed RFN networks, we design a new loss function (L_{RFN}) to preserve the detail information from visible image and maintain the salient features from infrared image, simultaneously.

The main contributions of RFN-Nest can be summarized as follows,

(1) A novel residual fusion network(RFN) is proposed to supersede handcrafted fusion strategies. Although many methods [22][24][34][25] now use deep features to achieve good performance, the heuristic approach to selecting a suitable fusion strategy is their weakness. The proposed RFN is a learnable fusion network that overcomes this weakness.

(2) A two-stage training strategy is developed to design our network. The feature extraction and feature reconstruction abilities are the key for the encoder and decoder networks. Us-

ing only one stage training strategy to simultaneously train the whole network (encoder, decoder and RFN networks) is insufficient. Inspired by [25], firstly, the encoder and the decoder network are trained as an auto-encoder. With the fixed encoder and decoder, the RFN networks are trained using an appropriate loss function.

(3) A loss function capable of preserving the image detail, together with a feature enhancing loss function are designed to train our RFN networks. We show that with these loss functions, more detail information and image salient features are preserved in the fused image.

(4) We show that, compared with the state-of-the-art fusion methods, the proposed RFN-Nest framework exhibits better fusion performance on public datasets in both subjective visual assessment and objective evaluation.

The rest of our paper is structured as follows. In Section 2, we briefly review the related work on deep learning-based fusion. The proposed fusion framework is described in detail in Section 3. The experimental results are presented in Section 4 and Section 5. Finally, we draw the paper to conclusion in Section 6.

2. Related Works

Recently, many deep learning methods have been developed for image fusion. Most of them are based on convolutional neural networks (CNN). These methods can be classified into the non end-to-end learning and end-to-end learning categories. In this section, we briefly overview the most representative deep learning based methods from these two categories.

2.1. Non End-to-end Methods

In the early days, deep learning neural networks were used to extract deep features as a bank of “decision” maps [22][23][24]. In [22], Li et al. proposed a fusion framework based on a pre-trained network (VGG-19 [35]). Firstly, the source images are decomposed into salient parts (texture and edges) and base parts (contour and luminance). Then, VGG-19 is used to extract multi-level deep features from the salient parts. At each level, the decision maps are computed from the deep features and a candidate fused salient part is generated. The fused image is reconstructed by combining the fused base parts and the fused salient parts using an appropriate fusion strategy. In [23], the pre-trained ResNet-50 [36] is utilized to extract deep features from the source images directly. A decision map is obtained by zero-phase component analysis(ZCA) and l_1 -norm. The PCANet-based fusion method [24] also follows this framework to generate the fused image, in which PCANet, instead of VGG-19 or ResNet-50, is used to extract the features.

In addition to pure feature extraction, in [27][28], the two key processes (feature extraction and fusion strategy) are implemented by a single network. In [28], a decision map is generated by a CNN trained on image patches of multiple blurred versions of the input image. In [27], the convolutional sparse representation instead of CNN is utilized to extract features and to generate a decision map. From the generated decision map, the fused image can easily be reconstructed.

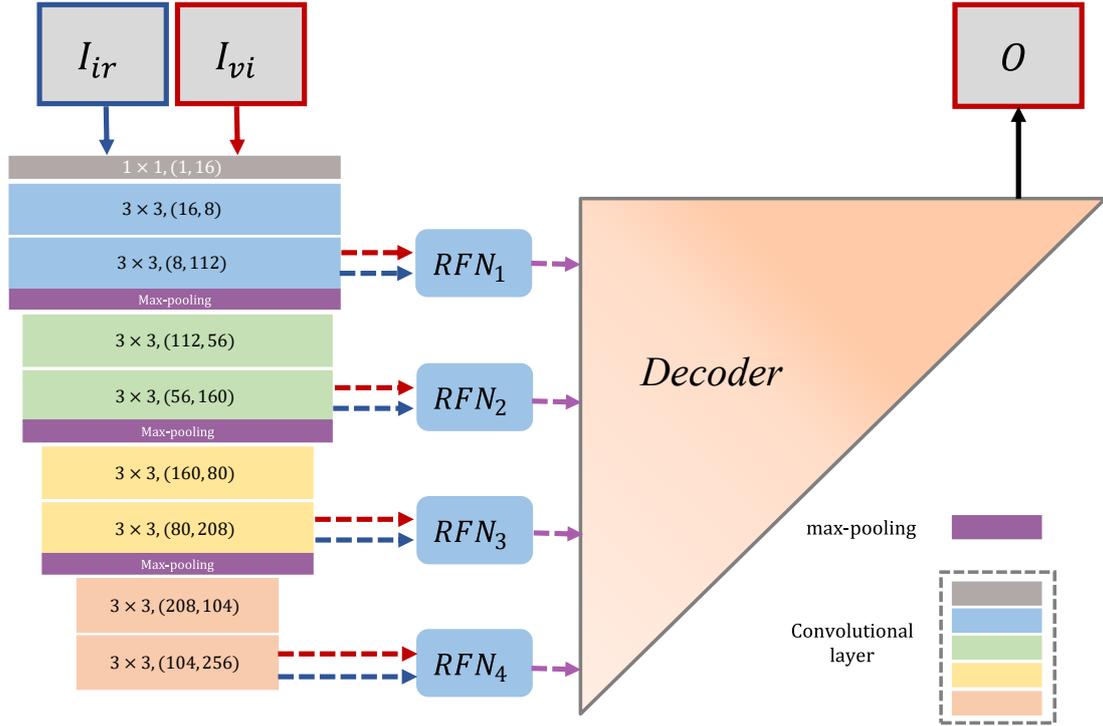


Figure 1: The framework of RFN-Nest. ‘ RFN_{1-4} ’ denote the residual fusion network. The nest connection-based decoder network(‘Decoder’) will be introduced later. ‘ $3 \times 3, (16, 8)$ ’ means the kernel size is 3×3 , input channel is 16 and output channel is 8 in a convolutional layer.

Besides the above methods, a deep auto-encoder network based fusion framework was proposed in [25]. Inspired by DeepFuse [34], the authors proposed a novel network architecture which contains an encoder, a fusion layer and a decoder. The dense block [37] based encoder network was adopted as it extracts more complementary deep features from the source images. In their framework, the fusion strategy becomes very important.

Inspired by DenseFuse[25] and the architecture in [33], Li et al. proposed NestFuse[32] to preserve more detail background information from visible images, while enhancing the salient features in infrared images. Additionally, a novel spatial/channel attention models are designed to fuse the multi-scale deep features. Although these frameworks achieve good fusion performance, it is very difficult to find an effective hand-crafted fusion strategy for image fusion.

2.2. End-to-end Methods

To eliminate the arbitrariness of handcrafted features and fusion strategies, several end-to-end fusion frameworks have been suggested [29] [30] [38] [39] [31] [40].

In [29], a GAN-based fusion framework (FusionGAN) was introduced to the infrared and visible image fusion field. The generator network is the engine which computes the fused image, while the discriminator network constrains the fused image to contain the detail information from the visible image. The loss function has two terms: content loss and discriminator loss. Due to the content loss, the fused image tends to become

similar to the infrared image, failing to preserve the image detail, in spite of the discriminator network.

To preserve more detail information from the visible images, the authors of [30] proposed a new version of FusionGAN which was named FusionGANv2. In this new version, the authors deepen the generator and discriminator networks, endowing them with more powerful feature representation ability. In addition, two new loss functions, namely detail loss and target edge-enhancement loss, were presented to preserve the detail information. With these improvements, the fused images reconstructed more scene details, with clearly highlighted edge-sharpened targets.

A general end-to-end image fusion network (IFCNN) [38] was also proposed, which is a simple yet effective fusion method. In IFCNN, two convolutional layers are utilized to extract deep features from the source images. Element-wise fusion rules (elementwise-maximum, elementwise-sum, elementwise-mean) are used to fuse the deep features. The fused image is generated from the fused deep features by two convolutional layers. Although IFCNN achieves a satisfactory fusion performance in multiple image fusion tasks, its architecture is too simplistic to extract powerful deep features, and the fusion strategies designed using a traditional way are not optimal.

3. The Proposed Fusion Framework

The proposed fusion network is introduced in this section. Firstly, the architecture of our network is presented in Section

3.1. The advocated two-stage training strategy is described in Section 3.2.

3.1. The Architecture of the Fusion Network

The RFN-Nest is an end-to-end fusion network, the architecture of which is shown in Fig.1. RFN-Nest contains three parts: encoder (left part), residual fusion network (RFN_{1-4}) and decoder (right part). For a convolutional layer, “ $k \times k, (in, out)$ ” means the kernel size is $k \times k$, input channel is *in* and output channel is *out*.

With the max pooling operation in the encoder network, multi-scale deep features can be extracted from the source images. The RFN is utilized to fuse multi-modal deep features extracted at each scale. While the shallow layer features preserve more detail information, the deeper layer features convey semantic information, which is important for reconstructing the salient features. Finally, the fused image is reconstructed by the nest connection-based decoder network, which fully exploits the multi-scale structure of the features.

As shown in Fig.1, I_{ir} and I_{vi} indicate the source images (infrared image and visible image). O denotes the output of RFN-Nest, that is the fused image. “ RFN_m ” means one residual fusion network for deep features at scale m . The architecture of the encoder in our framework is constituted by four RFN networks, $m \in \{1, 2, 3, 4\}$. These RFN networks share the same architecture but with different weights.

We now introduce the RFN and the decoder in detail.

3.1.1. Residual Fusion Network (RFN)

The RFN is based on the concept of residual block [36] which has been adapted to the task of image fusion. The RFN architecture is shown in Fig.2.

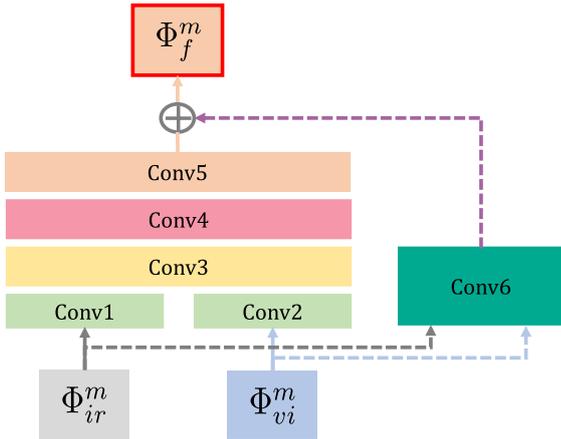


Figure 2: The architecture of RFN_m .

In Fig.2, Φ_{ir}^m and Φ_{vi}^m indicate the m -th scale deep features extracted by the encoder network, with $m \in \{1, 2, 3, 4\}$ indicating the index of the RFN network. “Conv1-6” denote six convolutional layers in RFN. In this residual architecture, the outputs of “Conv1” and “Conv2” are concatenated as the input of “Conv3”. “Conv6” is the first fusion layer to generate initial

fused features. With this architecture, RFN can easily be optimized by our training strategy. The convolutional operations produce the fused deep features Φ_f^m , which are fed into the decoder network.

Thanks to the multi-scale deep features and the proposed learning process, both the image detail and salient structures are preserved by the shallow RFN networks and deep RFN networks, respectively.

3.1.2. Decoder Network

The decoder network based on the nest connection architecture is shown in Fig.3. Compared to UNet++ [33], regarding the image fusion task, we simplify the network architecture to make it light yet effective to reconstruct fused images, this architecture was also utilized in [32].

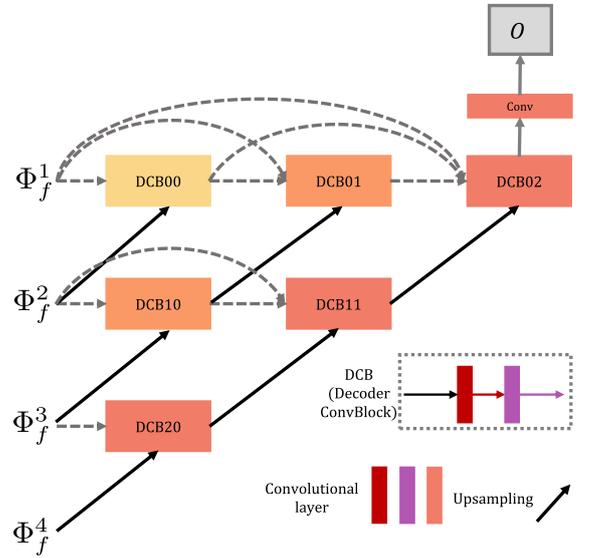


Figure 3: The architecture of the decoder.

Φ_f^m ($m \in \{1, 2, 3, 4\}$) denote the fused multi-scale features obtained by the RFN networks. “DCB” indicates a decoder convolutional block, which has two convolutional layers. In each row, these blocks are connected by short connections which are similar to the dense block architecture [37]. The cross-layer links connect multi-scale deep features in the decoder network.

The output of the network is the fused image reconstructed from the fused multi-scale features.

3.2. Two-stage Training Strategy

Note that the ability of the encoder in our network to perform feature extraction and that of the decoder to conduct feature reconstruction are absolutely crucial for successful operation. Accordingly, we develop a two-stage training strategy to make sure that each part in our network can achieve the expected performance.

Firstly, the encoder and the decoder are trained as an auto-encoder network to reconstruct the input image. After learning the encoder and decoder networks, in the second training stage,

several RFN networks are trained to fuse the multi-scale deep features.

In this section, a novel two-stage training strategy is introduced in detail.

3.2.1. Training of the Auto-encoder Network

Inspired by DenseFuse [25], in the first stage, the encoder network is trained to extract multi-scale deep features. The decoder network is trained to reconstruct the input image with multi-scale deep features. The auto-encoder network training framework is shown in Fig.4.

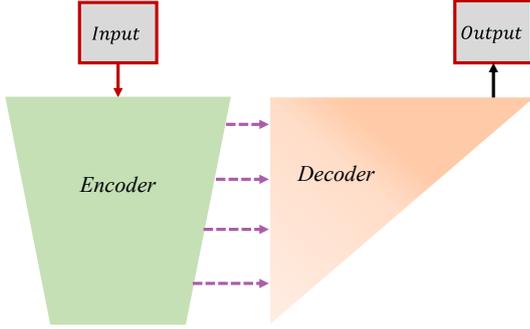


Figure 4: The training of the auto-encoder network.

In Fig.4, *Input* and *Output* denote the input image and the output image (both indicate one image), respectively. In contrast to [34][25][29][30], our feature extraction part includes a down sampling operation (max pooling), which extracts deep features at four scales. These multi-scale deep features are fed into the decoder network to reconstruct the input image. With short cross-layer connections, the multi-scale deep features are fully used to reconstruct the input image.

The auto-encoder network is trained using the loss function, L_{auto} defined as follows,

$$L_{auto} = L_{pixel} + \lambda L_{ssim} \quad (1)$$

where L_{pixel} and L_{ssim} denote the pixel loss and the structure similarity (SSIM) loss between the input image (*Input*) and the output image (*Output*). λ is the trade-off parameter between L_{pixel} and L_{ssim} .

The pixel loss (L_{pixel}) is calculated by Eq.2,

$$L_{pixel} = \|Output - Input\|_F^2 \quad (2)$$

where the $\|\cdot\|_F$ is the Frobenius norm. L_{pixel} constrains the reconstructed image to be like the input image at the pixel level.

The SSIM loss (L_{ssim}) is defined as,

$$L_{ssim} = 1 - SSIM(Output, Input) \quad (3)$$

where $SSIM(\cdot)^1$ is the structural similarity measure [41] which quantifies the structural similarity of the two images. The structural similarity between *Input* and *Output* is constrained by L_{ssim} .

¹The definition of $SSIM(\cdot)$ is introduced in our supplementary material (Section 1).

3.2.2. Training of the RFN

The RFN is proposed to implement a fully learnable fusion strategy. In the second stage, with the encoder and decoder fixed, the RFN is trained with an appropriate loss function. The training process is shown in Fig.5.

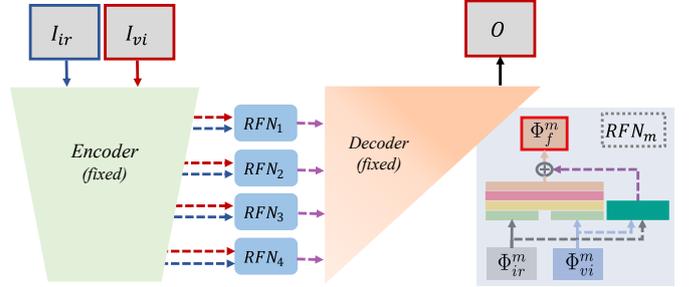


Figure 5: The training of RFN. In the RFN, inputs Φ_{ir}^m and Φ_{vi}^m denote infrared and visible deep features, respectively. Φ_f^m represents the fused deep features obtained by RFN_m . In our framework, $m \in \{1, 2, 3, 4\}$ indicates the scale of deep features.

The fixed encoder network is utilized to extract multi-scale deep features (Φ_{ir}^m and Φ_{vi}^m) from the source images. For each scale, an RFN is used to fuse these deep features. Then, the fused multi-scale features (Φ_f^m) are fed into the fixed decoder network.

To train our RFN, we propose a novel loss function L_{RFN} , which is defined as,

$$L_{RFN} = \alpha L_{detail} + L_{feature} \quad (4)$$

where L_{detail} and $L_{feature}$ indicate the background detail preservation loss function and the target feature enhancement loss function, respectively. α is a trade-off parameter.

In the case of infrared and visible image fusion, most of the background detail information comes from the visible image. L_{detail} aims to preserve the detail information and structural features from visible image, which is defined as

$$L_{detail} = 1 - SSIM(O, I_{vi}) \quad (5)$$

As the infrared image contains more salient target features than the visible image, the loss function $L_{feature}$ is designed to constrain the fused deep features so as to preserve the salient structures. The $L_{feature}$ is defined as,

$$L_{feature} = \sum_{m=1}^M w_1(m) \|\Phi_f^m - (w_{vi}\Phi_{vi}^m + w_{ir}\Phi_{ir}^m)\|_F^2 \quad (6)$$

In Eq.6, M is the number of the multi-scale deep features, which is set to 4. Owing to the magnitude difference between the scales, w_1 is a trade-off parameter vector for balancing the loss magnitudes. It assumes four values $\{1, 10, 100, 1000\}$. w_{vi} and w_{ir} control the relative influence of the visible and infrared features in the fused feature map Φ_f^m .

As the visible information is constrained by L_{detail} and the aim of $L_{feature}$ is to preserve salient features from the infrared image, in Eq.6, w_{ir} is usually greater than w_{vi} .



Figure 6: Five pairs of source images. The top row contains visible images, and the second row contains infrared images.

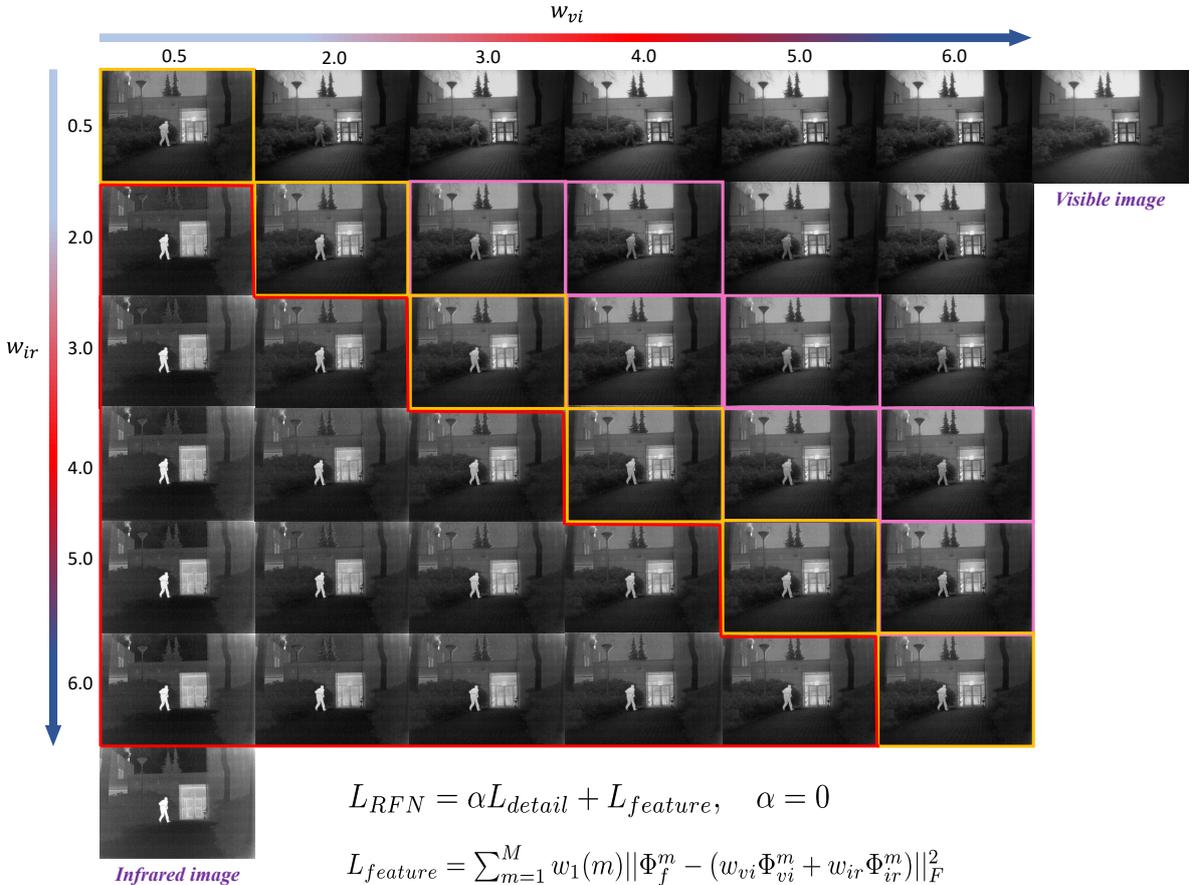


Figure 7: The fusion results obtained with different values of w_{vi} and w_{ir} , when $\alpha = 0$. In $L_{feature}$, $w_1 \in \{1, 10, 100, 1000\}$.

4. Experimental Validation

In this section, we conduct an experimental validation of the proposed fusion method. After detailing the experimental settings in the training phase and the test phase, we present several ablation studies to investigate the effect of different elements of the proposed fusion network. Finally, we compare our fusion framework with other existing algorithms qualitatively. For this purpose, we use several performance metrics to evaluate the fusion performance objectively.

Our network is implemented on the NVIDIA TITAN Xp

GPU using PyTorch as a programming environment.

4.1. Experimental Settings in the Training Phase

In this section, we introduce the training datasets used in our two-stage training strategy.

In the first stage, we use the dataset MS-COCO [42] to train our auto-encoder network. 80000 images are chosen to constitute the training set. These images are converted to gray scale and reshaped to 256×256 . In Eq. 1, the parameter λ is set to 100 to balance the magnitude difference between L_{pixel} and

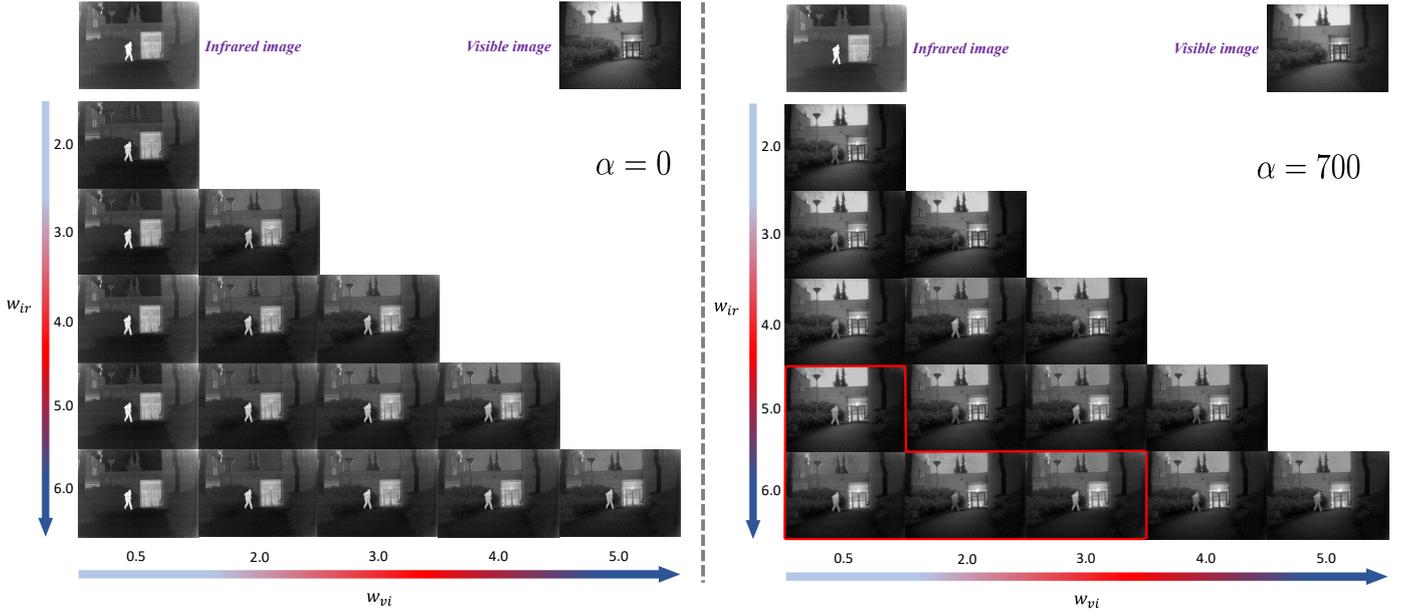


Figure 8: The fusion results obtained without ($\alpha = 0$) or with ($\alpha = 700$) L_{detail} . In $L_{feature}$, $w_1 \in \{1, 10, 100, 1000\}$.

L_{ssim} . The batch size and epoch are set to 4 and 2, respectively. The learning rate is set to 1×10^{-4} .

For the second training stage, we choose the KAIST [43] dataset to train our RFN networks. It contains almost 90000 pairs of images. In this dataset, 80000 pairs of infrared and visible images are chosen for training. These images are also converted to gray scale and resized to 256×256 . The batch size and epoch are set to 4 and 2, respectively. The learning rate was also set to 1×10^{-4} , as in the first stage.

4.2. Experimental Settings in the Test Phase

Our test images come from two datasets which were collected from TNO [44] and VOT2020-RGBT [45]. These images are available at [46]. Some samples of these images are shown in Fig.6. The first dataset contains 21 pairs of infrared and visible images collected from TNO. The second dataset contains 40 pairs of infrared and visible images, which were collected from TNO and VOT2020-RGBT.

We use six quality metrics² to evaluate our fusion algorithm objectively. These include: entropy(En) [47]; standard deviation (SD) [48]; mutual information(MI) [49]; modified fusion artifacts measure(N_{abf}) [50], which evaluates the noise information in fused images; the sum of the correlations of differences(SCD) [51]; and the multi-scale structural similarity(MS -SSIM) [52]. The fusion performance improves with the increasing numerical index of all these seven metrics.

4.3. Ablation Study for L_{detail} and $L_{feature}$

In this section, we discuss the effect of L_{detail} and tune the parameters in $L_{feature}$. Then, we investigate the impact of the

²The definitions of these metrics are introduced in our supplementary material.

relative weights of the visible and infrared features on the fusion performance.

Once the auto-encoder network is trained in the first stage, the parameters of the encoder and decoder are fixed and we use L_{RFN} to train four RFN networks. As discussed in Section 3.2.2, due to the magnitude difference between L_{detail} and $L_{feature}$, the value of the parameter α should be large. Furthermore, the role of L_{detail} is to preserve the detail information from the visible image. Based on the above considerations, in this experiment, α is set to 0 and 700 to analyze its influence on our network.

In $L_{feature}$, w_1 is a trade-off vector to balance the $L_{feature}$ values between the scales. To preserve the salient features from the infrared image, w_{vi} and w_{ir} should be set appropriately. In view of the role of L_{detail} , w_{vi} should be relatively small to reduce any redundancy in reconstructing the image detail information. In contrast, w_{ir} should be large to preserve the complementary salient features in the infrared image. However, if w_{vi} is set to 0, which constrains the fused features to mirror the infrared features, the network fails to converge due to the conflicting constraints of L_{detail} and $L_{feature}$. So, in our experiment, w_{vi} is set to a non-zero value.

As different combination of w_{vi} and w_{ir} can lead to different fusion results, we analyze the influence of these two parameters for different values from the range of [0.5, 6.0].

Firstly, when $\alpha = 0$, which means only $L_{feature}$ is utilized to train RFN networks, some of the fusion results with different w_{vi} and w_{ir} are shown in Fig.7.

In Fig.7, when w_{ir} is small, the fused images are similar to the visible image and the salient features in the infrared images are suppressed (as shown in first two rows). On the contrary, when w_{ir} is large (greater than 3.0), the salient features in the infrared image are retained. In contrast, the detail information in the visible image is not preserved.

To capture both types of information, for $\alpha = 0$, we choose

Table 1: The average values of the objective metrics obtained with different parameters (α , w_{vi} , w_{ir}) on 21 pairs of infrared and visible images.

α	w_{ir}	w_{vi}	En [47]	SD [48]	MI [49]	N_{abf} [50]	SCD [51]	MS -SSIM[52]
0	0.5	0.5	6.71845	67.66313	13.43690	0.09354	1.83520	0.92903
	2.0	2.0	6.71557	67.63524	13.43114	0.09252	1.83495	0.92887
		3.0	6.80410	74.73724	13.60821	0.09240	1.82712	0.92294
		4.0	6.83492	79.75125	13.66983	0.09419	1.78649	0.90543
	3.0	3.0	6.72263	67.83451	13.44526	0.09339	1.83713	0.92988
		4.0	6.78738	72.45840	13.57476	0.09230	1.83518	0.92715
		5.0	6.81292	76.36078	13.62583	0.09324	1.81501	0.91696
	4.0	4.0	6.72150	67.53190	13.44299	0.09355	1.83367	0.92842
		5.0	6.77188	70.98434	13.54376	0.09208	1.83538	0.92775
		6.0	6.80239	74.64694	13.60478	0.09218	1.82594	0.92263
5.0	5.0	6.71684	67.48675	13.43368	0.09218	1.83366	0.92847	
	6.0	6.76875	70.35820	13.53750	0.08944	1.83707	0.92870	
6.0	6.0	6.72585	67.82480	13.45170	0.09209	1.83665	0.92949	
700	5.0	0.5	6.95916	91.41847	13.9183	0.14375	1.58717	0.84109
	6.0	0.5	6.79112	68.28532	13.58224	0.07838	1.78391	0.88602
		3.0	6.84134	71.90131	13.68269	0.07288	1.83676	0.91456

a middle value (yellow and pink boxes in Fig.7) to perform the objective evaluation. The evaluation metrics for different w_{vi} and w_{ir} are presented in Table 1. The best values are indicated in **bold**.

When $\alpha = 700$, the detail information is preserved by L_{detail} . The aim of $L_{feature}$ is to promote the salient features conveyed by the source images. Accordingly, the values of w_{vi} must be smaller than w_{ir} . We choose different combinations of w_{vi} and w_{ir} (the red boxes in Fig.7) to find the best values of w_{vi} and w_{ir} . The fusion results obtained with ($\alpha = 700$) or without ($\alpha = 0$) L_{detail} in the same combinations of w_{vi} and w_{ir} are shown in Fig.8.

In Fig.8 (right part), the fusion results in the red boxes contain more detail information from the source images, yet the infrared features are still maintained. Compared with the left part, the fusion results on the right (red boxes) evidently preserve more detail information. When $\alpha = 700$, the objective values for different parameters (Fig.8 (right part), red boxes) are also presented in Table 1.

From Fig.7 and Table 1, the different values of α , w_{vi} and w_{ir} have a significant influence on the results. If the detail preservation loss function (L_{detail}) is not used ($\alpha = 0$) in the training phase, the proposed fusion network fails to obtain acceptable fusion results. Although the fusion performance appears to be comparable in subjective evaluation (Fig.7, yellow and pink boxes), the subjective and objective assessments indicates a notable degradation compared with the optimal parameter combination ($\alpha = 700$, $w_{ir} = 6.0$ and $w_{vi} = 3.0$).

When the detail preservation loss function (L_{detail}) is switched on $\alpha = 700$, our RFN-Nest fusion network scores the comparable metrics values of six metrics with $w_{ir} = 6.0$ and $w_{vi} = 3.0$. Based on this analysis, we set $w_{ir} = 6.0$ and $w_{vi} = 3.0$ in our next experiments.

In next section, we will analyze the impact of parameter α in our loss function.

4.4. Ablation Study for α in L_{RFN}

As discussed in Section 4.3, when the detail preserving loss function L_{detail} is discarded ($\alpha = 0$), both the subjective and objectively measured fusion performance will be poor. It is evident from Fig.8 and Table 1 that our fusion network can achieve better fusion performance when α is not 0. Thus, choosing an optimal value of α becomes an important issue.

In our study, the parameters of w_{vi} and w_{ir} are set to 3.0 and 6.0, respectively. w_1 is set to {1, 10, 100, 1000} to balance the discrepancy in the orders of magnitude of different scales. To find the optimal α , we set it to {10, 100, 200, 500, 700, 1000} and compute the results.

Some examples of the fusion results are shown in Fig.9. With the increase α (1000), the salient features (man in the yellow box) are not clear, even suppressed, which makes the fused image similar to the visible image. When α is set to 500 and 700, the fusion results contain more detail information and the salient features are also maintained.

Based on these observations, we objectively evaluate our fusion method with α set to 10, 100, 200, 500, 700, 1000. The metrics values of the fusion results with different α are shown in Table 2. The best values are indicated in **bold**.

As shown in Fig.9 and Table 2, when α is 700, the proposed fusion network achieves better fusion performance in both subjective and objective evaluation. In Table 2, the proposed network scores best in four out of seven metrics with $\alpha = 700$. Thus, in our next experiments, the parameter α is set to 700.

$$L_{RFN} = \alpha L_{detail} + L_{feature}$$

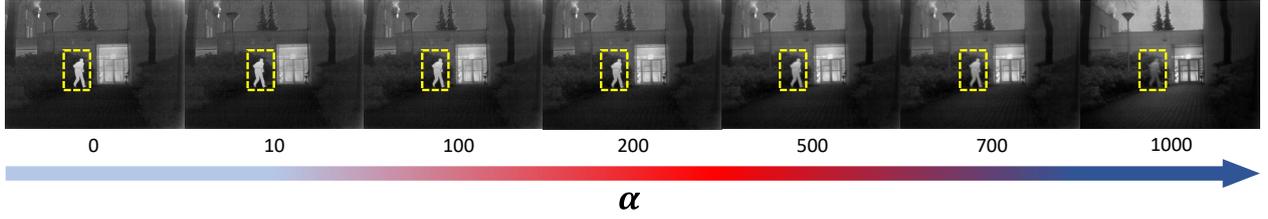


Figure 9: The fusion results obtained with different α .

Table 2: The average metrics values of the proposed fusion network with different α on 21 pairs of infrared and visible images.

α	En [47]	SD [48]	MI [49]	N_{abf} [50]	SCD [51]	$MS-SSIM$ [52]
10	6.66878	62.83593	13.33757	0.08012	1.77192	0.87593
100	6.70939	63.77416	13.41878	0.07129	1.79329	0.90096
200	6.75446	66.01632	13.50893	0.06680	1.80880	0.91177
500	6.82103	70.34117	13.64206	0.06768	1.83252	0.91453
700	6.84134	71.90131	13.68269	0.07288	1.83676	0.91456

4.5. Ablation Study for Training Strategy

The proposed two-stage training strategy is a critical operation in our training phase. In this section, we discuss why this strategy is effective, and show its relative merits compared to the one-stage strategy.

One-stage training strategy means the encoder, RFN and decoder are trained, simultaneously. The training framework is shown in Fig.10, where both the encoder and the decoder are free to adapt their weights. The loss function and the parameter settings are the same as L_{RFN} , which means $\alpha = 700$, $w_1 = \{1, 10, 100, 1000\}$, $w_{vi} = 3.0$ and $w_{ir} = 6.0$. The fusion results obtained by these two training strategies are shown in Fig.11.

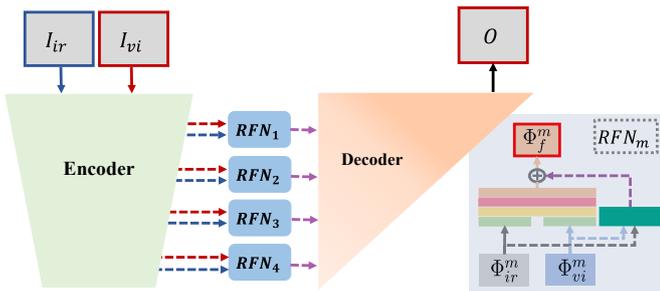


Figure 10: The training framework for the one-stage training strategy.

In Fig.11 (c), the visible spectrum detail information is enhanced with the one-stage training strategy. However, the salient objects in infrared image are lost. The premise of image fusion is not realised. In contrast, the two-stage training strategy (Fig.11, d) enables the fused image to preserve the salient infrared objects and contain more detail information from visible images.

The reason is that the encoder and the decoder may not have the desirable feature extraction and reconstruction ability when

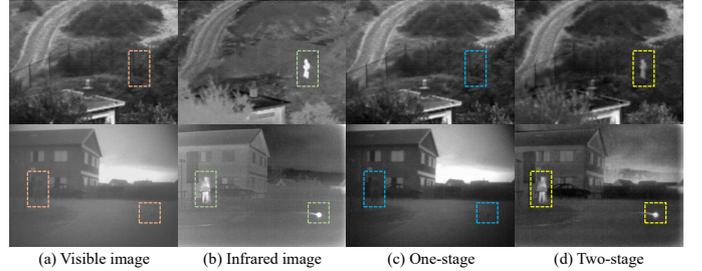


Figure 11: The fusion results obtained by one-stage and two-stage training strategy. (a) Visible images; (b) Infrared images; (c) Fused images obtained by one-stage strategy; (d) Fused images obtained by two-stage strategy.

designed using the one-stage training strategy. More importantly, as the RFN is the key in our fusion network, it should be trained carefully to obtain good fusion performance.

In conclusion, we use the two-stage training strategy to train our fusion network. In the first training stage, the encoder is trained to extract powerful multi-scale deep features, to be used by the decoder for image reconstruction. In the second stage, with the fixed encoder and decoder, the RFN networks are trained to fuse the multi-scale deep features, to enhance the detail information from the visible spectrum image and to preserve salient features from the infrared source image.

4.6. Ablation Study for Nest Connection in Decoder

In this section, we discuss the influence of the nest connection in the decoder. Fig.12 shows the decoder network structure without nest connection (remove the short connections between “DCB”). We train this new decoder architecture with the same training strategy and the same loss functions as discussed in Section 3.2.

The values of seven quality metrics are shown in Table 3. The “No-nest” denotes the decoder without nest-connection archi-

Table 3: The quality metrics values of two ablation studies. “No-nest” indicates the decoder without nest-connection architecture. “Encoder & Decoder” denotes that the encoder and the decoder are fixed, the fusion strategy is “add”, “max”, “l₁-norm”, “l_{*}-norm” or “SCA”. “RFN-Nest” means the proposed fusion network.

		<i>En</i> [47]	<i>SD</i> [48]	<i>MI</i> [49]	<i>N_{abf}</i> [50]	<i>SCD</i> [51]	<i>MS-SSIM</i> [52]
<i>No-nest</i>		6.75935	66.48558	13.51871	0.05278	1.80356	0.90172
Encoder & Decoder	<i>add</i>	6.68274	67.45593	13.36548	0.09209	<i>1.83367</i>	0.92831
	<i>max</i>	6.71760	<i>92.49952</i>	13.43519	0.21454	1.58628	0.77823
	<i>l₁-norm</i>	<i>6.83073</i>	93.21573	<i>13.66146</i>	0.20760	1.56378	0.76769
	<i>l_*-norm</i>	6.81192	73.66134	13.62385	<i>0.09010</i>	<i>1.80934</i>	<i>0.92628</i>
	<i>SCA</i>	6.91971	<i>82.75242</i>	13.83942	0.13405	1.73353	0.86248
RFN-Nest		<i>6.84134</i>	71.90131	<i>13.68269</i>	<i>0.07288</i>	1.83676	<i>0.91456</i>

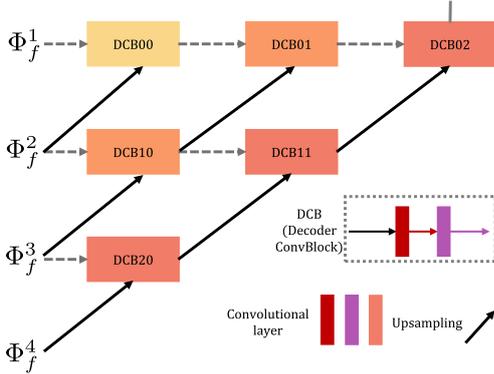


Figure 12: The decoder network without short connection (“No-nest”).

texture. The best values, the second-best values and the third-best values are indicated in **bold**, *red and italic* and *blue and italic*, respectively.

Compared with “No-nest”, the RFN-Nest (the decoder with nest connection) obtains one best metrics value, three second-best metrics values and one third-best metrics value. This indicates that the nest connection architecture plays an important role in boosting the reconstruction ability of the decoder network. With the nest connection, the decoder is able to preserve more image information conveyed by the multiscale deep features (*MI*, *FFMI_{dct}*, *FFMI_w*) and generate more natural and clearer fused image (*EN*, *SD*, *VIF*).

4.7. Ablation Study for Fusion Strategy

In this section, we analyze the importance of RFN as an adaptive fusion mechanism in our fusion network. We choose five classical handcrafted fusion strategies (“add”, “max”, “l₁-norm”, “l_{*}-norm” and “SCA”) which are used in existing fusion networks [25][38][32] to do the experiments.

The trained encoder and decoder are utilized to extract the multi-scale deep features and generate the final image from the fused features, respectively.

Let Φ_{ir}^m and Φ_{vi}^m denote the multi-scale deep features extracted by the trained encoder from the infrared and visible image, respectively. Φ_f^m are the fused deep features. m indicates the scale of the deep features. The formulas of these five strategies are shown in Table 4.

“add”, means the fused features are obtained by adding the source features, directly. In “max” strategy, $\max(\cdot)$ denotes an

Table 4: The formulas of different fusion strategies.

Fusion Strategy	Formula
<i>add</i> [25]	$\Phi_f^m = \Phi_{ir}^m + \Phi_{vi}^m$
<i>max</i> [38]	$\Phi_f^m = \max(\Phi_{ir}^m, \Phi_{vi}^m)$
<i>l₁-norm</i> [25]	$\Phi_f^m = l_1(\Phi_{ir}^m, \Phi_{vi}^m)$
<i>l_*-norm</i> [32]	$\Phi_f^m = l_*(\Phi_{ir}^m, \Phi_{vi}^m)$
<i>SCA</i> [32]	Spatial and channel attention fusion strategy

element wise choose-max strategy [38]. For the “l₁-norm” strategy, $l_1(\cdot)$, the weights are calculated based on l₁-norm. For details on how to calculate these weights, please refer to [25]. For “l_{*}-norm” (known as nuclear-norm), $l_*(\cdot)$ calculates the sum of singular values of a matrix involved in the global pooling operation of deep features to obtain the fusion weights.

The “SCA” indicates the Spatial/Channel Attention fusion strategy which was utilized in NestFuse [32]. In this experiment³, the l₁-norm is used to do the spatial attention fusion and the average pooling is utilized to calculate the channel attention.

Some examples of the images fused using different fusion strategies are shown in Fig.13. Compared with other handcrafted fusion strategies, the fused image obtained by the RFN-based network preserves more detail information from visible image (blue boxes) and the fused image contains less artefacts (red boxes).

The results of fusing 21 pairs of infrared and visible images have been evaluated in terms of the seven quality metrics. The metrics values are shown in Table 3. The table also reports the results obtained with other fusion strategies. The RFN-based network (RFN-Nest) achieves five best values. This indicates that when the learnable fusion network is used as a fusion strategy, the detail image information will be boosted (*En*, *SD*) thanks to the proposed loss function. Regarding feature preservation, the proposed strategy still obtains three best values (*MI*, *FFMI_{dct}* and *FFMI_w*) and two comparable results (*SSIM_a* and *VIF*).

In Section 5, we adopt this learnable fusion network (RFN) for the object tracking task to illustrate the effectiveness of RFN-based fusion strategy in other vision task.

³The “SCA” fusion strategy is used in NestFuse.

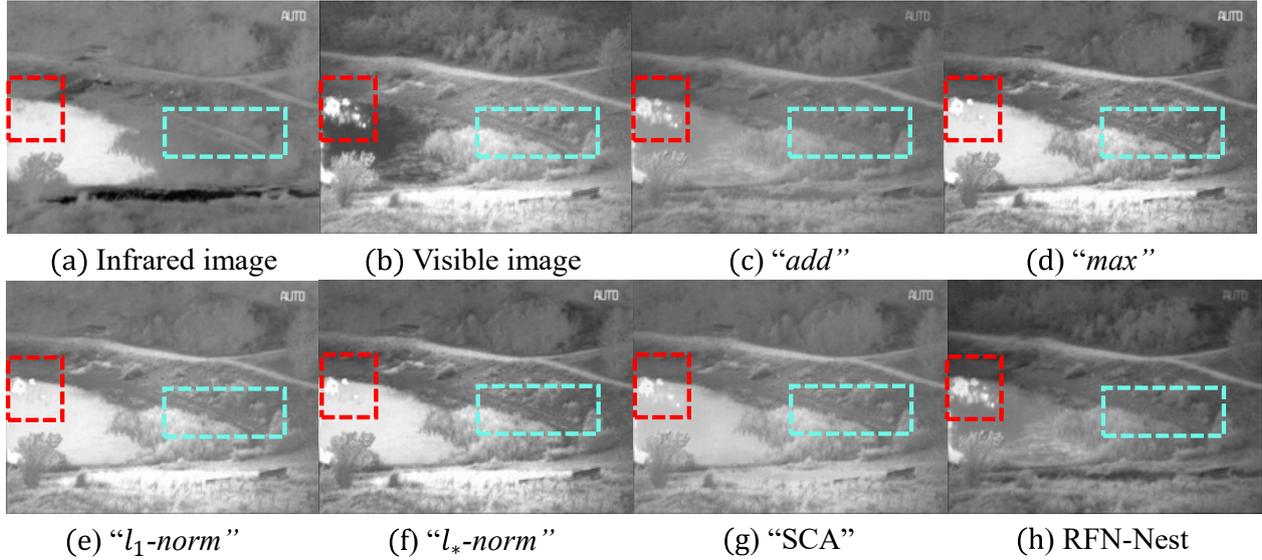


Figure 13: The fusion results with different fusion strategies. “RFN-Nest” means the adaptive RFN networks are utilized in the fusion operation.

4.8. Fusion Results Analysis on 21 pairs Images

To compare the fusion performance of the proposed method with the state-of-the-art algorithms, eleven representative fusion methods are chosen, including discrete cosine harmonic wavelet transform(DCHWT) [50], gradient transfer and total variation minimization(GTF) [53], convolutional sparse representation(ConvSR) [27], multi-layer deep features fusion method(VggML) [22], DenseFuse [25], FusionGAN [29], IFCNN [38] (elementwise-maximum), NestFuse [32], PMGI [39], DDcGAN [31] and U2Fusion [40].

For DenseFuse, we choose the sum strategy and set the trade-off parameter to 1×10^2 . For NestFuse, the average pooling is utilized for the channel attention fusion strategy. All these fusion methods are implemented using publicly available codes, and their parameters are set by referring to their original reports.

To evaluate the visual effects of the fusion results⁴, two pairs of visible and infrared images are selected, namely “man” and “umbrella”. The fused images obtained by the existing methods and our fusion method (RFN-Nest) are shown in Fig.14 and Fig.15, respectively.

In Fig.14 and Fig.15, the fused images obtained by DCHWT are more noisy and contain image artefacts. The fused images obtained by GTF and FusionGAN exhibit clearer features and more detailed background information. Although these fused images retain more complementary information, they look more like the infrared image, especially in the background. In view of the importance of the background information, ConvSR, VggML, DenseFuse, IFCNN, NestFuse, PMGI and U2Fusion are designed to preserve more detail information from the visible image. These methods appear to reduce the salient infrared features, compared with GTF and FusionGAN, producing acceptable fusion results. Although DDcGAN is also

designed to maintain more detail information from visible images, in Fig.14 (l), it injects more noise into the fused image and the infrared targets are blurred.

Although the target features are not enhanced too much in the fused image, the contrast is better than in the visible image. Moreover, for the detail preservation, in Fig.14, compared with the other fusion methods, the ‘tree’ and ‘street lamp’ (red box) are clearer in the fused image obtained by our proposed method. The detail textures of ‘bushes’ (green box) are also preserved into the fused image.

In Fig.15, in the green box, many fusion methods are unable to preserve the salient features of ‘pavilion’ from the visible image except IFCNN, NestFuse and the proposed method, which means these fusion methods fuse too much background information from the infrared image. Compared with all these fusion methods, in the red box, the detail information of the ‘tree’ reconstructed by the proposed method is clearer in the fused image (Fig.15 (n)).

The background and the context around salient parts are not very clear and sometime even invisible because of the difficulty in extracting salient features from source images, as shown in Fig.14 and Fig.15. This drawback will cause a performance degradation when the image fusion algorithms are used in other computer vision tasks, such as RGB-T visual object tracking. In contrast, our RFN-Nest fusion network is able to preserve more detail information and to maintain the contrast of infrared parts.

Compared with all the above fusion methods, the fused image obtained by the proposed method appears to retain a better balance between the visible background information and the infrared features.

We evaluate the fusion performance objectively using the seven quality metrics to compare the seventeen existing fusion methods and our proposed fusion framework. The values of these metrics averaged over all fused images are shown in Ta-

⁴More experimental results are shown in our supplementary material.

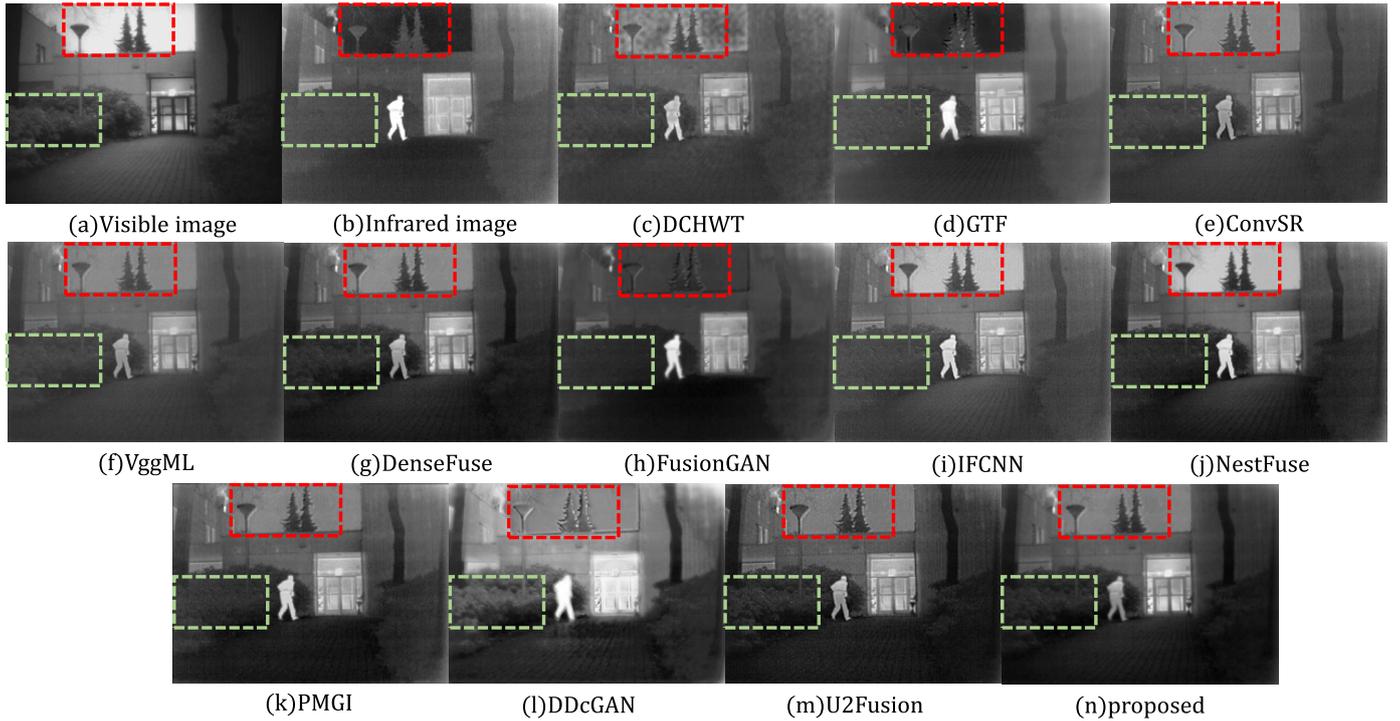


Figure 14: The experimental results on “man” images. (a) Visible; (b) Infrared; (c) DCHWT; (d) GTF; (e) ConvSR; (f) VggML; (g) DenseFuse; (h) FusionGAN; (i) IFCNN; (j) NestFuse; (k) PMGI; (l) DDcGAN; (m) U2Fusion; (n) proposed.

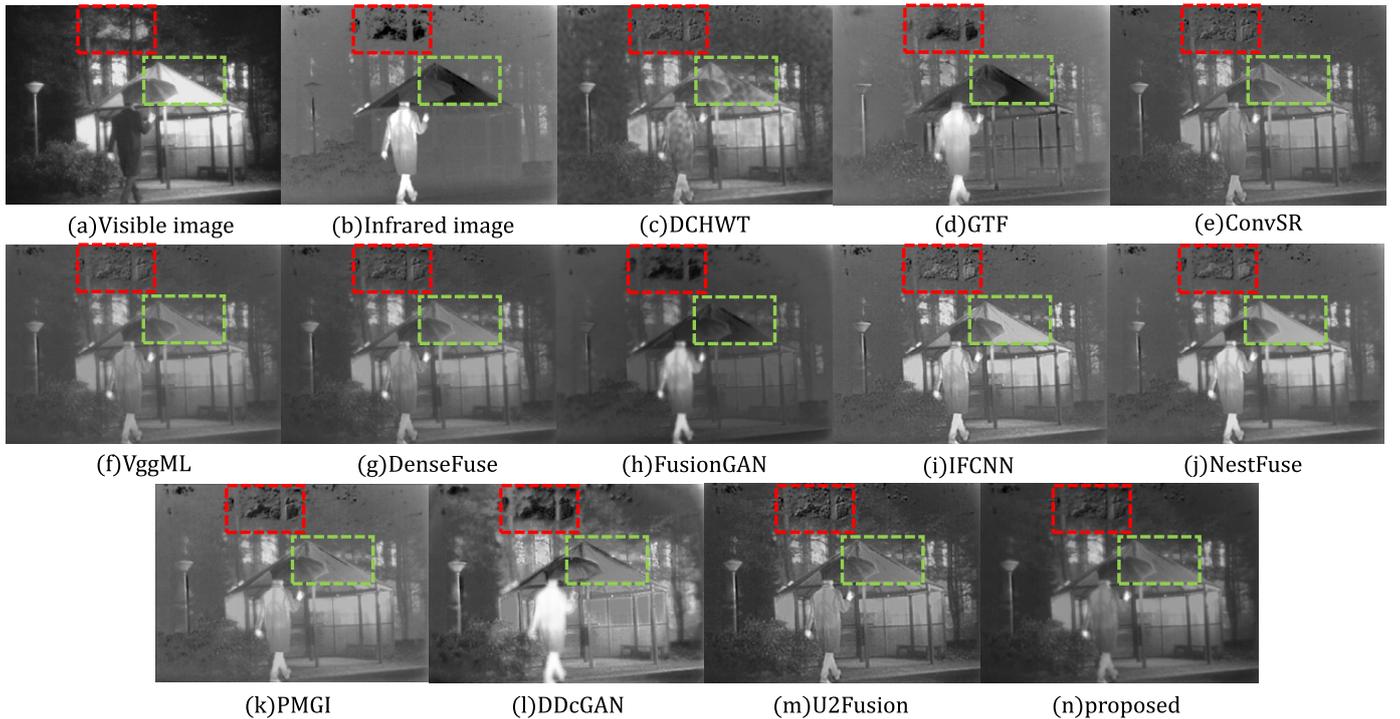


Figure 15: The experimental results on “umbrella” images. (a) Visible; (b) Infrared; (c) DCHWT; (d) GTF; (e) ConvSR; (f) VggML; (g) DenseFuse; (h) FusionGAN; (i) IFCNN; (j) NestFuse; (k) PMGI; (l) DDcGAN; (m) U2Fusion; (n) proposed.

ble 5. The best values, the second-best values and the third-best values are indicated in **bold**, *red and italic* and *blue and italic*, respectively.

From Table 5, the proposed fusion framework (RFN-Nest) obtains one best values (SCD) and three third-best values (SD , N_{abf} , $MS\text{-}SSIM$) compared to the other methods. The reason

Table 5: The values of seven quality metrics averaged over the fused images on 21 pairs of infrared and visible images.

	En [47]	SD [48]	MI [49]	N_{abf} [50]	SCD [51]	MS-SSIM[52]
DCHWT[50]	6.56777	64.97891	13.13553	0.12295	1.60993	0.84326
GTF[53]	6.63433	67.54361	13.26865	0.07951	1.00488	0.80844
ConvSR[27]	6.25869	50.74372	12.51737	0.01958	1.64823	0.90281
VggML[22]	6.18260	48.15779	12.36521	0.00120	1.63522	0.87478
DenseFuse[25]	6.67158	67.57282	13.34317	0.09214	<i>1.83502</i>	0.92896
FusionGan[29]	6.36285	54.35752	12.72570	<i>0.06706</i>	1.45685	0.73182
IFCNN[38]	6.59545	66.87578	13.19090	0.17959	1.71375	0.90527
NestFuse[32]	<i>6.91971</i>	<i>82.75242</i>	<i>13.83942</i>	0.13405	1.73353	0.86248
PMGI[39]	<i>6.93391</i>	71.54806	<i>13.86783</i>	0.13525	1.78242	0.88934
DDcGAN[31]	7.47310	100.34809	14.94620	0.33784	1.60926	0.76636
U2Fusion[40]	6.75708	64.91158	13.51416	0.29088	<i>1.79837</i>	<i>0.92533</i>
<i>proposed</i>	6.84134	<i>71.90131</i>	13.68269	<i>0.07288</i>	1.83676	<i>0.91456</i>

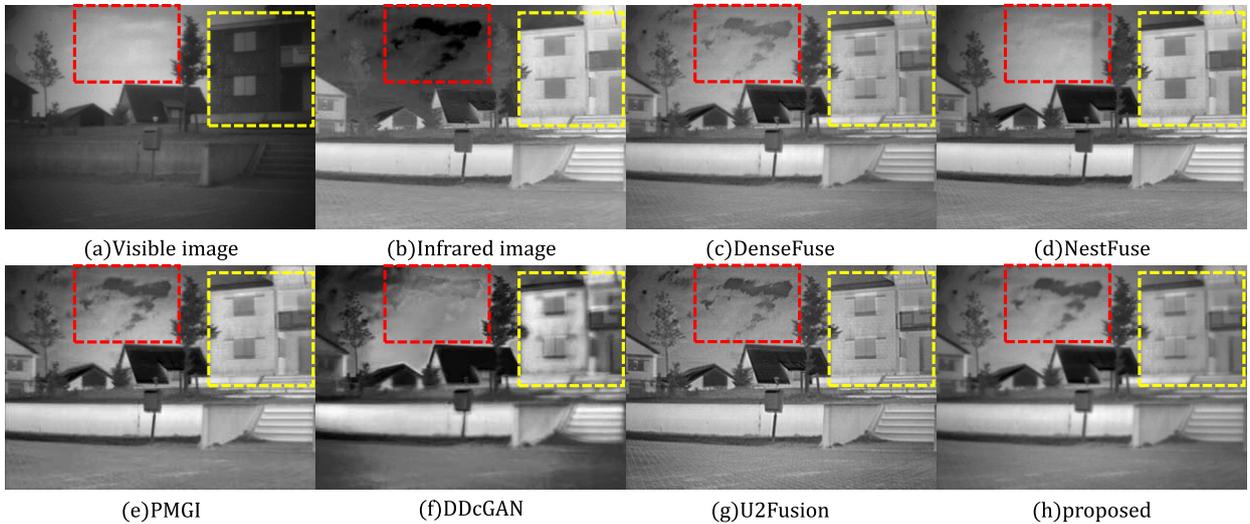


Figure 16: The experimental results on “street” images. (a) Visible; (b) Infrared; (c) DenseFuse; (d) NestFuse; (e) PMGI; (f) DDcGAN; (g) U2Fusion; (h) proposed.

why DDcGAN obtains larger values of En , SD and MI is that DDcGAN introduces more noise and artefacts into the fused image. Our fusion network achieves good fusion performance, producing sharper content and exhibiting more visual information fidelity.

4.9. Further Analysis on 40 Pairs Images

The previous ablation studies and experiments are conducted on one test dataset which contains 21 pairs of infrared and visible images. To verify the generalization performance of the proposed fusion network, a new test dataset is created. It contains 40 pairs of infrared and visible images which are collected from TNO [44] and VOT2020-RGBT [45].

In this section, we choose several state-of-the-art deep learning based fusion methods to perform comparative experiments. These methods include DenseFuse [25] which is a classical autoencoder-based fusion method, NestFuse [32] which has the same backbone (encoder and decoder) with the proposed method, and three latest fusion methods (PMGI [39], DDcGAN [31] and U2Fusion [40]).

An example of the fused images obtained by these fusion methods and the proposed network is shown in Fig.16. The unnatural textures in the sky are introduced by the infrared image (Fig.16 (b), red box) should be looked more natural in the fused image. It is observed that DenseFuse, PMGI and the proposed method generate fused images of natural appearance. Moreover, compared with the existing fusion methods, our network also preserves more detail information from both infrared and visible images (Fig.16, the house in yellow box). Note, the DDcGAN again introduces noise into the fused image and blurs the salient feature content.

The same six quality metrics are used for comparative evaluation. The average values of these metrics are shown in Table 6. The best values are indicated in **bold**, the second-best values are denoted in *red and italic* and the third-best values are denoted in *blue and italic*.

Compared with the results on the 21 pairs of images, the proposed network exhibits even better performance on the 40 image pairs. The method achieves one best value (SCD), one second-best values (N_{abf}) and two third-best values (SD , MS-

Table 6: The values of seven quality metrics averaged over the fused images on 40 pairs of infrared and visible images which collected from TNO and VOT2020.

	En [47]	SD [48]	MI [49]	N_{abf} [50]	SCD [51]	MS-SSIM[52]
DenseFuse[25]	6.77630	73.63462	13.55261	0.06346	<i>1.74862</i>	<i>0.92944</i>
NestFuse[32]	<i>6.99347</i>	<i>90.28951</i>	<i>13.98693</i>	<i>0.11138</i>	1.67540	0.88611
PMGI[39]	<i>6.96974</i>	77.25462	<i>13.93948</i>	0.11434	1.68523	0.88830
DDcGAN[31]	7.50173	106.99113	15.00346	0.30998	1.55359	0.78419
U2Fusion[40]	6.94970	76.80347	13.89939	0.28363	<i>1.74780</i>	0.93141
<i>proposed</i>	6.92952	<i>78.22247</i>	13.85904	<i>0.06357</i>	1.76116	<i>0.90894</i>

SSIM). Even compared with DDcGAN, our fusion network performance is comparable. This confirms that our fusion network trained by the two-stage fusion strategy and the novel loss function demonstrates better generalization.

5. Experiments on RGBT Object Tracking

Over the past two years, multi-modality object tracking has been of interest in many vision applications. In Vision Object Tracking challenge (VOT) 2019 [54], for the first time, the committee introduced two new sub-challenges (RGBD and RGBT), in which each sequence in the dataset of RGBD or RGBT contains two modalities (RGB image and depth image, RGB image and infrared image) as the input. As we focus on the fusion of infrared and visible images, the RGBT sub-challenge data is used to evaluate the performance of the proposed learnable fusion network (RFN) and the novel loss functions.

In VOT2020 [45], the video sequences are the same as VOT2019 [54], but a new performance evaluation protocol is introduced for short-term tracker evaluation (includes RGBT sub-challenge). The new protocol avoids tracker-dependent resets and reduces the variance of the performance evaluation measures.

A state-of-the-art siamese-based tracker AFAT [55] is chosen to be the base tracker. In AFAT, a failure-aware system, realized by a Quality Prediction Network (QPN), based on convolutional and LSTM modules was proposed and obtained better tracking performance in many datasets. For RGBT object tracking, the proposed fusion strategy network (RFN) and the proposed loss function are incorporated into AFAT.

5.1. The RFN and The Loss Function For RGBT Tracking

As we discussed, the proposed residual fusion network (RFN) is a learnable fusion strategy. Thus, ideally, when RFN is applied into AFAT [55], it needs a sufficient quantity of data to train the whole model.

However, due to the lack of labeled training data, we were forced to simplify the architecture of RFN, by reducing the number of convolutional layers, as shown in Fig.17. In the training phase, we only train the RFN module, the AFAT modules are fixed to reduce the number of learnable parameters ⁵.

⁵The framework of RFN-base AFAT is shown in our supplementary material.

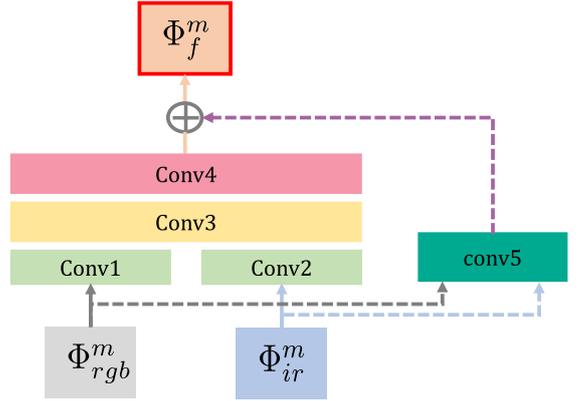


Figure 17: The RFN architecture for RGBT tracking.

Three RGBT datasets are used to train our RFN module, namely GTOT [56], VT821 [57], VT1000 [58]. These datasets only contain 17.6k frames in total. GTOT is a dataset for RGBT tracking, the VT821 and the VT1000 are built for RGBT salient object detection.

To train the RFN module, the proposed loss function (L_{RFN} , Section 3.2.2) is used in the AFAT training. As RGBT tracking does not involve image generation, it is inevitable that the background detail preservation loss function (L_{detail}) needs to be modified to become applicable to the tracking task. L_{detail} is defined as follows,

$$L_{detail} = 1 - SSIM(\Phi_f^m, l_1(\Phi_{rgb}^m, \Phi_{ir}^m)) \quad (7)$$

where Φ_f^m denotes the fused deep feature obtained by the RFN module, and l_1 indicates the “ l_1 -norm” based fusion strategy discussed in Section 4.7. The target feature enhancement loss function ($L_{feature}$) is the same as in Section 3.2.2.

5.2. The Tracking Results on VOT-RGBT

The video sequences in VOT2020-RGBT are the same as in VOT2019-RGBT. Thus, we only present a few tracking results on VOT2020-RGBT in Fig.18. The ‘RFN’ denotes the RFN-based AFAT.

To evaluate the tracking performance, three measures [54] were selected: Accuracy (A), Robustness (R) and Expected Average Overlap (EAO). (1) Accuracy denotes the average overlap between the ground truth and the predicted bounding boxes; (2) Robustness evaluates how many times the tracker loses the

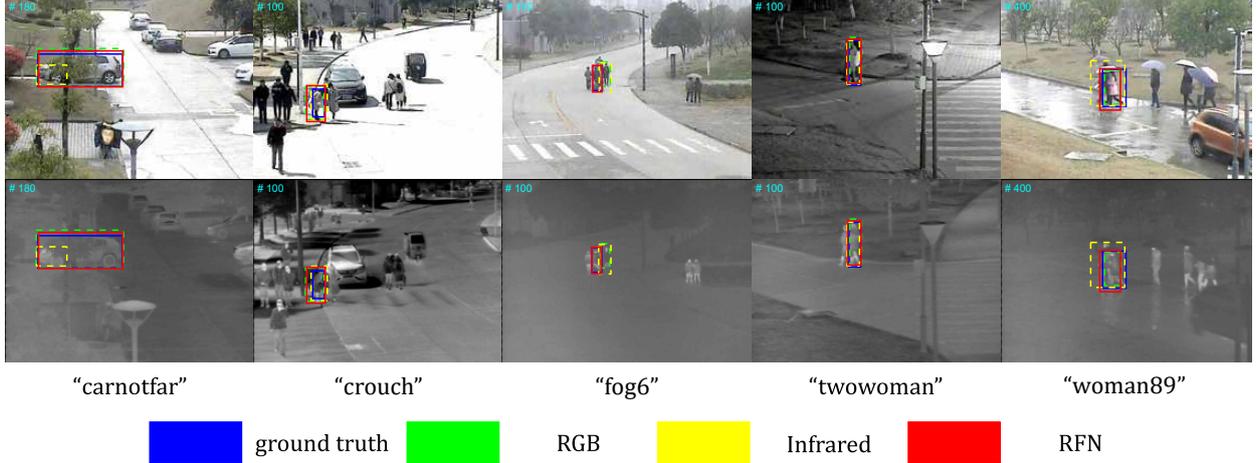


Figure 18: Some tracking results on VOT2020-RGBT. The frames in first row and second row are RGB frames and infrared frames, respectively. Five pairs of RGB and infrared frames in columns are collected from different video sequences (“carnotfar”, “crouch”, “fog6”, “twowoman” and “woman89”). AFAT is the base tracker. In these frames, blue boxes denote the ground truth, green boxes and yellow boxes are the tracking results obtained by AFAT with only RGB frames or infrared frames fed as input. Red boxes indicate the tracking results obtained by the RFN-based tracker.

target (fails) during tracking (3) EAO is an estimator of the average overlap of a tracker. For the detail of Accuracy, Robustness and EAO please refer to [59].

In VOT2020-RGBT [45]⁶, for the Accuracy, Robustness and EAO , they have the same meanings but their calculation methods are re-defined by the committee. Thus, these metrics are indicated as EAO_{new} , A_{new} and R_{new} . The higher values of EAO , EAO_{new} , A , A_{new} , R and R_{new} , the better the tracker.

In addition to the base tracker (AFAT), we choose two further trackers for each dataset (VOT2019, VOT2020) to analyze the tracking performance of RFN-based AFAT. In the VOT2019-RGBT competition, mfDiMP and FSRPN won the third and fourth place on the public dataset, respectively. For VOT2020-RGBT, DFAT and M2C2Frgbt won third and seventh place on the public dataset, respectively. Note that DFAT is the winner on VOT2020-RGBT challenge. All the metrics values are provided by the VOT committee and available on the VOT reports [54] [45].

The tracking results of RFN-based AFAT and other trackers are shown in Table 7 and Table 8. *RGB* and *Infrared* indicate that only one modality (RGB or infrared) is fed into AFAT.

Table 7: The tracking results obtained on the VOT2019-RGBT dataset. AFAT is used as the base tracker to evaluate the RFN network fusion strategy.

VOT2019		EAO	A	R
FSRPN[54]		0.3553	<i>0.6362</i>	<i>0.7069</i>
mfDiMP[54]		0.3879	0.6019	0.8036
AFAT[55]	<i>RGB</i>	0.32590	0.61130	0.5700
	<i>Infrared</i>	0.18120	0.56740	0.1800
RFN-based		<i>0.35840</i>	0.64470	0.6500

From these two tables, compared with just feeding one modality into AFAT, the RFN-based AFAT delivers better tracking performance in all measures both on VOT-RGBT2019 and

Table 8: The tracking results obtained on the VOT2020-RGBT dataset. AFAT is set as the base tracker to evaluate the RFN network fusion strategy.

VOT2020		EAO_{new}	A_{new}	R_{new}
M2C2Frgbt[45]		0.332	0.636	0.722
DFAT[45]		0.390	0.672	0.779
AFAT[55]	<i>RGB</i>	0.329	0.635	0.669
	<i>Infrared</i>	0.265	0.573	0.588
RFN-based		<i>0.371</i>	<i>0.668</i>	<i>0.726</i>

on VOT-RGBT2020. On VOT2019-RGBT, although mfDiMP achieves the best performance, the results produced by the RFN-based tracker are comparable (EAO) and the accuracy is better. On VOT2020-RGBT, even compared with the winning tracker, DFAT, our tracker is also competitive.

These experiments demonstrate that even with insufficient training data, the tracker performance is improved by incorporating the proposed residual fusion network (RFN) into the AFAT tracking framework. When more training data becomes available, we believe the RFN-based tracker will achieve even better tracking performance.

6. Conclusions

Motivated by the weakness of the existing fusion methods in preserving image detail, in this paper, we proposed a novel end-to-end fusion framework (RFN-Nest) which is based on the nest connection incorporated into a residual fusion network. To design our RFN-Nest, a two-stage training strategy was presented. In the proposed scheme, an auto-encoder network is trained using the SSIM loss function (L_{ssim}) and the pixel loss function (L_{pixel}). The trained encoder is utilized to extract multi-scale features from the source images and the nest connection-based decoder network is designed to reconstruct the fused images using the fused multi-scale features. The key component of RFN-

⁶The toolkit version of VOT2020-RGBT is 0.2.0.

Nest is the residual fusion network (RFN). In the second stage of the training strategy, four residual fusion networks (RFN) are trained to preserve the image detail, and preserve the salient features using L_{detail} and $L_{feature}$, respectively. Once the two-stage training is accomplished, the fused image is reconstructed using the encoder, the RFN networks and the decoder. Compared to seventeen existing fusion methods, the RFN-Nest achieves the best fusion performance in both subjective and objective evaluation.

To validate the generality of the fusion network, we also applied the proposed RFN and the novel loss functions to a state-of-the-art tracker to perform a multimodal tracking task (RGBT tracking). Compared with single modality, the RFN-based tracker delivers better tracking performance in all measures on VOT2019 and VOT2020. Even compared with the state of the art RGBT trackers, the RFN-based tracker achieves very good performance. This demonstrates that with this proposed innovations, the RFN-Nest network has a wide applicability, extending beyond image fusion.

Acknowledgement

This work was supported by the National Natural Science Foundation of China (62020106012, U1836218, 61672265), the 111 Project of Ministry of Education of China (B12018), and the Engineering and Physical Sciences Research Council (EPSRC) (EP/N007743/1, MURI/EPSRC/DSTL, EP/R018456/1).

References

- [1] C. Li, X. Liang, Y. Lu, N. Zhao, J. Tang, RGB-T object tracking: benchmark and baseline, *Pattern Recognition* 96 (2019) 106977.
- [2] C. Li, C. Zhu, J. Zhang, B. Luo, X. Wu, J. Tang, Learning Local-Global Multi-Graph Descriptors for RGB-T Object Tracking, *IEEE Transactions on Circuits and Systems for Video Technology* (2018).
- [3] C. Luo, B. Sun, K. Yang, T. Lu, W.-C. Yeh, Thermal infrared and visible sequences fusion tracking based on a hybrid tracking framework with adaptive weighting scheme, *Infrared Physics & Technology* 99 (2019) 265–276.
- [4] V. Shrinidhi, P. Yadav, N. Venkateswaran, IR and Visible Video Fusion for Surveillance, in: 2018 International Conference on Wireless Communications, Signal Processing and Networking (WiSPNET), IEEE, 2018, pp. 1–6.
- [5] J. Ma, Y. Ma, C. Li, Infrared and visible image fusion methods and applications: A survey, *Information Fusion* 45 (2019) 153–178.
- [6] S. Li, X. Kang, L. Fang, J. Hu, H. Yin, Pixel-level image fusion: A survey of the state of the art, *Information Fusion* 33 (2017) 100–112.
- [7] Y. Liu, X. Chen, Z. Wang, Z. J. Wang, R. K. Ward, X. Wang, Deep learning for pixel-level image fusion: Recent advances and future prospects, *Information Fusion* 42 (2018) 158–173.
- [8] G. Pajares, J. M. De La Cruz, A wavelet-based image fusion tutorial, *Pattern recognition* 37 (2004) 1855–1872.
- [9] A. Ben Hamza, Y. He, H. Krim, A. Willisky, A multiscale approach to pixel-level image fusion, *Integrated Computer-Aided Engineering* 12 (2005) 135–146.
- [10] S. Yang, M. Wang, L. Jiao, R. Wu, Z. Wang, Image fusion based on a new contourlet packet, *Information Fusion* 11 (2010) 78–84.
- [11] S. Li, X. Kang, J. Hu, Image fusion with guided filtering, *IEEE Transactions on Image processing* 22 (2013) 2864–2875.
- [12] J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry, Y. Ma, Robust face recognition via sparse representation, *IEEE transactions on pattern analysis and machine intelligence* 31 (2008) 210–227.
- [13] G. Liu, Z. Lin, Y. Yu, Robust subspace segmentation by low-rank representation., in: *ICML*, volume 1, 2010, p. 8.
- [14] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, Y. Ma, Robust recovery of subspace structures by low-rank representation, *IEEE transactions on pattern analysis and machine intelligence* 35 (2012) 171–184.
- [15] Q. Zhang, Y. Fu, H. Li, J. Zou, Dictionary learning method for joint sparse representation-based image fusion, *Optical Engineering* 52 (2013) 057006.
- [16] C. Liu, Y. Qi, W. Ding, Infrared and visible image fusion method based on saliency detection in sparse domain, *Infrared Physics & Technology* 83 (2017) 94–102.
- [17] R. Gao, S. A. Vorobyov, H. Zhao, Image fusion with cosparsity analysis operator, *IEEE Signal Processing Letters* 24 (2017) 943–947.
- [18] H. Li, X.-J. Wu, Multi-focus image fusion using dictionary learning and low-rank representation, in: *International Conference on Image and Graphics*, Cham, Switzerland: Springer, 2017, pp. 675–686.
- [19] X. Lu, B. Zhang, Y. Zhao, H. Liu, H. Pei, The infrared and visible image fusion algorithm based on target separation and sparse representation, *Infrared Physics & Technology* 67 (2014) 397–407.
- [20] M. Yin, P. Duan, W. Liu, X. Liang, A novel infrared and visible image fusion algorithm based on shift-invariant dual-tree complex shearlet transform and sparse representation, *Neurocomputing* 226 (2017) 182–191.
- [21] C. Liu, Y. Qi, W. Ding, Infrared and visible image fusion method based on saliency detection in sparse domain, *Infrared Physics & Technology* 83 (2017) 94–102.
- [22] H. Li, X.-J. Wu, J. Kittler, Infrared and Visible Image Fusion using a Deep Learning Framework, in: 2018 24th International Conference on Pattern Recognition (ICPR), IEEE, 2018, pp. 2705–2710.
- [23] H. Li, X.-J. Wu, T. S. Durrani, Infrared and Visible Image Fusion with ResNet and zero-phase component analysis, *Infrared Physics & Technology* 102 (2019) 103039.
- [24] X. Song, X.-J. Wu, Multi-focus Image Fusion with PCA Filters of PCANet, in: *IAPR Workshop on Multimodal Pattern Recognition of Social Signals in Human-Computer Interaction*, Springer, 2018, pp. 1–17.
- [25] H. Li, X.-J. Wu, DenseFuse: A Fusion Approach to Infrared and Visible Images, *IEEE Transactions on Image Processing* 28 (2019) 2614–2623.
- [26] H. Li, X.-J. Wu, J. Kittler, MDLatLRR: A novel decomposition method for infrared and visible image fusion, *IEEE Transactions on Image Processing* (2020). Doi: 10.1109/TIP.2020.2975984.
- [27] Y. Liu, X. Chen, R. K. Ward, Z. J. Wang, Image fusion with convolutional sparse representation, *IEEE signal processing letters* 23 (2016) 1882–1886.
- [28] Y. Liu, X. Chen, H. Peng, Z. Wang, Multi-focus image fusion with a deep convolutional neural network, *Information Fusion* 36 (2017) 191–207.
- [29] J. Ma, W. Yu, P. Liang, C. Li, J. Jiang, FusionGAN: A generative adversarial network for infrared and visible image fusion, *Information Fusion* 48 (2019) 11–26.
- [30] J. Ma, P. Liang, W. Yu, C. Chen, X. Guo, J. Wu, J. Jiang, Infrared and visible image fusion via detail preserving adversarial learning, *Information Fusion* 54 (2020) 85–98.
- [31] J. Ma, H. Xu, J. Jiang, X. Mei, X.-P. Zhang, DDcGAN: A Dual-Discriminator Conditional Generative Adversarial Network for Multi-Resolution Image Fusion, *IEEE Transactions on Image Processing* 29 (2020) 4980–4995.
- [32] H. Li, X.-J. Wu, T. Durrani, NestFuse: An Infrared and Visible Image Fusion Architecture based on Nest Connection and Spatial/Channel Attention Models, *IEEE Transactions on Instrumentation and Measurement* (2020). Doi: 10.1109/TIM.2020.3005230.
- [33] Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh, J. Liang, Unet++: A nested u-net architecture for medical image segmentation, in: *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*, Granada, Spain, Springer, 2018, pp. 3–11.
- [34] K. Ram Prabhakar, V. Sai Srikar, R. Venkatesh Babu, Deepfuse: a deep unsupervised approach for exposure fusion with extreme exposure image pairs, in: *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 4714–4722.
- [35] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, *arXiv preprint arXiv:1409.1556* (2014).
- [36] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

- [37] G. Huang, Z. Liu, L. Van Der Maaten, K. Q. Weinberger, Densely connected convolutional networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 4700–4708.
- [38] Y. Zhang, Y. Liu, P. Sun, H. Yan, X. Zhao, L. Zhang, IFCNN: A general image fusion framework based on convolutional neural network, *Information Fusion* 54 (2020) 99–118.
- [39] H. Zhang, H. Xu, Y. Xiao, X. Guo, J. Ma, Rethinking the image fusion: A fast unified image fusion network based on proportional maintenance of gradient and intensity, in: Proceedings of the AAAI Conference on Artificial Intelligence, 2020, pp. 12797–12804.
- [40] H. Xu, J. Ma, J. Jiang, X. Guo, H. Ling, U2Fusion: A Unified Unsupervised Image Fusion Network, *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2020). Doi: 10.1109/TPAMI.2020.3012548.
- [41] Z. Wang, A. C. Bovik, H. R. Sheikh, E. P. Simoncelli, et al., Image quality assessment: from error visibility to structural similarity, *IEEE transactions on image processing* 13 (2004) 600–612.
- [42] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, C. L. Zitnick, Microsoft coco: Common objects in context, in: European conference on computer vision, Zurich, Switzerland, Springer, 2014, pp. 740–755.
- [43] S. Hwang, J. Park, N. Kim, Y. Choi, I. So Kweon, Multispectral pedestrian detection: Benchmark dataset and baseline, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 1037–1045.
- [44] A. Toet, TNO Image Fusion Dataset, 2014. https://figshare.com/articles/TN_Image_Fusion_Dataset/1008029.
- [45] M. Kristan, J. Matas, A. Leonardis, M. Felsberg, et al., The eighth visual object tracking VOT2020 challenge results, in: Proc. 16th Eur. Conf. Comput. Vis. Workshop, 2020.
- [46] H. Li, Code of RFN-Nest, 2020. <https://github.com/hli1221/imagefusion-rfn-nest>.
- [47] J. W. Roberts, J. A. Van Aardt, F. B. Ahmed, Assessment of image fusion procedures using entropy, image quality, and multispectral classification, *Journal of Applied Remote Sensing* 2 (2008) 023522.
- [48] Y.-J. Rao, In-fibre Bragg grating sensors, *Measurement science and technology* 8 (1997) 355.
- [49] G. Qu, D. Zhang, P. Yan, Information measure for performance of image fusion, *Electronics letters* 38 (2002) 313–315.
- [50] B. S. Kumar, Multifocus and multispectral image fusion based on pixel significance using discrete cosine harmonic wavelet transform, *Signal, Image and Video Processing* 7 (2013) 1125–1143.
- [51] V. Aslantas, E. Bendes, A new image quality metric for image fusion: The sum of the correlations of differences, *Aeu-international Journal of electronics and communications* 69 (2015) 1890–1896.
- [52] K. Ma, K. Zeng, Z. Wang, Perceptual quality assessment for multi-exposure image fusion, *IEEE Transactions on Image Processing* 24 (2015) 3345–3356.
- [53] J. Ma, C. Chen, C. Li, J. Huang, Infrared and visible image fusion via gradient transfer and total variation minimization, *Information Fusion* 31 (2016) 100–109.
- [54] M. Kristan, J. Matas, A. Leonardis, M. Felsberg, R. Pflugfelder, J.-K. Kamarainen, L. Cehovin Zajc, O. Drbohlav, A. Lukezic, A. Berg, et al., The seventh visual object tracking vot2019 challenge results, in: Proceedings of the IEEE International Conference on Computer Vision Workshops, 2019, pp. 1–36.
- [55] T. Xu, Z.-H. Feng, X.-J. Wu, J. Kittler, AFAT: Adaptive Failure-Aware Tracker for Robust Visual Object Tracking, *arXiv preprint arXiv:2005.13708v1* (2020).
- [56] C. Li, H. Cheng, S. Hu, X. Liu, J. Tang, L. Lin, Learning collaborative sparse representation for grayscale-thermal tracking, *IEEE Transactions on Image Processing* 25 (2016) 5743–5756.
- [57] J. Tang, D. Fan, X. Wang, Z. Tu, C. Li, RGBT Salient Object Detection: Benchmark and A Novel Cooperative Ranking Approach, *IEEE Transactions on Circuits and Systems for Video Technology* (2019).
- [58] Z. Tu, T. Xia, C. Li, X. Wang, Y. Ma, J. Tang, RGB-T Image Saliency Detection via Collaborative Graph Learning, *IEEE Transactions on Multimedia* 22 (2019) 160–173.
- [59] M. Kristan, J. Matas, A. Leonardis, M. Felsberg, L. Cehovin, G. Fernandez, T. Vojir, G. Hager, G. Nebehay, R. Pflugfelder, The visual object tracking vot2015 challenge results, in: Proceedings of the IEEE international conference on computer vision workshops, 2015, pp. 1–23.