

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

# DSTFuse: Enhancing Deblurring via Style Transfer for Visible and Infrared Image Fusion

Anonymous WACV Algorithms Track submission

Paper ID \*\*\*\*\*

## Abstract

Infrared and visible image fusion aims at obtaining fused images that keep advantages of source images, e.g., detailed textures and clear edge structures. To tackle the challenge in modeling features from visible image under motion blur and low light conditions, we propose a novel fusion framework, DSTFuse, which aims to leverage infrared image as the style image and enable it to perform style transfer on the visible image to efficiently eliminate motion blur. Specifically, DSTFuse contains a Cross-Modality Style Transfer Module (CST-module) that collect appropriate style information from the infrared image and guide the transformation of blurry objects into the corresponding style while preserve all other elements without alteration. The output of CST-module is integrated with the image with a multitude of visible features from another module and mapped into final image. Extensive experiments show that DSTFuse achieves promising results in infrared-visible image fusion task. And it is also shown that DSTFuse can boost the performance in downstream infrared-visible object detection. Code will be released at <https://anonymous.4open.science/r/DSTFuse-0C1D>

## 1. Introduction

Image fusion is a fundamental image enhancement technique. It aims to combine images with distinct modality features into a image that retains the advantage of the source images [1, 33, 34, 48, 49, 53]. One prevalent application of image fusion is the infrared and visible image fusion (IVIF) [30, 38, 39, 41]. Proverbially, visible images can reflect the appearance and color information of objects, while infrared images provide thermal radiation information, characterized by a high contrast between the target and its surroundings. By integrating the complementary information from both visible and infrared images, IVIF generates a fused image that that overcomes the limitations of visible images under environmental constraints and the lack

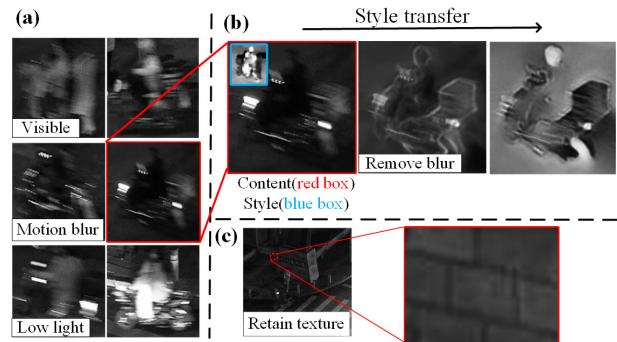


Figure 1. The effect of style transfer on suboptimal visible data for IVIF tasks. (a) The suboptimal visible data with motion blur and low light condition. (b) The blur removing process of transferring visible images into infrared style. The red box and blue box are visible content image and infrared style image, respectively. (c) The effect of texture details retaining by DSTFuse.

of detail in infrared images. Therefore, IVIF has broad applications across various fields such as military [18], security [8], and medical image processing [49].

To tackle the shortcomings of conventional IVIF methods [15, 23, 35], numerous deep learning-based techniques have been developed. These method can be categorized into two main classes: Generative Adversarial (GANs)-based network [3, 31, 32, 45] and the Auto-Encoder-based network [12, 19, 20, 47]. The GAN-based methods typically consist of a generator responsible for generating the fused image and a discriminator that evaluates the fusion performance. And Auto-Encoder-based methods extract features from infrared-visible images through an encoder and then map these features to a new representation space through a decoder. To leverage the the multimodal features, numerous previous works have attempted to map the most representative features of images from different modalities into the final image [10, 26, 39, 46, 54].

Despite a lot of researches on IVIF, there are few studies

108 concentrating on utilizing suboptimal data, especially for  
 109 data containing a significant number of blurry objects. The  
 110 vast majority of studies utilize high-quality datasets such as  
 111 TNO [42] and MSRS [40], which typically do not exhibit  
 112 motion blur (Fig. 1(a)). However, the previous works using  
 113 high-quality datasets have a limitation. Due to the variability  
 114 of real-world environments, motion blur in source images  
 115 is inevitable in practical applications of downstream  
 116 tasks such as detection and segmentation. Therefore, it is  
 117 crucial to mitigate the impact of a large volume of suboptimal  
 118 data on IVIF tasks. Moreover, due to the significantly  
 119 longer exposure time of RGB cameras compared to infrared  
 120 cameras, the quality of infrared images for blurred objects  
 121 in the same scene is superior to that of visible images. It  
 122 is also a significant challenge to utilize the higher tolerance  
 123 to blur that infrared images inherently possess due to differ-  
 124 ences in shooting equipment.  
 125

126 For the source visible and infrared images, the content  
 127 information is intensely correlative. This is attributed to the  
 128 high degree of coincidence in both the scene and the time  
 129 of capture for each pair of infrared-visible images. It is in-  
 130 tuitive that visible images, often prone to blurring due to  
 131 equipment and target movement, have the potential to be  
 132 transformed into consistently sharp infrared images. Previous  
 133 studies on style transfer task have closely aligned with  
 134 this concept.

135 In this paper, we present DSTFuse – a conceptually sim-  
 136 ple framework that aims to enhance deblurring via style  
 137 transfer for IVIF. In DSTFuse, the blurry visible image is  
 138 transformed into an image that combines infrared style with  
 139 visible features by an Auto-Encoder-based cross-modality  
 140 style transfer module (CST-module). Specifically, it aims  
 141 to utilize infrared images as a reference to impose fea-  
 142 ture constraints on the blurry visible images, thus reduc-  
 143 ing motion-induced artifacts and enhancing details. Sub-  
 144 sequently, DSTFuse utilizes the visible-infrared images to  
 145 generate a fused image with rich background information  
 146 and seamlessly integrates it with the output of CST-module  
 147 into a meticulously crafted mapping function. As shown in  
 148 Fig. 1(b), the contours of the blurred object in visible im-  
 149 ages under low-light conditions are gradually outlined, and  
 150 details are filled in as the style transfer process. Moreover,  
 151 the details in the fused image are also remarkably retained  
 152 (Fig. 1(c)). This approach effectively harnesses the strong  
 153 correlation between cross-modal images and the capability  
 154 of style transfer to adapt to different modalities. The con-  
 155 tributions of this work can be summarized in three aspects:

- 156 • We propose a dual-branch CNN-based framework for  
 157 deblurring local blurry target and extracting and fus-  
 158 ing global information, which better reflects the cor-  
 159 respondence between modalities.
- 160 • We propose a style transfer module for the IVIF task to

162 deblur the blurry target and retain visible information.  
 163 • Our method achieves promising image fusion results  
 164 and also performs more superior in downstream tasks  
 165 such as detection and segmentation.  
 166

## 2. Related Work

### 2.1. Infrared-visible fusion

With the development of deep learning, numerous work on IVIF task have emerged [10, 26, 39, 46, 54]. Ma *et al.* proposed a GAN for IVIF task [32], conceptualizing the fusion algorithm as an adversarial game between retaining infrared thermal radiation information and maintaining visible appearance texture information, and achieved substantial breakthroughs. Then, Zhao *et al.* pioneered the exploration of the two-scale decomposition in IVIF task [54], utilizing an encoder to decompose the images into background feature maps and detail feature maps, followed by a decoder used to reconstruct the original image. Following this, Tang *et al.* utilize illumination probability to construct an illumination-aware loss, which guides the training of the fusion network, allowing it to adaptively integrate meaningful information based on lighting conditions. Recently, considering the combination of fusion and downstream pattern recognition tasks, Sun *et al.* and Tang *et al.* proposed the network driven by the downstream task and achieved promising results [37, 39]. Additionally, incorporating a pre-processing registration module before the fusion module has been shown to effectively address the misregistration of source images [10]. Zhao *et al.* introduced a dual-branch Transformer-CNN network to correlate global and local features, achieving a fusion process where low-frequency features are related and high-frequency features are unrelated [51].

### 2.2. Style transfer

Style transfer, initially proposed by Leon *et al.* [4], aims to transfer the artistic style of one image onto another, creating an image with a unique artistic flair. Due to its innovative nature, this technique has attracted significant attention, then numerous style transfer models are implemented and utilized in various field [13, 14], particularly in image restoration and video processing. For image transformation problems, where an input image is converted into an output image, perceptual loss [16] has been designed and utilized for style transfer tasks. Then, Xun *et al.* achieved arbitrary style transfer in real-time by introducing a novel adaptive instance normalization [9]. To tackle the chanllenge of versatile style transfer, Wu *et al.* implemented video style transfer without video in training process [44] through InfoNCE loss [43]. Recently, Kwon *et al.* proposed a network called CLIPstyler, capable of performing style transfer with just a single text condition, achieving results comparable to

216 other models that use more complex inputs [17]. The fundamental principle of classical style transfer methods is to  
 217 generate an image that preserves the content of the original image while seamlessly incorporating the distinctive characteristics of the target style. This ensures that visible images retain more visually detailed texture during the deblurring process.  
 218  
 219  
 220  
 221  
 222  
 223

### 224 3. Method

225 The DSTFuse mainly consists of three modules, which are detailed in Fig. 2. In the cross-modality style transfer  
 226 module (CST-module), the original visible image is combined with the edge information generated by the Sobel algorithm as input. This concatenated input is then fed into the Auto-Encoder-like module to generate a structure-clear image that is similar in style to an infrared image. Finally, the infrared style image re-enters the CST-module as input to generate a new infrared-like image with more visible features. In the fusion module, the pipeline aims to train a Auto-Encoder-based structure for extracting features and reconstructing original images (in reconstruction stage) or generating fusion images initially (in fusion stage). In the mapping module, the output images from the fusion module and the CST-module are merged through an attention block to generate the final output image. The detailed workflow is illustrated in Fig. 2.  
 227  
 228  
 229  
 230  
 231  
 232  
 233  
 234  
 235  
 236  
 237  
 238  
 239  
 240  
 241  
 242  
 243

#### 244 3.1. Cross-modality style transfer module

245 The CST-module aims to retain the visual effect of the visible image while minimizing motion blur as much as possible. To achieve this, the CST-module divides the training  
 246 into two stages, focusing more on the infrared information in the first stage and the visible information in the second. In order to deblur efficiently, it adds feature constraints similar to style transfer to guide training of model and incorporated edge information  $\mathcal{D}_S$  obtained from the Sobel algorithm  $\mathcal{S}$  in both stages:  
 247  
 248  
 249  
 250  
 251  
 252  
 253  
 254

$$\mathcal{D}_S = \mathcal{S}(I) \oplus V, \quad (1)$$

255 where  $\mathcal{S}(I)$  means the result of the Sobel algorithm on  
 256 the infrared image, which only retains the structure of the objects.  $\oplus$  means element-wise addition.  
 257  
 258  
 259  
 260

261 **Stage-1.** Considering the focus of the first stage is the information of infrared image, the input of the first EBlock in  
 262 encoder is designed as the concatenation of visible image  $V$  and edge information  $\mathcal{S}$  to obtain more structural information. In addition, the edge information  $\mathcal{S}$  is extracted as shallow features  $\phi_S$  through a convolution. Then, the  $\phi_S$  is used as the input, together with the second-to-last skip connection, into the final layer of the decoder.  
 263  
 264  
 265  
 266  
 267  
 268  
 269

270 **Stage-2.** After obtaining the image with more functional  
 271 highlight and less motion blur, the output of the first stage  
 272 serves as the input for the second stage. Different from the  
 273 first stage,  $\mathcal{S}$  is no longer used as an input so that the final  
 274 output image will not contain highlighted edges.  
 275

276 The CST-module eliminate the motion blur of the target by introducing the style of infrared images. At the  
 277 same time, the output image should not retain an excessive  
 278 amount of visual information from the input image. Therefore,  
 279 the perceptual loss [16] perfectly meets the requirements. The CST loss is:  
 280  
 281

$$L_{CST} = \alpha_1 L_{perceptual}(D, I, i) + \alpha_2 L_{SSIM}(F, V), \quad (2)$$

282 where  $L_{perceptual}(D, I, i) = \|\phi(D, i) - \phi(I, i)\|_2$ , and  $D$   
 283 is the output of CST-module,  $\phi(\cdot, i)$  is the first  $i$  layers of a  
 284 simple model extractor similar to VGG. As the  $i$  increases,  
 285 the features become more abstract and the style becomes  
 286 more biased towards the infrared image. In first stage of  
 287 CST-module, it's need to retaining the structure of targets  
 288 and reduce blurriness, while in the second stage, the focus  
 289 is on retaining color, texture. Therefore, the layer of model  
 290 extractor in the first stage is more than the second stage.  
 291  
 292  
 293

#### 294 3.2. Fusion module

295 **Reconstruction stage.** The key to make Auto-Encoder  
 296 perform better in image fusion is to extract the most repre-  
 297 sentative features from source images. Capturing accurately  
 298 feature that precisely reflects the advantage of visible and  
 299 infrared images poses a significant challenge. And directly  
 300 extract such features using a randomly initialized encoder  
 301 instead of a well-pretrained one is not feasible.  
 302

303 To address this issue, the reconstruction stage is sched-  
 304 uled before the fusion stage. In this stage, a encoder is  
 305 trained to extract features and a decoder to reconstruct them  
 306 into original images for the subsequent fusion stage. Specif-  
 307 ically, for the input image, it will pass through the encoder  
 308 containing three EBlocks and the decoder with two DBlock  
 309 and one OutBlock to reconstruct itself. The block struct  
 310 can be seen in Fig. 2. And for each block, the residual-  
 311 connection is used to accelerate convergence. In addition,  
 312 skip connections between the first and last layers, and be-  
 313 tween the second and second-to-last layers, prevent gradient  
 314 vanishing.  
 315

316 Since the aim of the reconstruction stage is to minimize  
 317 the information loss of source image, the loss of reconstruc-  
 318 tion can be defined as:  
 319

$$L_{reconstruct} = \alpha_1 f(I, \hat{I}) + \alpha_2 f(V, \hat{V}), \quad (3)$$

320 where  $I$  and  $\hat{I}$ ,  $V$  and  $\hat{V}$  represent the input and output of  
 321 infrared and visible images, respectively. And  
 322

$$f(X, \hat{X}) = \|X - \hat{X}\|_2 + \lambda L_{SSIM}(X, \hat{X}), \quad (4)$$

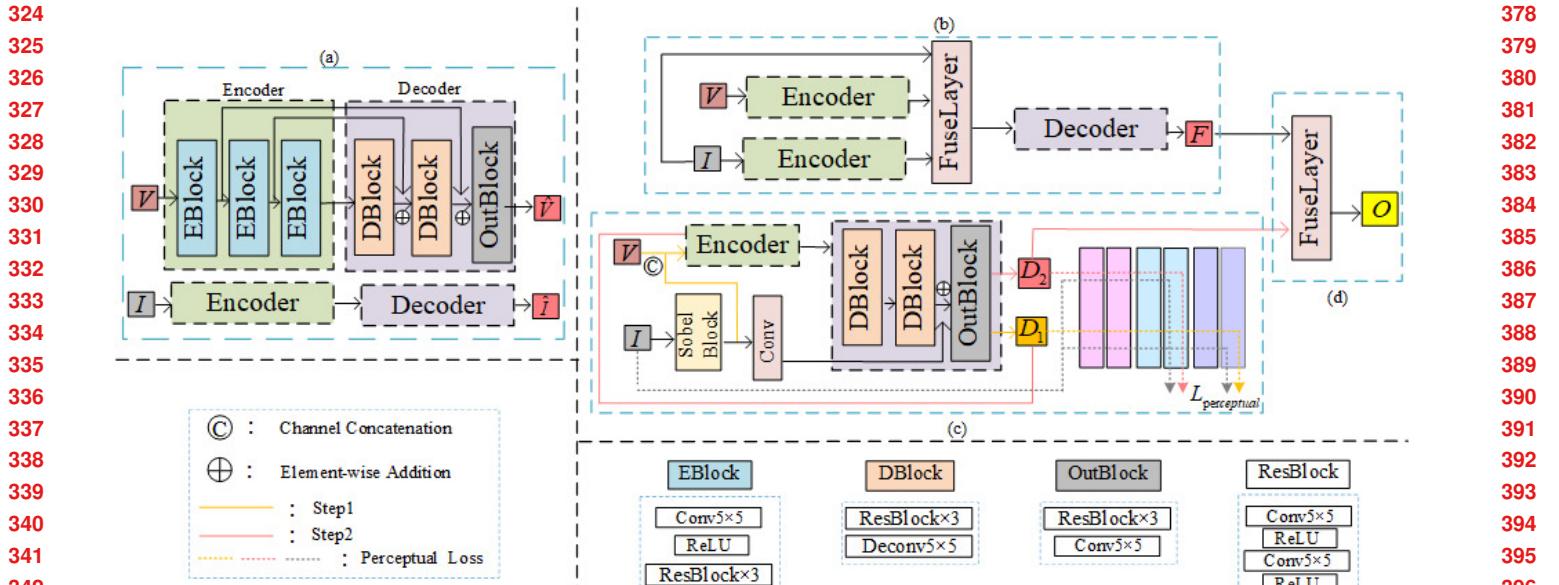


Figure 2. The architecture of DSTFuse, (a) The reconstruction stage of fusion module. (b) The fusion stage of fusion module. (c) The cross-modality style transfer module. (d) The mapping module.

where  $X$  and  $\hat{X}$  represent the above input and output image, and  $L_{SSIM}(X, \hat{X}) = 1 - SSIM(X, \hat{X})$ . SSIM is the structural similarity index, which is a measure of the similarity between two pictures.

**Fusion stage.** After the reconstruction stage, the well-trained feature extractor  $\mathcal{E}(\cdot)$  can be obtained. And the feature  $\{\phi_V, \phi_I\}$  can be extracted from visible and infrared input  $\{V, I\}$  by:

$$\phi_V = \mathcal{E}(V), \quad \phi_I = \mathcal{E}(I). \quad (5)$$

In previous studies, the neglect of suboptimal data has resulted in poor performance on datasets containing blurry images. In contrast to these work, the fusion module in DSTFuse is designed to prioritize the incorporation of detailed background information into the fused image, while deliberately disregarding the target object, which is instead the central focus of the CST module. Considering the high correlation between source visible and infrared images, it can be assumed that objects which appear motion-blurred in the visible image correspond to high-contrast and distinct targets in the infrared image. Therefore, the decoder should be prompted to learn the environmental information excluding the high-contrast targets. To this end, a fusion layer with attention block is added to highlight the background. And the mapping function is described as follow:

$$\phi_A = (\phi_V \oplus \phi_I) \oplus (\phi_V \oplus \phi_I) \otimes (1 - \mathcal{A}(I)), \quad (6)$$

where  $\phi_V$  and  $\phi_I$  are the features extracted from visible and infrared input, respectively.  $\oplus$  and  $\otimes$  means element-wise addition and element-wise multiplication.  $\mathcal{A}(\cdot)$  is attention map matrix.

Finally, the output image  $F$  will preserve more detailed textures which is constrained by the Sobel algorithm and the gradient information. Additionally, the output should be similar to the visible image, so the loss function is:

$$L_{\text{fuse}} = \alpha_1 \text{Sobel}(F, V, I) + \alpha_2 \| F - \max(V, I) \|_1 + \alpha_3 L_{SSIM}(F, V), \quad (7)$$

$$\text{Sobel}(F, V, I) = \text{Sobel}(F) - \max(\text{Sobel}(V), \text{Sobel}(I)), \quad (8)$$

where  $\text{Sobel}(\cdot)$  is the Sobel algorithm.

In addition, the fusion module and the CST-module can be trained simultaneously.

### 3.3. Mapping module

After training through the fusion module and the CST-module, it is possible to obtain a fused image with detailed environmental information and a small amount of functional highlights, as well as a infrared style image with a clear target structure and no motion blur. To integrate the benefits of both images into a final composite, the mapping module generates an attention map matrix derived from the infrared input. This matrix emphasizes the edges of all targets present in the scene. The mapping function is:

$$O = (D_2 \oplus V) \otimes \mathcal{A}(I) \oplus (F \oplus V) \otimes (1 - \mathcal{A}(I)), \quad (9)$$

432 where  $D_2$  and  $F$  are the outputs of CST-module and fusion  
 433 module,  $V$  and  $I$  is the input of visible and infrared image,  
 434 respectively.  $\oplus$  and  $\otimes$  means element-wise addition and  
 435 element-wise multiplication.  $\mathcal{A}(\cdot)$  is the attention block.  
 436

437 After mapping, blurry parts of the final image are com-  
 438 posed of the deblurred image, while the rest is composed of  
 439 the fused image. The attention loss prompts the mapping  
 440 matrix to focus only on the edges of the image, similar to  
 441 fusion loss, which should be constrained by gradient and  
 442 edge information:

$$443 L_{map} = \alpha_1 \text{Sobel}(F, V, I) + \alpha_2 \| F - \max(V, I) \|_1. \quad (10)$$

## 4. Experiment

### 4.1. Settings

444 **Dataset and metrics.** To verify the performance of model  
 445 on deblurring, we select images with motion blur from the  
 446 LLVIP dataset [11] as training set (317 pairs) and test set  
 447 (60 pairs).

448 There are eight metrics used to quantitatively measure  
 449 the fusion results: spatial frequency (SF), average gradi-  
 450 ent (AG), mean square error (MSE), peak signal to noise  
 451 ratio (PSNR), mutual information (MI), visual informa-  
 452 tion fidelity (VIF), correlation coefficient (CC), and struc-  
 453 tural similarity index measure (SSIM). The details of these  
 454 metrics can be found in [29].

455 **Implement details.** DSTFuse is trained by Pytorch on  
 456 single NVIDIA GeForce RTX 3090 GPU and Intel Xeon  
 457 Gold 6330 CPU. The training samples are converted to  
 458 grayscale images and resized to  $640 \times 640$  in the prepro-  
 459 cessing stage. In the training process, the Adam optimizer  
 460 is employed, initializing the learning rate at  $10^{-4}$ . The total  
 461 number of training epochs is set to 15. During the first  
 462 ten epochs, both the fusion module and the CST-module  
 463 undergo concurrent training, with each of the reconstruc-  
 464 tion and fusion stages receiving training for precisely three  
 465 epochs. In the final five epochs, the training is solely di-  
 466 rected at the mapping module. For the tuning parameters in  
 467 loss function, in Eq. (2),  $\alpha_1$  and  $\alpha_2$  are set to 100 and 1. In  
 468 Eq. (3),  $\alpha_1$ ,  $\alpha_2$  and  $\lambda$  are set to 1, 1 and 5. In Eq. (7),  $\alpha_1$  to  
 469  $\alpha_3$  are set to 10, 5 and 1. In Eq. (10),  $\alpha_1$  and  $\alpha_2$  are set to  
 470 10 and 1.

### 4.2. Comparison with SOTA methods

471 In this section, DSTFuse is tested on the test set and  
 472 compare the fusion results with the state-of-the-art meth-  
 473 ods including DIDFuse [54], RFN-Nest [21], MFEIF [26],  
 474 ReCoNet [10], SeAFusion [39], DeFusion [24], MetaFu-  
 475 sion [50], LLRNet [22], EMMA [52].

476 **Qualitative comparison.** It has been shown the qualita-  
 477 tive comparison in Fig. 3. Obviously, the proposed method  
 478 more effectively integrates thermal radiation information  
 479 from infrared images with detailed textures from visible  
 480 images. As show in visual comparison result, the back-  
 481 ground information that was easily overlooked in previous  
 482 methods due to the prominence of infrared images is per-  
 483 fectly retained in DSTFuse. This can be attributed to the  
 484 CST-module, which does not forcibly merge visible images  
 485 with infrared image , but rather performs only style conver-  
 486 sion, thereby preserving most of the visible details. Conse-  
 487 quently, for the objects in dark regions, DSTFuse appropri-  
 488 ately highlights them for identification in downstream task.  
 489 For blurry object, DSTFuse providing details that conform  
 490 to human visual perception.

491 **Quantitative comparison.** Afterward, we follow the pre-  
 492 vious IVIF works by reporting eight metrics for visual eval-  
 493 uation criterion. There are excellent performance across  
 494 most metrics, demonstrating that it is suitable for the hu-  
 495 man visual perception without bias from observers or inter-  
 496 preters. Specifically, the optimal results on MI and CC [29]  
 497 show that the fused image contain the most amount of infor-  
 498 mation and the strongest correlation between source images  
 499 and fused image, respectively. Besides, the promising result  
 500 on SF, AG, MSE, PSNR and VIF [29] indicates show that  
 501 the proposed fusion method produces the most texture de-  
 502 tails, least distortion and best matches to the human visual  
 503 system.

504 **Visualization of CST-module.** Fig. 4 visualizes the effec-  
 505 tiveness of perceptual loss in CST-module. Obviously, with  
 506 training goes on, more detail texture of target are activated  
 507 and more background information are inactivate. As the in-  
 508 put of CST-module, the visible image contains the abundant  
 509 details of the target and exhibit significant perceptual differ-  
 510 ences compared to the infrared images which is regarded as  
 511 style image in style transfer. In the group of CST output, the  
 512 CST-module firstly focus on the profile of target, showing  
 513 that the deblurring function works well. As the perceptual  
 514 loss reaches convergence, an increasing amount of detail is  
 515 incorporated.

### 4.3. Ablation studies

516 The ablation studies are conducted on the LLVIP  
 517 dataset [11] to prove the rationality of DSTFuse, with the  
 518 results shown in Tab. 2 and Fig. 5.

519 **Essential module in DSTFuse.** To independently vali-  
 520 date the efficacy of the fusion module and the mapping  
 521 module, two comparative experiments have been devised.  
 522 In Exp. I, the mapping module is removed to ascertain its  
 523

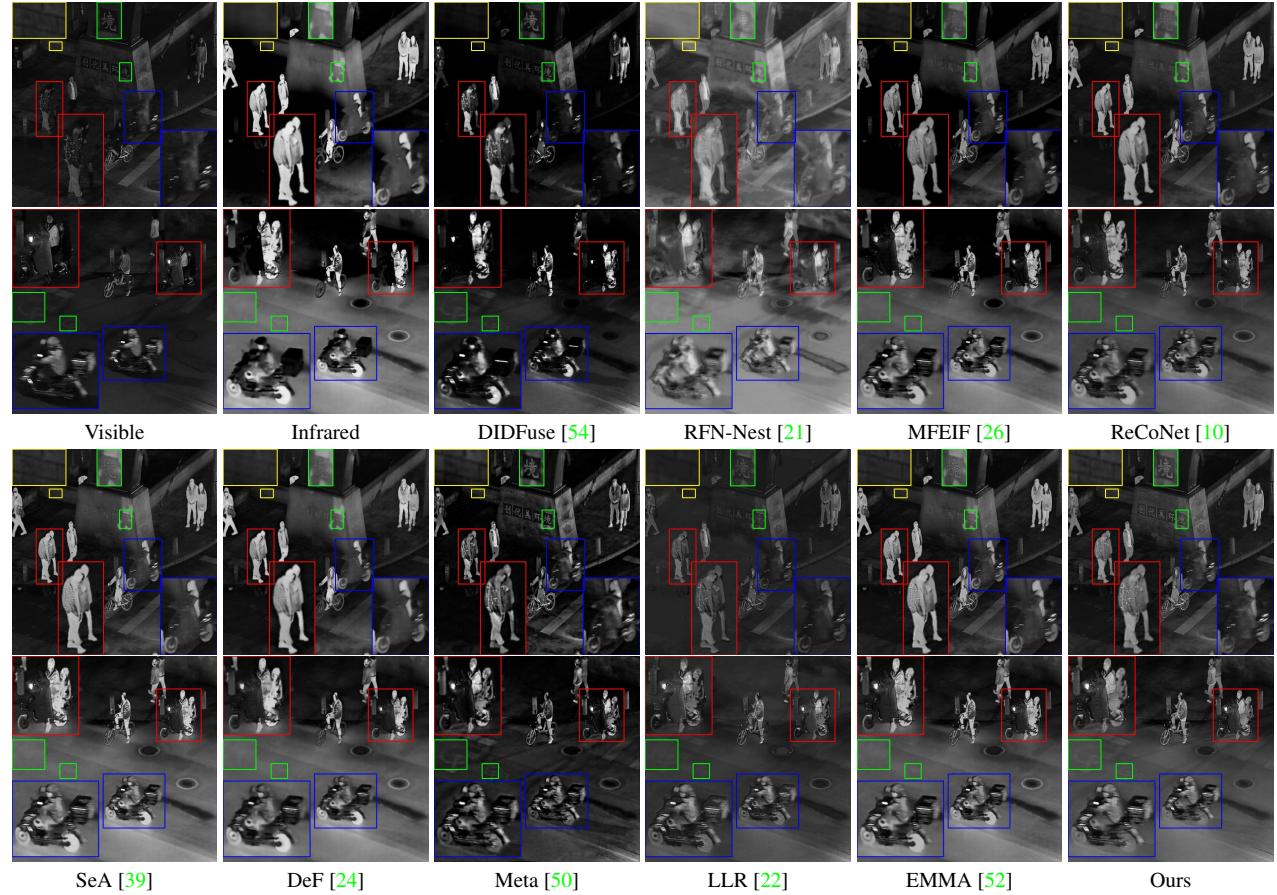


Figure 3. Visual comparison for “010018” (up) and “050131” (down) in LLVIP IVIF dataset.

	DID [54]	RFN [21]	MFE [26]	ReC [10]	SeA [39]	DeF [24]	Meta [50]	LLR [22]	EMM [52]	Ours	
SF	4.943	4.818	4.567	4.712	5.440	4.808	<b>6.137</b>	4.528	5.469	5.582	625
AG	2.036	2.251	1.882	2.031	2.432	2.148	<b>2.985</b>	1.787	2.430	2.540	626
MSE	0.043	0.060	<b>0027</b>	0.030	0.037	0.036	0.037	0.031	0.030	0.029	627
PSNR	13.71	12.22	<b>15.78</b>	14.26	14.59	14.40	14.46	14.83	15.42	<b>15.56</b>	628
MI	1.581	<u>1.601</u>	1.181	1.357	1.506	1.286	1.214	1.494	1.595	<b>1.650</b>	629
VIF	0.746	0.753	0.814	0.813	<b>0.934</b>	0.850	0.898	0.593	0.903	0.919	630
CC	0.686	0.674	<u>0.712</u>	0.707	0.696	0.647	0.684	0.693	0.703	<b>0.734</b>	631
SSIM	1.044	0.999	1.298	<b>1.398</b>	1.358	<u>1.367</u>	1.206	1.304	1.306	1.316	632

Table 1. Quantitative results of the IVIF task. The **Bold** and underline show the best, second-best value, respectively.

capability in accurately mapping cross-modal information. As an alternative, the summation method is used to integrate output of CST-module with that of the fusion module. In formula, the summation method can be described as:

$$O = (D_2 \oplus F), \quad (11)$$

where  $D_2$  and  $F$  are the outputs of CST-module and fusion module, respectively.

When removing the mapping module, although the net-

work retains the ability to execute feature mapping, it falls short in precisely selecting the requisite information from distinct images. In Exp. II, the fusion module is eliminated to confirm the proficiency in extracting background information. To substitute for the output of this module, the original visible image is utilized to provide background information. Results in Exp. II illustrate that the absence of effective feature extraction results in a lack of detail and texture, particularly in darker regions, thereby causing a de-

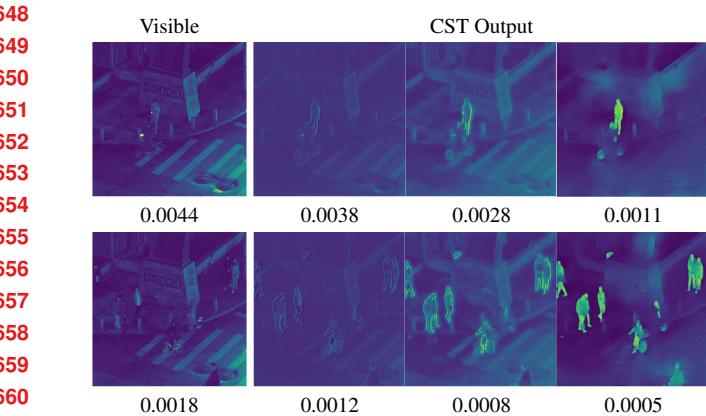


Figure 4. Visualization of the CST-module for “010018” (up), “010054” (down) in LLVIP IVIF dataset. The values represent the results of the perceptual loss.

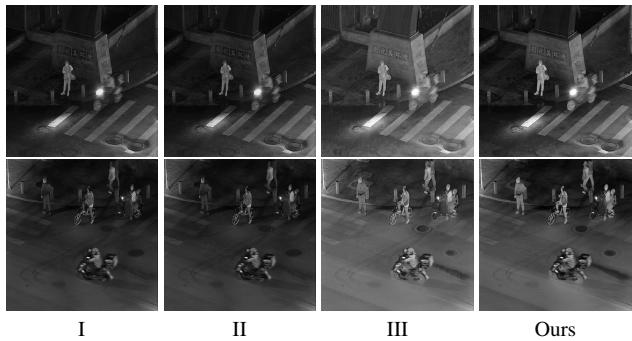


Figure 5. Ablation experiment for “010562” (up) and “050131” (down) in LLVIP IVIF dataset.

cline in overall performance.

**Term in loss function.** Then, in Exp. III, it separately removes the perceptual loss from CST-module and modifies it to adopt the conventional loss function used in other fusion tasks, denoted as  $\mathcal{L}_2 = \|x - \hat{x}\|_2$ . And in the first step of CST-module, the  $x$  represents the infrared image, while in the second step, it represents the visible image. The perceptual loss ensures that, during the style transfer process within the CST-module, the content image adequately inherits information from the style image, thereby making the generated image perceptually more similar to the source image. In contrast, the conventional loss function merely enforces the image to be similar to the source image. Results in Exp. III demonstrate the necessity of perceptual loss.

#### 4.4. Application in the downstream tasks

To evaluate the promoting effect of fused image and its improved performances on downstream task, further exten-

	Configurations	SF	PSNR	CC	VIF	702
I	w/o Mapping Module	5.364	14.36	0.637	0.645	703
II	w/o Fusion Module	4.947	14.24	0.680	0.835	704
III	w/o Perceptual Loss	5.474	14.15	0.730	0.835	705
	Ours	<b>5.582</b>	<b>15.56</b>	<b>0.734</b>	<b>0.919</b>	706
						707

Table 2. Ablation experiment results. **Bold** indicates the best value.

nal validation is conducted. For infrared-visible object detection, the fused images generated by state-of-the-art models are evaluated using five classic detectors by comparing the AP value for person detection. The selected detectors include Faster R-CNN [5], YOLOv5 [36], SSD [27], RetinaNet [25] and Mask R-CNN [6]. For infrared-visible semantic segmentation, the segmentation network includes FCN [28], DeeplabV3 [2] and LSR-APP [7]. And the performance is evaluated using Intersection over Union (IoU) for person segmentation.

**Object detection.** As shown in Tab. 3, DSTFuse plays a significantly positive role in detection. In comparison to direct predictions made on the source images, the fused images generated by DSTFuse substantially enhance prediction accuracy across all five detection models. Compared to previous work, DSTFuse exhibits the promising superior detection capabilities, which can be contributed to its ability to preserve information that aligns closely with the human visual system. To obtain a more intuitive comparison, the detection results are compared using YOLOv5 [36] as the detector and the visual results are shown in Fig. 6. In the first example, when the infrared source images have already been effectively detected, only the fused image generated by DSTFuse can retain the infrared features and be detected. And for those infrared images that perform poorly in detection due to overly prominent functional highlights (e.g., the second example in Fig. 6), the fused images generated by DSTFuse appropriately balance the high contrast of the targets with the real pixel intensity. This allows the detector to accurately identify each target.

**Semantic segmentation.** To evaluate the performance of DSTFuse on infrared-visible semantic segmentation, we selected 42 pairs of infrared and visible images from the LLVIP dataset [11], and proceeded to annotate the person category within these images. The result in Tab. 4 show that DSTFuse effectively integrates the contour details from the source images, thereby enhancing the model’s ability to recognize object boundaries and achieving more accurate segmentation.

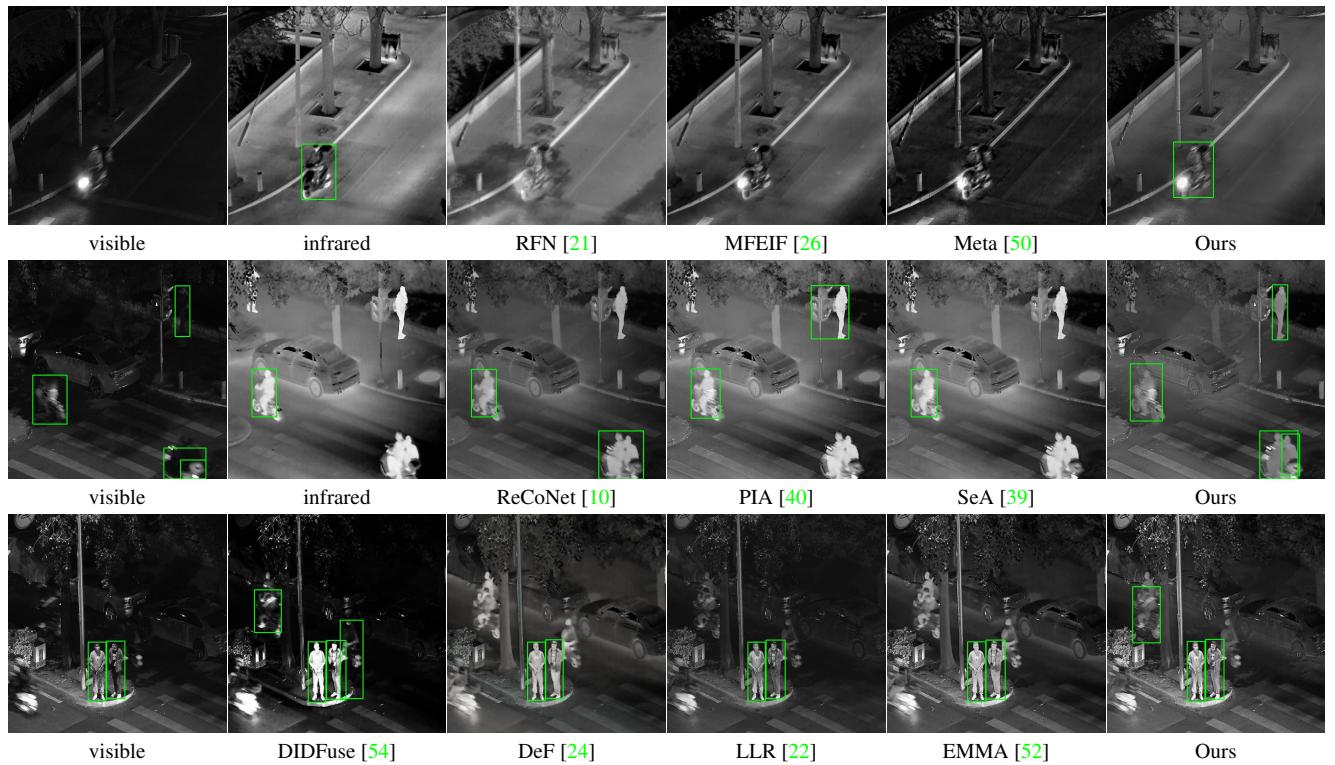


Figure 6. Detection results on source images and different fused image for “040145” (up), “080246” (middle) and “080786” (down) in LLVIP IVIF dataset.

	Faster RCNN [5]	YOLOv5 [36]	SSD [27]	Retina [25]	Mask RCNN [6]
VI	30.52	39.26	40.40	43.08	55.20
IR	34.03	37.80	39.29	41.28	48.97
DID [54]	32.05	39.11	44.97	43.56	54.38
RFN [21]	15.15	27.65	16.52	27.63	33.88
MFE [26]	32.48	42.57	39.67	41.00	50.94
ReC [10]	<b>37.59</b>	46.10	<u>46.65</u>	42.23	54.68
SeA [39]	36.25	44.54	46.13	46.84	55.24
DeF [24]	<u>37.57</u>	44.31	43.45	46.60	52.82
Meta [50]	36.43	43.49	<b>49.36</b>	<b>52.19</b>	56.42
LLR [22]	33.55	<b>53.37</b>	45.02	45.90	52.83
EMMA [52]	34.34	45.91	45.02	43.67	53.16
Ours	36.58	<u>46.87</u>	42.85	<u>48.86</u>	<b>57.24</b>

Table 3. AP(%) values of person for detection on LLVIP dataset. The **bold** and underline show the best and second-best value, respectively.

## 5. Conclusion

This paper presents a infrared-visible fusion framework through introducing the style transfer. With the cross-modality style transfer module, target with motion blur in visible image are more clearly outlined and more easily recognized. Experiments demonstrate the fusion effect of

	FCN [28]	DeeplabV3 [2]	LSR-APP [7]
DID [54]	47.91	48.12	48.07
RFN [21]	44.69	46.49	47.99
MFE [26]	48.23	48.51	48.37
ReC [10]	48.32	<b>48.70</b>	<b>48.57</b>
SeA [39]	48.34	48.63	<u>48.53</u>
DeF [24]	<u>48.37</u>	48.52	48.52
Meta [50]	47.99	48.32	48.26
LLR [22]	47.97	48.29	48.13
EMMA [52]	48.10	48.37	48.41
Ours	<b>48.48</b>	<u>48.64</u>	48.48

Table 4. IoU(%) values of person for semantic segmentation on LLVIP dataset. The **bold** and underline show the best and second-best value, respectively.

DSTFuse, and the performance on downstream detection and segmentation can be also improved.

## References

- [1] Njuod Alsudays, Jing Wu, Yu-Kun Lai, and Ze Ji. Afpsnet: Multi-class part parsing based on scaled attention and feature fusion. In WACV, pages 4033–4042, 2023. 1

- 864 [2] Liang-Chieh Chen, George Papandreou, Florian Schroff, and  
865 Hartwig Adam. Rethinking atrous convolution for semantic  
866 image segmentation. *arXiv preprint arXiv:1706.05587*,  
867 2017. 7, 8
- 868 [3] Yu Fu, Xiao-Jun Wu, and Tariq Durrani. Image fusion based  
869 on generative adversarial network consistent with perception.  
870 *Information Fusion*, 72:110–125, 2021. 1
- 871 [4] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge.  
872 Image style transfer using convolutional neural networks. In  
873 *CVPR*, June 2016. 2
- 874 [5] Ross Girshick. Fast r-cnn. In *ICCV*, pages 1440–1448, 2015.  
875 7, 8
- 876 [6] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick.  
877 Mask r-cnn. In *ICCV*, pages 2961–2969, 2017. 7,  
878 8
- 879 [7] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh  
880 Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu,  
881 Ruoming Pang, Vijay Vasudevan, Quoc V. Le, and Hartwig  
882 Adam. Searching for mobilenetv3. In *ICCV*, October 2019.  
883 7, 8
- 884 [8] Hai-Miao Hu, Jiawei Wu, Bo Li, Qiang Guo, and Jin Zheng.  
885 An adaptive fusion algorithm for visible and infrared videos  
886 based on entropy and the cumulative distribution of gray levels.  
887 *IEEE T MULTIMEDIA*, 19(12):2706–2719, 2017. 1
- 888 [9] Xun Huang and Serge Belongie. Arbitrary style transfer in  
889 real-time with adaptive instance normalization. In *ICCV*,  
890 2017. 2
- 891 [10] Zhanbo Huang, Jinyuan Liu, Xin Fan, Risheng Liu, Wei  
892 Zhong, and Zhongxuan Luo. Reconet: Recurrent correction  
893 network for fast and efficient multi-modality image fusion.  
894 In *ECCV*, pages 539–555. Springer, 2022. 1, 2, 5, 6, 8
- 895 [11] Xinyu Jia, Chuang Zhu, Minzhen Li, Wenqi Tang, and Wenli  
896 Zhou. Llivip: A visible-infrared paired dataset for low-light  
897 vision. In *ICCV*, pages 3496–3504, 2021. 5, 7
- 898 [12] Lihua Jian, Xiaomin Yang, Zheng Liu, Gwanggil Jeon, Min-  
899 gliang Gao, and David Chisholm. Sedrfuse: A symmetric  
900 encoder-decoder with residual block network for infrared  
901 and visible image fusion. *IEEE T INSTRUM MEAS*, 70:1–  
902 15, 2020. 1
- 903 [13] Yongcheng Jing, Xiao Liu, Yukang Ding, Xinchao Wang,  
904 Errui Ding, Mingli Song, and Shilei Wen. Dynamic instance  
905 normalization for arbitrary style transfer. In *AAAI*, 2020. 2
- 906 [14] Yongcheng Jing, Yang Liu, Yezhou Yang, Zunlei Feng,  
907 Yizhou Yu, Dacheng Tao, and Mingli Song. Stroke con-  
908 trollable fast style transfer with adaptive receptive fields. In  
909 *ECCV*, 2018. 2
- 910 [15] Jing jing Zong and Tian shuang Qiu. Medical image fusion  
911 based on sparse representation of classified image patches.  
912 *Biomedical Signal Processing and Control*, 34:195–205,  
913 2017. 1
- 914 [16] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual  
915 losses for real-time style transfer and super-resolution. In  
916 *ECCV*, pages 694–711. Springer, 2016. 2, 3
- 917 [17] Gihyun Kwon and Jong Chul Ye. Clipstyler: Image style  
918 transfer with a single text condition. In *CVPR*, pages 18062–  
919 18071, 2022. 3
- 920 [18] Fayed Lahoud and Sabine Susstrunk. Ar in vr: Simulating  
921 infrared augmented vision. In *ICIP*, pages 3893–3897. IEEE,  
922 2018. 1
- 923 [19] Hui Li and Xiao-Jun Wu. Densefuse: A fusion approach to  
924 infrared and visible images. *IEEE TIP*, 28(5):2614–2623,  
925 May 2019. 1
- 926 [20] Hui Li, Xiao-Jun Wu, and Tariq Durrani. NestFuse: An In-  
927frared and Visible Image Fusion Architecture based on Nest  
928 Connection and Spatial/Channel Attention Models. *IEEE T  
929 INSTRUM MEAS*, 69(12):9645–9656, 2020. 1
- 930 [21] Hui Li, Xiao-Jun Wu, and Josef Kittler. Rfn-nest: An end-to-  
931 end residual fusion network for infrared and visible images.  
932 *Information Fusion*, 73:72–86, March 2021. 5, 6, 8
- 933 [22] Hui Li, Tianyang Xu, Xiao-Jun Wu, Jiwen Lu, and Josef Kit-  
934 tler. LRRNet: A novel representation learning guided fusion  
935 framework for infrared and visible images. *IEEE TPAMI*,  
936 45(9):11040–11052, 2023. 5, 6, 8
- 937 [23] Shutao Li, Bin Yang, and Jianwen Hu. Performance com-  
938 parison of different multi-resolution transforms for image fu-  
939 sion. *Information Fusion*, 12(2):74–84, 2011. 1
- 940 [24] Pengwei Liang, Junjun Jiang, Xianming Liu, and Jiayi Ma.  
941 Fusion from decomposition: A self-supervised decomposi-  
942 tion approach for image fusion. In *ECCV*, 2022. 5, 6, 8
- 943 [25] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and  
944 Piotr Dollár. Focal loss for dense object detection. In *ICCV*,  
945 pages 2980–2988, 2017. 7, 8
- 946 [26] Jinyuan Liu, Xin Fan, Ji Jiang, Risheng Liu, and Zhongx-  
947 uan Luo. Learning a deep multi-scale feature ensemble  
948 and an edge-attention guidance for image fusion. *TCSVT*,  
949 32(1):105–119, 2021. 1, 2, 5, 6, 8
- 950 [27] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian  
951 Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C  
952 Berg. Ssd: Single shot multibox detector. In *ECCV*, pages  
953 21–37. Springer, 2016. 7, 8
- 954 [28] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully  
955 convolutional networks for semantic segmentation. In  
956 *CVPR*, pages 3431–3440, 2015. 7, 8
- 957 [29] Jiayi Ma, Yong Ma, and Chang Li. Infrared and visible im-  
958 age fusion methods and applications: A survey. *Information  
959 Fusion*, 45:153–178, 2019. 5
- 960 [30] Jiayi Ma, Linfeng Tang, Fan Fan, Jun Huang, Xiaoguang  
961 Mei, and Yong Ma. Swinfusion: Cross-domain long-range  
962 learning for general image fusion via swin transformer. *JAS*,  
963 9(7):1200–1217, 2022. 1
- 964 [31] Jiayi Ma, Han Xu, Junjun Jiang, Xiaoguang Mei, and Xiao-  
965 Ping Zhang. Ddcgan: A dual-discriminator conditional gen-  
966 erative adversarial network for multi-resolution image fu-  
967 sion. *IEEE TIP*, 29:4980–4995, 2020. 1
- 968 [32] Jiayi Ma, Wei Yu, Pengwei Liang, Chang Li, and Junjun  
969 Jiang. Fusiongan: A generative adversarial network for in-  
970frared and visible image fusion. *Information Fusion*, 48:11–  
971 26, 2019. 1, 2
- 972 [33] Bikash Meher, Sanjay Agrawal, Rutuparna Panda, and Ajith  
973 Abraham. A survey on region based image fusion methods.  
974 *Information Fusion*, 48:119–132, 2019. 1
- 975 [34] Lukas Mehl, Azin Jahedi, Jenny Schmalfuss, and Andrés  
976 Bruhn. M-fuse: Multi-frame fusion for scene flow estima-  
977 tion. In *WACV* 2020–2029, 2023. 1

- 972 [35] Ujwala Patil and Uma Mudengudi. Image fusion using hierarchical pca. In *ICIP*, pages 1–6. IEEE, 2011. 1 1026
- 973 [36] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, pages 779–788, 2016. 7, 8 1027
- 974 [37] Yiming Sun, Bing Cao, Pengfei Zhu, and Qinghua Hu. Detfusion: A detection-driven infrared and visible image fusion network. In *ACM MM*, pages 4003–4011, 2022. 2 1028
- 975 [38] Linfeng Tang, Yuxin Deng, Yong Ma, Jun Huang, and Jiayi Ma. Superfusion: A versatile image registration and fusion network with semantic awareness. *JAS*, 9(12):2121–2137, 2022. 1 1029
- 976 [39] Linfeng Tang, Jiteng Yuan, and Jiayi Ma. Image fusion in the loop of high-level vision tasks: A semantic-aware real-time infrared and visible image fusion network. *Information Fusion*, 82:28–42, 2022. 1, 2, 5, 6, 8 1030
- 977 [40] Linfeng Tang, Jiteng Yuan, Hao Zhang, Xingyu Jiang, and Jiayi Ma. Piafusion: A progressive infrared and visible image fusion network based on illumination aware. *Information Fusion*, 83-84:79–92, 2022. 2, 8 1031
- 978 [41] Linfeng Tang, Hao Zhang, Han Xu, and Jiayi Ma. Deep learning-based image fusion: A survey. *Journal of Image and Graphics*, 28(1):3–36, 2023. 1 1032
- 979 [42] Alexander Toet and Maarten A. Hogervorst. Progress in color night vision. *Optical Engineering*, 51:010901 – 010901, 2012. 2 1033
- 980 [43] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *CoRR*, abs/1807.03748, 2018. 2 1034
- 981 [44] Zijie Wu, Zhen Zhu, Junping Du, and Xiang Bai. Ccpl: Contrastive coherence preserving loss for versatile style transfer. In *ECCV*, pages 189–206. Springer, 2022. 2 1035
- 982 [45] Han Xu, Pengwei Liang, Wei Yu, Junjun Jiang, and Jiayi Ma. Learning a generative model for fusing infrared and visible images via conditional generative adversarial network with dual discriminators. In *IJCAI*, pages 3954–3960, 2019. 1 1036
- 983 [46] Han Xu, Xinya Wang, and Jiayi Ma. Drf: Disentangled representation for visible and infrared image fusion. *IEEE T INSTRUM MEAS*, 70:1–13, 2021. 1, 2 1037
- 984 [47] Han Xu, Hao Zhang, and Jiayi Ma. Classification saliency-based rule for visible and infrared image fusion. *IEEE TCI*, 7:824–836, 2021. 1 1038
- 985 [48] Mingde Yao, Zhiwei Xiong, Lizhi Wang, Dong Liu, and Xuejin Chen. Spectral-depth imaging with deep learning based reconstruction. *Optics express*, 27(26):38312–38325, 2019. 1 1039
- 986 [49] Hao Zhang, Han Xu, Xin Tian, Junjun Jiang, and Jiayi Ma. Image fusion meets deep learning: A survey and perspective. *Information Fusion*, 76:323–336, 2021. 1 1040
- 987 [50] Wenda Zhao, Shigeng Xie, Fan Zhao, You He, and Huchuan Lu. Metafusion: Infrared and visible image fusion via meta-feature embedding from object detection. In *CVPR*, pages 13955–13965, 2023. 5, 6, 8 1041
- 988 [51] Zixiang Zhao, Haowen Bai, Jiangshe Zhang, Yulun Zhang, Shuang Xu, Zudi Lin, Radu Timofte, and Luc Van Gool. Cddfuse: Correlation-driven dual-branch feature decomposition for multi-modality image fusion. In *CVPR*, pages 5906–5916, June 2023. 2 1042
- 989 [52] Zixiang Zhao, Haowen Bai, Jiangshe Zhang, Yulun Zhang, Kai Zhang, Shuang Xu, Dongdong Chen, Radu Timofte, and Luc Van Gool. Equivariant multi-modality image fusion. In *CVPR*, June 2024. 5, 6, 8 1043
- 990 [53] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, and Jiangshe Zhang. Bayesian fusion for infrared and visible images. *Signal Processing*, 177:107734, Dec. 2020. 1 1044
- 991 [54] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1045
- 992 [55] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1046
- 993 [56] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1047
- 994 [57] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1048
- 995 [58] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1049
- 996 [59] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1050
- 997 [60] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1051
- 998 [61] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1052
- 999 [62] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1053
- 1000 [63] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1054
- 1001 [64] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1055
- 1002 [65] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1056
- 1003 [66] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1057
- 1004 [67] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1058
- 1005 [68] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1059
- 1006 [69] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1060
- 1007 [70] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1061
- 1008 [71] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1062
- 1009 [72] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1063
- 1010 [73] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1064
- 1011 [74] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1065
- 1012 [75] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1066
- 1013 [76] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1067
- 1014 [77] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1068
- 1015 [78] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1069
- 1016 [79] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1070
- 1017 [80] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1071
- 1018 [81] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1072
- 1019 [82] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1073
- 1020 [83] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1074
- 1021 [84] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1075
- 1022 [85] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1076
- 1023 [86] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1077
- 1024 [87] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1078
- 1025 [88] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image decomposition for infrared and visible image fusion. In *PRI-CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8 1079