

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 1326

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, 34, 35, 49, 50, 54]. One prevalent application of image fusion is the infrared and visible image fusion (IVIF) [31, 39, 40, 42]. 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

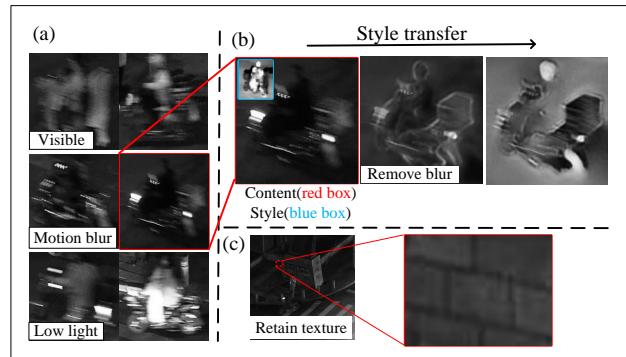


Figure 1. The effect of style transfer on suboptimal visible data for IVIF tasks. (a) The suboptimal visible data with motion blur and low light condition. (b) The blur removing process of transferring visible images into infrared style. The red box and blue box are visible content image and infrared 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 [19], security [9], and medical image processing [50].

To tackle the shortcomings of conventional IVIF methods [16, 24, 36], numerous deep learning-based techniques have been developed. These method can be categorized into two main classes: Generative Adversarial (GANs)-based network [3, 32, 33, 46] and the Auto-Encoder-based network [13, 20, 21, 48]. 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 [11, 27, 40, 47, 55].

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 [43] and MSRS [41], 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 [4, 14, 17, 45] have closely
 134 aligned with 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 deblurring local blurry target and extracting and fusing global information, which better reflects the correspondence between modalities.
- 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 [11, 27, 40]. Ma *et al.* [33] proposed a GAN for IVIF task, 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.* [55] pioneered the exploration of the two-scale decomposition in IVIF task, 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. Recently, considering the combination of fusion and downstream pattern recognition tasks, Sun *et al.* [38] and Tang *et al.* [40] proposed the network driven by the downstream task and achieved promising results. Additionally, incorporating a pre-processing registration module before the fusion module has been shown to effectively address the misregistration of source images [11]. Zhao *et al.* [52] 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.

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 [14, 15], particularly in image restoration and video processing. For image transformation problems, where an input image is converted into an output image, perceptual loss [17] has been designed and utilized for style transfer tasks. Then, Xun *et al.* [10] achieved arbitrary style transfer in real-time by introducing a novel adaptive instance normalization. To tackle the chanllenge of versatile style transfer, Wu *et al.* [45] implemented video style transfer without video in training process through InfoNCE loss [44]. Recently, Kwon *et al.* [18] proposed a network called CLIPstyler, capable of performing style transfer with just a single text condition, achieving results comparable to other models that use more complex inputs. The fundamental principle of classical style transfer methods is to generate an image that preserves the content of the original image while seamlessly incorporating the distinctive charac-

216 teristics of the target style. This ensures that visible images
 217 retain more visually detailed texture during the deblurring
 218 process.
 219

220 3. Method

221 The DSTFuse mainly consists of three modules, which
 222 are detailed in Fig. 2. In the cross-modality style transfer
 223 module (CST-module), the original visible image is com-
 224 bined with the edge information generated by the Sobel al-
 225 gorithm [6] as input. This concatenated input is then fed
 226 into the Auto-Encoder-like module to generate a structure-
 227 clear image that is similar in style to an infrared image. Fi-
 228 nally, the infrared style image re-enters the CST-module as
 229 input to generate a new infrared-like image with more visi-
 230 ble features. In the fusion module, the pipeline aims to train
 231 a Auto-Encoder-based structure for extracting features and
 232 reconstructing original images (in reconstruction stage) or
 233 generating fusion images initially (in fusion stage). In the
 234 mapping module, the output images from the fusion module
 235 and the CST-module are merged through an attention block
 236 to generate the final output image. The detailed workflow is
 237 illustrated in Fig. 2.
 238

239 3.1. Cross-modality style transfer module

240 The CST-module aims to retain the visual effect of the
 241 visible image while minimizing motion blur as much as pos-
 242 sible. To achieve this, the CST-module divides the training
 243 into two stages, focusing more on the infrared information
 244 in the first stage and the visible information in the second.
 245 In order to deblur efficiently, it adds feature constraints sim-
 246 ilar to style transfer to guide training of model and incorpo-
 247 rated edge information \mathcal{D}_S obtained from the Sobel algo-
 248 rithm \mathcal{S} in both stages:
 249

$$250 \quad \mathcal{D}_S = \mathcal{S}(I) \oplus V, \quad (1)$$

251 where $\mathcal{S}(I)$ means the result of the Sobel algorithm [6]
 252 on the infrared image, which only retains the structure of
 253 the objects. \oplus means element-wise addition.
 254

255 **Stage-1.** Considering the focus of the first stage is the in-
 256 formation of infrared image, the input of the first EBlock in
 257 encoder is designed as the concatenation of visible image
 258 V and edge information \mathcal{S} to obtain more structural infor-
 259 mation. In addition, the edge information \mathcal{S} is extracted as
 260 shallow features ϕ_S through a convolution. Then, the ϕ_S
 261 is used as the input, together with the second-to-last skip
 262 connection, into the final layer of the decoder.
 263

264 **Stage-2.** After obtaining the image with more functional
 265 highlight and less motion blur, the output of the first stage
 266 serves as the input for the second stage. Different from the
 267

268 first stage, \mathcal{S} is no longer used as an input so that the final
 269 output image will not contain highlighted edges.
 270

271 The CST-module eliminate the motion blur of the tar-
 272 get by introducing the style of infrared images. At the
 273 same time, the output image should not retain an excessive
 274 amount of visual information from the input image. There-
 275 fore, the perceptual loss [17] perfectly meets the require-
 276 ments. The CST loss is:
 277

$$278 \quad L_{CST} = \alpha_1 L_{perceptual}(D, I, i) + \alpha_2 L_{SSIM}(F, V), \quad (2)$$

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

289 3.2. Fusion module

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

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

311 Since the aim of the reconstruction stage is to minimize
 312 the information loss of source image, the loss of reconstruc-
 313 tion can be defined as:
 314

$$315 \quad L_{reconstruct} = \alpha_1 f(I, \hat{I}) + \alpha_2 f(V, \hat{V}), \quad (3)$$

316 where I and \hat{I} , V and \hat{V} represent the input and output of
 317 infrared and visible images, respectively.
 318

$$319 \quad f(X, \hat{X}) = \|X - \hat{X}\|_2 + \lambda L_{SSIM}(X, \hat{X}), \quad (4)$$

320 where X and \hat{X} represent the above input and output im-
 321 age, and $L_{SSIM}(X, \hat{X}) = 1 - SSIM(X, \hat{X})$. SSIM is the
 322 structural similarity index, which is a measure of the simi-
 323 larity between two pictures.
 324

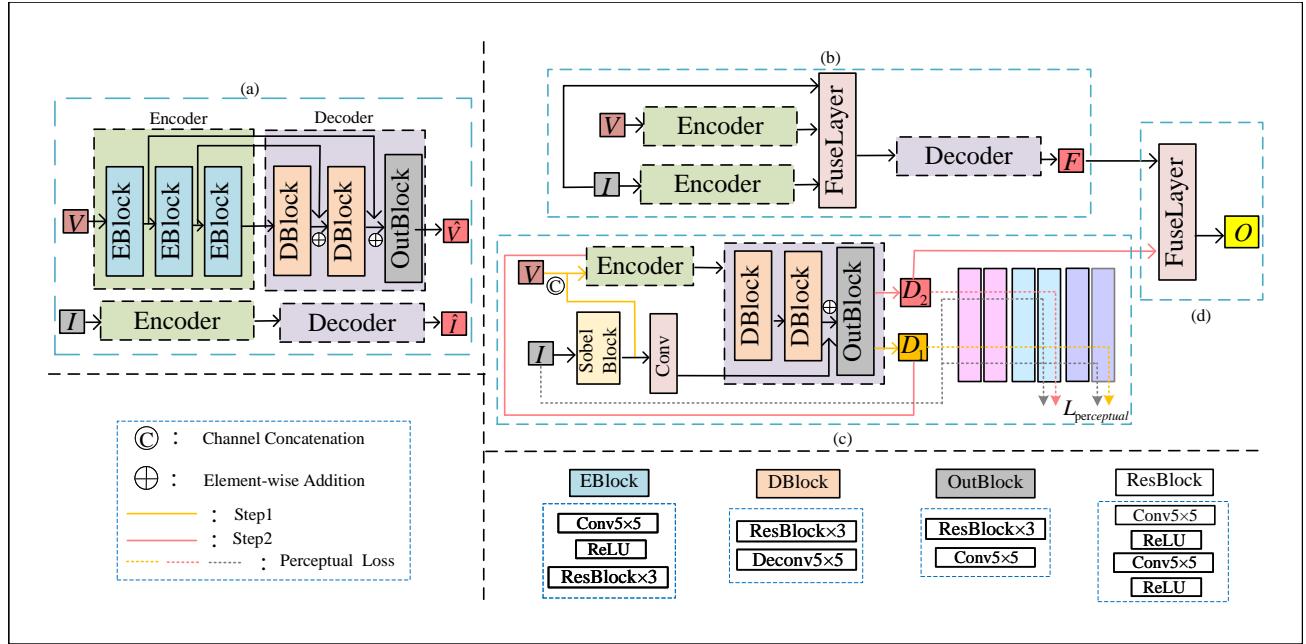


Figure 2. The architecture of DSTFuse, (a) The reconstruction stage of fusion module. (b) The fusion stage of fusion module. (c) The cross-modality style transfer module. (d) The mapping module.

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 ϕ_V and ϕ_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 [6] and the gradient information. Additionally, the output should be similar to the visible image, so the loss function is:

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

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

where $\text{Sobel}(\cdot)$ is the Sobel algorithm [6].

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:

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

where D_2 and F are the outputs of CST-module and fusion module, V and I is the input of visible and infrared image,

432 respectively. \oplus and \otimes means element-wise addition and
 433 element-wise multiplication. $\mathcal{A}(\cdot)$ is the attention block.
 434

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

$$439 L_{map} = \alpha_1 \text{Sobel}(F, V, I) + \alpha_2 \| F - \max(V, I) \|_1. \quad (10)$$

440 4. Experiment

441 4.1. Settings

442 **Dataset and metrics.** To verify the performance of model
 443 on deblurring, we select images with motion blur from the
 444 LLVIP dataset [12] as training set (317 pairs) and test set
 445 (60 pairs).

446 There are eight metrics used to quantitatively measure
 447 the fusion results: spatial frequency (SF), average gradient (AG),
 448 mean square error (MSE), peak signal to noise ratio (PSNR),
 449 mutual information (MI), visual information fidelity (VIF),
 450 correlation coefficient (CC), and structural similarity index measure (SSIM). The details of these met-
 451 riques can be found in [30].

452 **Implement details.** DSTFuse is trained by Pytorch on
 453 single NVIDIA GeForce RTX 3090 GPU and Intel Xeon
 454 Gold 6330 CPU. The training samples are converted to
 455 grayscale images and resized to 640×640 in the prepro-
 456 cessing stage. In the training process, the Adam optimizer
 457 is employed, initializing the learning rate at 10^{-4} . The total
 458 number of training epochs is set to 15. During the first
 459 ten epochs, both the fusion module and the CST-module
 460 undergo concurrent training, with each of the reconstruc-
 461 tion and fusion stages receiving training for precisely three
 462 epochs. In the final five epochs, the training is solely di-
 463 rected at the mapping module. For the tuning parameters in
 464 loss function, in Eq. (2), α_1 and α_2 are set to 100 and 1. In
 465 Eq. (3), α_1 , α_2 and λ are set to 1, 1 and 5. In Eq. (7), α_1 to
 466 α_3 are set to 10, 5 and 1. In Eq. (10), α_1 and α_2 are set to
 467 10 and 1.

468 4.2. Comparison with SOTA methods

469 In this section, DSTFuse is tested on the test set and
 470 compare the fusion results with the state-of-the-art meth-
 471 ods including DIDFuse [55], RFN-Nest [22], MFEIF [27],
 472 ReCoNet [11], SeAFusion [40], DeFusion [25], MetaFu-
 473 sion [51], LLRNet [23], EMMA [53].

474 **Qualitative comparison.** It has been shown the qualita-
 475 tive comparison in Fig. 3. Obviously, the proposed method
 476 more effectively integrates thermal radiation information

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

489 **Quantitative comparison.** Afterward, we follow the pre-
 490 vious IVIF works by reporting eight metrics for visual eval-
 491 uation criterion. There are excellent performance across
 492 most metrics, demonstrating that it is suitable for the hu-
 493 man visual perception without bias from observers or inter-
 494 preters. Specifically, the optimal results on MI and CC [30]
 495 show that the fused image contain the most amount of infor-
 496 mation and the strongest correlation between source images
 497 and fused image, respectively. Besides, the promising result
 498 on SF, AG, MSE, PSNR and VIF [30] indicates show that
 499 the proposed fusion method produces the most texture de-
 500 tails, least distortion and best matches to the human visual
 501 system.

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

514 4.3. Ablation studies

515 The ablation studies are conducted on the LLVIP
 516 dataset [12] to prove the rationality of DSTFuse, with the
 517 results shown in Tab. 2 and Fig. 5.

518 **Essential module in DSTFuse.** To independently vali-
 519 date the efficacy of the fusion module and the mapping
 520 module, two comparative experiments have been devised.
 521 In Exp. I, the mapping module is removed to ascertain its
 522 capability in accurately mapping cross-modal information.
 523 As an alternative, the summation method is used to inte-
 524 grate output of CST-modlue with that of the fusion module.

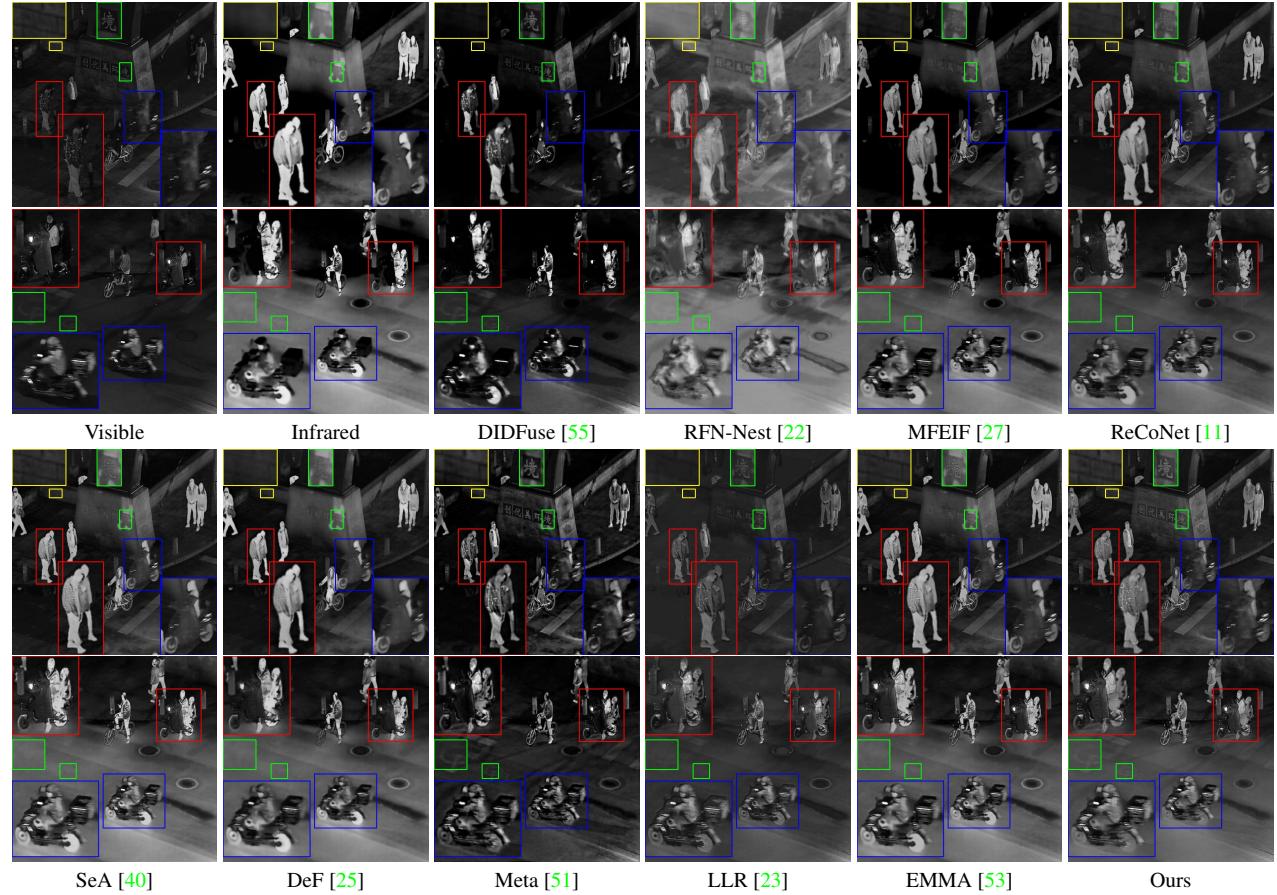


Figure 3. Visual comparison for “010018” (up) and “050131” (down) in LLVIP IVIF dataset.

	DID [55]	RFN [22]	MFE [27]	ReC [11]	SeA [40]	DeF [25]	Meta [51]	LLR [23]	EMM [53]	Ours	
SF	4.943	4.818	4.567	4.712	5.440	4.808	6.137	4.528	5.469	5.582	625
AG	2.036	2.251	1.882	2.031	2.432	2.148	2.985	1.787	2.430	2.540	626
MSE	0.043	0.060	0027	0.030	0.037	0.036	0.037	0.031	0.030	0.029	627
PSNR	13.71	12.22	15.78	14.26	14.59	14.40	14.46	14.83	15.42	15.56	628
MI	1.581	<u>1.601</u>	1.181	1.357	1.506	1.286	1.214	1.494	1.595	1.650	629
VIF	0.746	0.753	0.814	0.813	0.934	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	0.734	631
SSIM	1.044	0.999	1.298	1.398	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.

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 network 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 decline in overall performance.

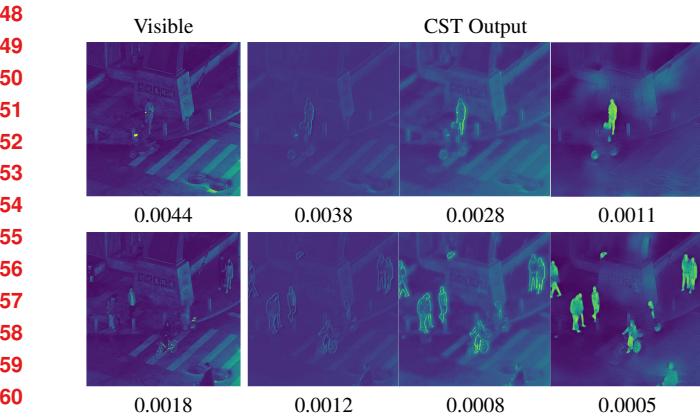


Figure 4. Visualization of the CST-module for “010018” (up), “010054” (down) in LLVIP IVIF dataset. The values represent the results of the perceptual loss.

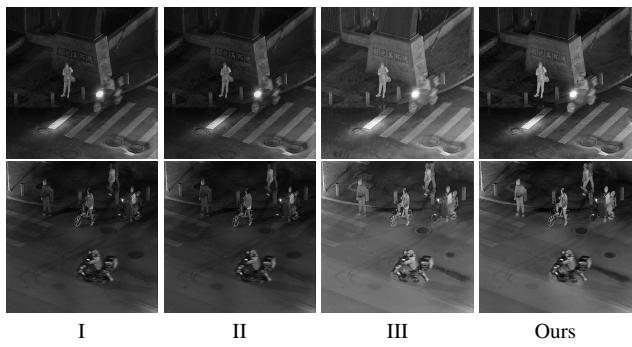


Figure 5. Ablation experiment for “010562” (up) and “050131” (down) in LLVIP IVIF dataset.

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 external validation is conducted. For infrared-visible object detection, the fused images generated by state-of-the-art mod-

	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	5.582	15.56	0.734	0.919	706
						707

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

els are evaluated using five classic detectors by comparing the AP value for person detection. The selected detectors include Faster R-CNN [5], YOLOv5 [37], SSD [28], RetinaNet [26] and Mask R-CNN [7]. For infrared-visible semantic segmentation, the segmentation network includes FCN [29], DeeplabV3 [2] and LSR-APP [8]. 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 [37] 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 [12], 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.

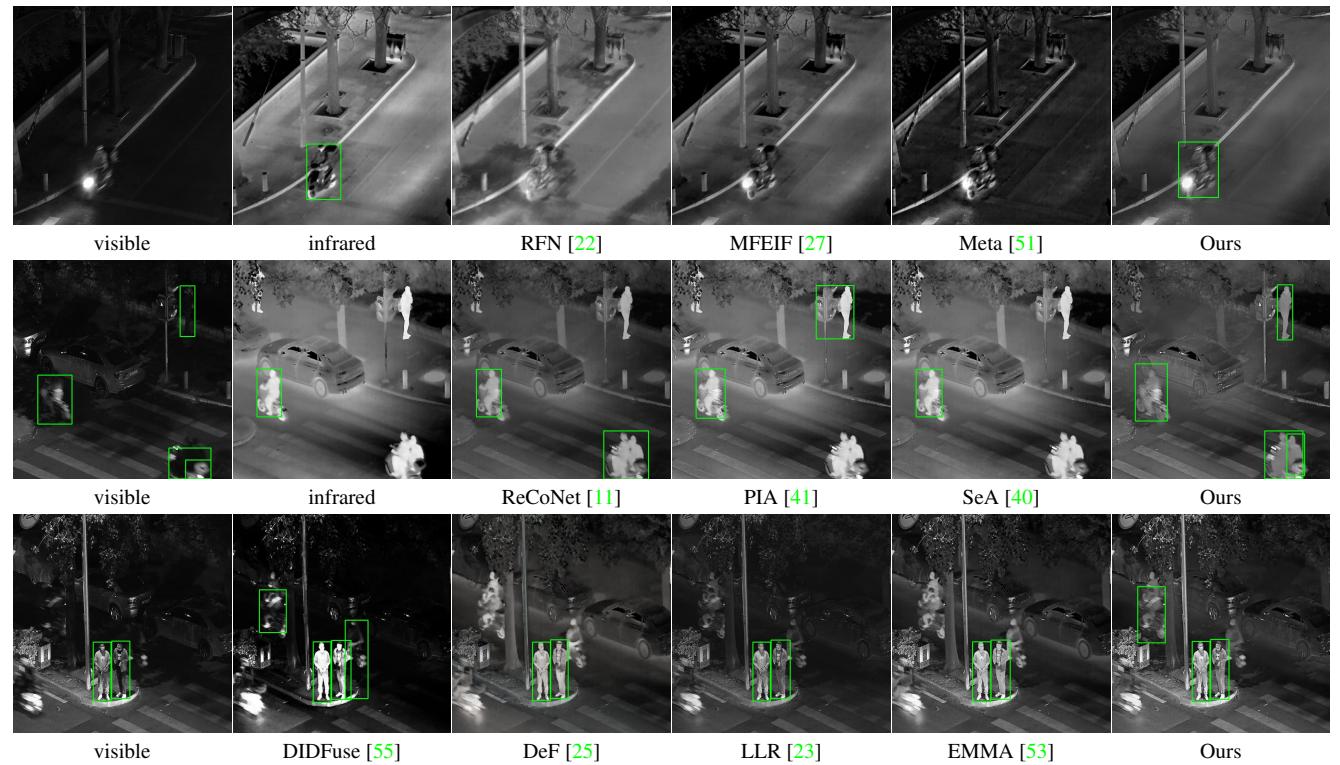


Figure 6. Detection results on source images and different fused image for “040145” (up), “080246” (middle) and “080786” (down) in LLVIP IVIF dataset.

	Faster RCNN [5]	YOLOv5 [37]	SSD [28]	Retina [26]	Mask RCNN [7]
VI	30.52	39.26	40.40	43.08	55.20
IR	34.03	37.80	39.29	41.28	48.97
DID [55]	32.05	39.11	44.97	43.56	54.38
RFN [22]	15.15	27.65	16.52	27.63	33.88
MFE [27]	32.48	42.57	39.67	41.00	50.94
ReC [11]	37.59	46.10	<u>46.65</u>	42.23	54.68
SeA [40]	36.25	44.54	46.13	46.84	55.24
DeF [25]	<u>37.57</u>	44.31	43.45	46.60	52.82
MFEIF [27]	36.43	43.49	49.36	52.19	56.42
LLR [23]	33.55	53.37	45.02	45.90	52.83
EMMA [53]	34.34	45.91	45.02	43.67	53.16
Ours	36.58	<u>46.87</u>	42.85	<u>48.86</u>	57.24

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 [29]	DeeplabV3 [2]	LSR-APP [8]
DID [55]	47.91	48.12	48.07
RFN [22]	44.69	46.49	47.99
MFE [27]	48.23	48.51	48.37
ReC [11]	48.32	48.70	48.57
SeA [40]	48.34	48.63	<u>48.53</u>
DeF [25]	<u>48.37</u>	48.52	48.52
Meta [51]	47.99	48.32	48.26
LLR [23]	47.97	48.29	48.13
EMMA [53]	48.10	48.37	48.41
Ours	48.48	<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 918
 865 Hartwig Adam. Rethinking atrous convolution for semantic 919
 866 image segmentation. *arXiv preprint arXiv:1706.05587*, 920
 867 2017. 7, 8
 868 [3] Yu Fu, Xiao-Jun Wu, and Tariq Durrani. Image fusion based 921
 869 on generative adversarial network consistent with perception. 922
 870 *Information Fusion*, 72:110–125, 2021. 1
 871 [4] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. 923
 872 Image style transfer using convolutional neural networks. In 924
 873 *CVPR*, June 2016. 2
 874 [5] Ross Girshick. Fast r-cnn. In *ICCV*, pages 1440–1448, 925
 875 2015. 7, 8
 876 [6] Rafael C Gonzalez, Richard E Woods, and Steven L Eddins. 926
 877 *Digital Image Processing Using MATLAB*. Prentice Hall, 927
 878 2009. 3, 4
 879 [7] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 928
 880 Mask r-cnn. In *ICCV*, pages 2961–2969, 2017. 7, 8
 881 [8] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh 929
 882 Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, 930
 883 Ruoming Pang, Vijay Vasudevan, Quoc V. Le, and Hartwig 931
 884 Adam. Searching for mobilenetv3. In *ICCV*, October 2019. 932
 885 7, 8
 886 [9] Hai-Miao Hu, Jiawei Wu, Bo Li, Qiang Guo, and Jin Zheng. 933
 887 An adaptive fusion algorithm for visible and infrared videos 934
 888 based on entropy and the cumulative distribution of gray levels. 935
 889 *IEEE T MULTIMEDIA*, 19(12):2706–2719, 2017. 1
 890 [10] Xun Huang and Serge Belongie. Arbitrary style transfer in 936
 891 real-time with adaptive instance normalization. In *ICCV*, 937
 892 2017. 2
 893 [11] Zhanbo Huang, Jinyuan Liu, Xin Fan, Risheng Liu, Wei 938
 894 Zhong, and Zhongxuan Luo. Recone: Recurrent correction 939
 895 network for fast and efficient multi-modality image fusion. 940
 896 In *ECCV*, pages 539–555. Springer, 2022. 1, 2, 5, 6, 8
 897 [12] Xinyu Jia, Chuang Zhu, Minzhen Li, Wenqi Tang, and Wenli 941
 898 Zhou. Llivip: A visible-infrared paired dataset for low-light 942
 899 vision. In *ICCV*, pages 3496–3504, 2021. 5, 7
 900 [13] Lihua Jian, Xiaomin Yang, Zheng Liu, Gwanggil Jeon, 943
 901 Minglei Gao, and David Chisholm. Sedrfuse: A symmetric 944
 902 encoder-decoder with residual block network for infrared 945
 903 and visible image fusion. *IEEE T INSTRUM MEAS*, 70:1– 946
 904 15, 2020. 1
 905 [14] Yongcheng Jing, Xiao Liu, Yukang Ding, Xinchao Wang, 947
 906 Errui Ding, Mingli Song, and Shilei Wen. Dynamic instance 948
 907 normalization for arbitrary style transfer. In *AAAI*, 2020. 2
 908 [15] Yongcheng Jing, Yang Liu, Yezhou Yang, Zunlei Feng, 949
 909 Yizhou Yu, Dacheng Tao, and Mingli Song. Stroke controllable 950
 910 fast style transfer with adaptive receptive fields. In *ECCV*, 951
 911 2018. 2
 912 [16] Jing jing Zong and Tian shuang Qiu. Medical image fusion 952
 913 based on sparse representation of classified image patches. 953
 914 *Biomedical Signal Processing and Control*, 34:195–205, 954
 915 2017. 1
 916 [17] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual 955
 917 losses for real-time style transfer and super-resolution. In 956
 918 *ECCV*, pages 694–711. Springer, 2016. 2, 3
 919 [18] Gihyun Kwon and Jong Chul Ye. Clipstyler: Image style 920
 920 transfer with a single text condition. In *CVPR*, pages 18062– 921
 921 18071, 2022. 2
 922 [19] Fayed Lahoud and Sabine Susstrunk. Ar in vr: Simulating 922
 923 infrared augmented vision. In *ICIP*, pages 3893–3897. IEEE, 923
 924 2018. 1
 925 [20] Hui Li and Xiao-Jun Wu. Densefuse: A fusion approach to 924
 926 infrared and visible images. *IEEE TIP*, 28(5):2614–2623, 925
 927 May 2019. 1
 928 [21] Hui Li, Xiao-Jun Wu, and Tariq Durrani. NestFuse: An 926
 929 Infrared and Visible Image Fusion Architecture based on Nest 927
 930 Connection and Spatial/Channel Attention Models. *IEEE T INSTRUM MEAS*, 69(12):9645–9656, 2020. 1
 931 [22] Hui Li, Xiao-Jun Wu, and Josef Kittler. Rfn-nest: An end-to- 931
 932 end residual fusion network for infrared and visible images. 932
 933 *Information Fusion*, 73:72–86, March 2021. 5, 6, 8
 934 [23] Hui Li, Tianyang Xu, Xiao-Jun Wu, Jiwen Lu, and Josef Kittler. 933
 935 LRRNet: A novel representation learning guided fusion 934
 936 framework for infrared and visible images. *IEEE TPAMI*, 45(9):11040–11052, 2023. 5, 6, 8
 937 [24] Shutao Li, Bin Yang, and Jianwen Hu. Performance 937
 938 comparison of different multi-resolution transforms for image 938
 939 fusion. *Information Fusion*, 12(2):74–84, 2011. 1
 940 [25] Pengwei Liang, Junjun Jiang, Xianming Liu, and Jiayi Ma. 940
 941 Fusion from decomposition: A self-supervised decomposition 941
 942 approach for image fusion. In *ECCV*, 2022. 5, 6, 8
 943 [26] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and 942
 944 Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 943
 945 pages 2980–2988, 2017. 7, 8
 946 [27] Jinyuan Liu, Xin Fan, Ji Jiang, Risheng Liu, and Zhongxuan 945
 947 Luo. Learning a deep multi-scale feature ensemble 946
 948 and an edge-attention guidance for image fusion. *TCSVT*, 947
 949 32(1):105–119, 2021. 1, 2, 5, 6, 8
 950 [28] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian 949
 951 Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. 950
 952 Ssd: Single shot multibox detector. In *ECCV*, pages 951
 953 21–37. Springer, 2016. 7, 8
 954 [29] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully 953
 955 convolutional networks for semantic segmentation. In *CVPR*, 954
 956 pages 3431–3440, 2015. 7, 8
 957 [30] Jiayi Ma, Yong Ma, and Chang Li. Infrared and visible 956
 958 image fusion methods and applications: A survey. *Information 957
 959 Fusion*, 45:153–178, 2019. 5
 960 [31] Jiayi Ma, Linfeng Tang, Fan Fan, Jun Huang, Xiaoguang 958
 961 Mei, and Yong Ma. Swinfusion: Cross-domain long-range 959
 962 learning for general image fusion via swin transformer. *JAS*, 960
 963 9(7):1200–1217, 2022. 1
 964 [32] Jiayi Ma, Han Xu, Junjun Jiang, Xiaoguang Mei, and Xiao- 962
 965 Ping Zhang. Ddcgan: A dual-discriminator conditional 963
 966 generative adversarial network for multi-resolution image 964
 967 fusion. *IEEE TIP*, 29:4980–4995, 2020. 1
 968 [33] Jiayi Ma, Wei Yu, Pengwei Liang, Chang Li, and Junjun 966
 969 Jiang. Fusiongan: A generative adversarial network for 967
 970 infrared and visible image fusion. *Information Fusion*, 48:11– 968
 971 26, 2019. 1, 2
 972 [34] Bikash Meher, Sanjay Agrawal, Rutuparna Panda, and Ajith 969
 973 Abraham. A survey on region based image fusion methods. 970
 974 *Information Fusion*, 48:119–132, 2019. 1

- 972 [35] Lukas Mehl, Azin Jahedi, Jenny Schmalfuss, and Andrés
973 Bruhn. M-fuse: Multi-frame fusion for scene flow estima-
974 tion. In *WACV*, pages 2020–2029, 2023. 1
- 975 [36] Ujwala Patil and Uma Mudengudi. Image fusion using hier-
976 archical pca. In *ICIP*, pages 1–6. IEEE, 2011. 1
- 977 [37] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali
978 Farhadi. You only look once: Unified, real-time object de-
979 tection. In *CVPR*, pages 779–788, 2016. 7, 8
- 980 [38] Yiming Sun, Bing Cao, Pengfei Zhu, and Qinghua Hu. Det-
981 fusion: A detection-driven infrared and visible image fusion
982 network. In *ACM MM*, pages 4003–4011, 2022. 2
- 983 [39] Linfeng Tang, Yuxin Deng, Yong Ma, Jun Huang, and Jiayi
984 Ma. Superfusion: A versatile image registration and fusion
985 network with semantic awareness. *JAS*, 9(12):2121–2137,
986 2022. 1
- 987 [40] Linfeng Tang, Jiteng Yuan, and Jiayi Ma. Image fusion in
988 the loop of high-level vision tasks: A semantic-aware real-
989 time infrared and visible image fusion network. *Information
990 Fusion*, 82:28–42, 2022. 1, 2, 5, 6, 8
- 991 [41] Linfeng Tang, Jiteng Yuan, Hao Zhang, Xingyu Jiang, and
992 Jiayi Ma. Piafusion: A progressive infrared and visible im-
993 age fusion network based on illumination aware. *Information
994 Fusion*, 83-84:79–92, 2022. 2, 8
- 995 [42] Linfeng Tang, Hao Zhang, Han Xu, and Jiayi Ma. Deep
996 learning-based image fusion: A survey. *Journal of Image
997 and Graphics*, 28(1):3–36, 2023. 1
- 998 [43] Alexander Toet and Maarten A. Hogervorst. Progress in
999 color night vision. *Optical Engineering*, 51:010901 –
0100 010901, 2012. 2
- 1001 [44] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Repre-
1002 sentation learning with contrastive predictive coding. *CoRR*,
1003 abs/1807.03748, 2018. 2
- 1004 [45] Zijie Wu, Zhen Zhu, Junping Du, and Xiang Bai. Ccpl: Con-
1005 trastive coherence preserving loss for versatile style transfer.
1006 In *ECCV*, pages 189–206. Springer, 2022. 2
- 1007 [46] Han Xu, Pengwei Liang, Wei Yu, Junjun Jiang, and Jiayi Ma.
1008 Learning a generative model for fusing infrared and visible
1009 images via conditional generative adversarial network with
1010 dual discriminators. In *IJCAI*, pages 3954–3960, 2019. 1
- 1011 [47] Han Xu, Xinya Wang, and Jiayi Ma. Drf: Disentangled rep-
1012 resentation for visible and infrared image fusion. *IEEE T
INSTRUM MEAS*, 70:1–13, 2021. 1
- 1013 [48] Han Xu, Hao Zhang, and Jiayi Ma. Classification saliency-
1014 based rule for visible and infrared image fusion. *IEEE TCI*,
1015 7:824–836, 2021. 1
- 1016 [49] Mingde Yao, Zhiwei Xiong, Lizhi Wang, Dong Liu, and
1017 Xuejin Chen. Spectral-depth imaging with deep learning
1018 based reconstruction. *Optics express*, 27(26):38312–38325,
1019 2019. 1
- 1020 [50] Hao Zhang, Han Xu, Xin Tian, Junjun Jiang, and Jiayi Ma.
1021 Image fusion meets deep learning: A survey and perspective.
1022 *Information Fusion*, 76:323–336, 2021. 1
- 1023 [51] Wenda Zhao, Shigeng Xie, Fan Zhao, You He, and Huchuan
1024 Lu. Metafusion: Infrared and visible image fusion via meta-
1025 feature embedding from object detection. In *CVPR*, pages
13955–13965, 2023. 5, 6, 8
- 1026 [52] Zixiang Zhao, Haowen Bai, Jiangshe Zhang, Yulun Zhang,
1027 Shuang Xu, Zudi Lin, Radu Timofte, and Luc Van Gool.
Cddfuse: Correlation-driven dual-branch feature decomposi-
1028 tion for multi-modality image fusion. In *CVPR*, pages 5906–
1029 5916, June 2023. 2
- 1030 [53] Zixiang Zhao, Haowen Bai, Jiangshe Zhang, Yulun Zhang,
1031 Kai Zhang, Shuang Xu, Dongdong Chen, Radu Timofte, and
1032 Luc Van Gool. Equivariant multi-modality image fusion. In
1033 *CVPR*, June 2024. 5, 6, 8
- 1034 [54] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu, and
1035 Jiangshe Zhang. Bayesian fusion for infrared and visible im-
1036 ages. *Signal Processing*, 177:107734, Dec. 2020. 1
- 1037 [55] Zixiang Zhao, Shuang Xu, Chunxia Zhang, Junmin Liu,
1038 Jiangshe Zhang, and Pengfei Li. Didfuse: Deep image de-
1039 composition for infrared and visible image fusion. In *PRI-
1040 CAI, IJCAI-PRICAI-2020. IJCAI*, July 2020. 1, 2, 5, 6, 8
- 1041
- 1042
- 1043
- 1044
- 1045
- 1046
- 1047
- 1048
- 1049
- 1050
- 1051
- 1052
- 1053
- 1054
- 1055
- 1056
- 1057
- 1058
- 1059
- 1060
- 1061
- 1062
- 1063
- 1064
- 1065
- 1066
- 1067
- 1068
- 1069
- 1070
- 1071
- 1072
- 1073
- 1074
- 1075
- 1076
- 1077
- 1078
- 1079