# IG-GAN: Interactive Guided Generative Adversarial Networks for Multimodal Image Fusion

Chenhong Sui, *Member, IEEE*, Guobin Yang, Danfeng Hong, *Senior Member, IEEE*,
Haipeng Wang, *Member, IEEE*, Jing Yao, *Member, IEEE*, Peter M. Atkinson,
and Pedram Ghamisi, *Senior Member, IEEE*

*Abstract*—Multimodal image fusion has recently garnered increasing interest in the field of remote sensing. By leveraging the complementary information in different modalities, the fused results may be more favorable in characterizing objects of interest, thereby increasing the chance of a more comprehensive and accurate perception of the scene. Unfortunately, most existing fusion methods tend to extract modality-specific features independently without considering intermodal alignment and complementarity, leading to a suboptimal fusion process. To address this issue, we propose a novel interactive generative adversarial network (IG-GAN), for the task of multimodal image fusion. IG-GAN comprises guided dual streams tailored for enhanced learning of details and content, as well as cross-modal consistency. Specifically, a details-guided interactive running-in module ($GIR_1$) and a content-guided interactive running-in module ($GIR_2$) are developed, with the stronger modality serving as guidance for detail richness or content integrity, and the weaker one assisting. To fully integrate multigranularity features from dual-modality, a hierarchical fusion and reconstruction branch is established. Specifically, a shallow interactive fusion (SIF) module followed by a multilevel interactive fusion (MIF) module is designed to aggregate multilevel local and long-range features. Concerning feature decoding and fused image generation, a high-level interactive fusion and reconstruction module (HRM) is further developed. In addition, to empower the fusion network to generate fused images with complete content, sharp edges, and high fidelity without supervision, a loss function facilitating the mutual game between the generator and two discriminators is also formulated. Comparative experiments with 14 state-of-the-art methods are conducted on three datasets. Qualitative and quantitative results indicate that IG-GAN exhibits obvious superiority in terms of both visual effect and quantitative metrics. Moreover, experiments on two RGB-IR object detection datasets are also conducted, which demonstrate that IG-GAN can enhance the accuracy of object detection by integrating complementary information from different modalities. The code will be available at https://github.com/flower6top.

*Index Terms*—Generative adversarial fusion, multimodal image fusion, transformer.

Chenhong Sui is with the Aerospace Information Research Institute, Chinese Academy of Sciences, Beijing 100094, China, and also with the School of Physics and Electronic Information, Yantai University, Yantai 264005, China (e-mail: sui6662008@163.com).

Guobin Yang is with the School of Physics and Electronic Information, Yantai University, Yantai 264005, China, and also with the Weifang Vocational College, Weifang 261041, China (e-mail: yangguobing@s.ytu.edu.cn).

Danfeng Hong is with the Aerospace Information Research Institute, Chinese Academy of Sciences, Beijing 100094, China, and also with the School of Electronic, Electrical and Communication Engineering, University of Chinese Academy of Sciences, Beijing 100094, China (e-mail: hongdf@aircas.ac.cn).

Haipeng Wang is with the School of Computer Science and Technology, Harbin Institute of Technology, Harbin 150001, China, and also with the Institute of Information Fusion, Naval Aviation University, Yantai 530014, China (e-mail: whp5691@163.com).

Jing Yao is with the Aerospace Information Research Institute, Chinese Academy of Sciences, Beijing 100094, China (e-mail: yaojing@aircas.ac.cn).

Peter M. Atkinson is with the Faculty of Science and Technology, Lancaster University, LA1 4YW Lancaster, U.K. (e-mail: pma@lancaster.ac.uk).

Pedram Ghamisi is with the Helmholtz-Zentrum Dresden-Rossendorf (HZDR), 09599 Dresden, Germany, and also with Lancaster University, LA1 4YW Lancaster, U.K. (e-mail: p.ghamisi@gmail.com).

Digital Object Identifier 10.1109/TGRS.2024.3433619

## I. INTRODUCTION

REMARKABLE progress in sensor technology makes it possible to acquire multimodal images of the same scene [1]. However, influenced by sensor imaging mechanisms and the complex ground environment, single-mode images often cannot provide sufficient and detailed scene information [2], [3], [4], [5], [6], [7]. For example, thermal infrared images contain the radiation signal from objects thus providing additional information to the visible spectrum but weak texture information [8], [9]. Synthetic aperture radar (SAR) images possess rich polarimetric scattering information of ground objects regardless of cloud cover but are seriously deficient in object details [10]. In comparison with infrared and SAR data, visible images are highly susceptible to weather and illumination variations despite providing great texture detail [8]. Therefore, multimodal image fusion is of great significance to simultaneously compensate for the content or detail deficiency of a single-mode sensor and enhance the information provided by images [11]. Fig. 1 presents an example of optical and SAR image fusion. As shown in Fig. 1, the SAR image can perceive objects in the red box, while its depiction of object details is weak. The optical image lacks complete perception of the objects and can only perceive the objects partially. For objects within the green box, the SAR image is unable to perceive objects that are subject to human interference, and the edges of the perceived objects are blurred, with unclear detailed features. In comparison, the optical image presents clearer object contours, and there is a significant difference in contrast between the object and the background.
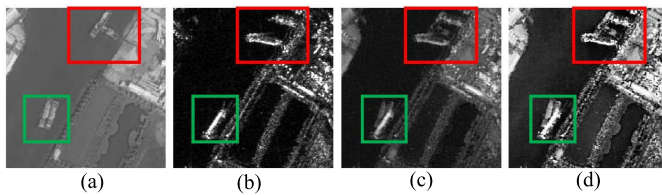
Fig. 1. Example of optical and SAR image fusion. (a) Optical image. (b) SAR image. (c) Fusion result of GANMcC [12]. (d) Fusion result from IG-GAN. It is evident that compared to GANMAC, IG-GAN can leverage the advantages of SAR and optical images to achieve a more complete perception of objects within red and green boxes, and fused images have significant advantages in sharpening edges and contrast.

Image fusion has benefitted from tremendous progress in the past decade [13]. Traditional fusion methods mainly depend on fixed transformations of source images and manually designed features for fusion [14]. These methods can be roughly divided into five categories (i.e., the multiscale transform-based [15], the sparse representation-based [16], the subspace-based [17], the saliency-based [18], and the hybrid method [19]). Despite the high efficiency of most traditional methods, artificially designed feature extraction or fusion rules cannot fully capture the characteristics of different modalities, resulting in limited fusion performance.

Deep fusion methods, as a new paradigm beyond traditional methods, rely on the powerful modeling capabilities of deep neural networks for adaptive information extraction and learning the fusion rules [20]. Examples include convolutional neural networks (CNNs) [21] and autoencoders (AEs) [22], [23] for modality-specific feature extraction and fusion. For example, Masi et al. [24] presented the CNN-based pansharpening methods with three convolutional layers, which can integrate the unique advantages of a panchromatic (PAN) image and a multispectral (MS) image. Zheng et al. [25] introduced the dual-attention to inject the spatial details of PAN into a hyperspectral (HS) image for superresolution. Yao et al. [26] adopted the encoder–decoder-based architecture for HS and Lidar fusion. U2Fusion [27] proposes an unsupervised fusion method based on a densely connected network for various image fusion tasks. RFN-Nest [28] introduces an end-to-end residual architecture-based fusion network for infrared and visible image fusion. For the complementary fusion of texture and intensity information, PMGI [29] advocates constructing a dual-path-based unified fusion network for better information extraction. In addition, since a generative adversarial network (GAN) is capable of generating high-fidelity fusion images through the min–max game between the generator and discriminator, it is naturally utilized for multimodal image fusion in various studies [4], [30], [31], [32]. For example, Zhang et al. [33] developed a GAN-based two-stage framework for spatiotemporal image fusion. Gao et al. [34] established a GAN-based network for SAR-optical fusion and cloud removal. In addition, in [35], GAN is successfully applied to HS pansharpening, which adopts the 3-D convolutional network to capture desirable high-frequency residuals. FusionGAN [36] establishes an adversarial game for the fusion of infrared and visible images, effectively avoiding the artificial design of fusion rules. DDcGAN [37]

extends FusionGAN by adding a discriminator, which helps the fused image fully retain information from multisource images. Multiclassification GAN (GANMcC) [12] recognizes that the details and contrast in the infrared image are not as significant as in the visible image. Therefore, it employs multiple classifiers as discriminators to output the probability that the input is an infrared image or a visible image, producing a visually appealing fused image. These deep fusion methods have significantly advanced the field of image fusion. Nevertheless, they are limited by the inherent constraints of convolutional operators in capturing intermodality or intramodality global context.

To address the above issue, some researchers have turned to the transformer structure [38], which embraces long-term modeling capabilities. For example, IFT [39] introduces the transformer to image fusion and achieves performance similar to the CNN architecture. Hu et al. [40] introduced the transformer to HS-MS fusion, where the structured embedding matrix is injected into the transformer encoder to learn the residual map. Bandara and Patel [41] leveraged the transformer for HS pansharpening, in which the modality-specific feature extractors are designed to capture textural details. Liu et al. [42] proposed an attention-based multiscale transformer network to model contextual information in bi-temporal images for change detection. Feng et al. [43] proposed the center attention transformer (CAT) with a stratified spatial–spectral token for HSI classification. Li et al. [44] learned a coupling model-driven and data-driven paradigm to distinguish between the background and anomalies for HS anomaly detection. TGFuse [45] embeds the transformer into GAN for global visible and infrared image fusion. SwinFusion [46] utilizes the Swin transformer [47] to extract features from different sources and leverages cross-domain attention for feature fusion.

For the transformer-based methods, the construction of a global feature association is beneficial leading to better-fused images [48]. However, the approach generally conducts semantic mining for each modal independently, with no perception of the intermodal discrepancy and consistency. This could result in the underutilization of the modality-specific advantages and cross-modality commonalities while the overusing of invalid or noisy information.

To the above end, this article introduces an interactive generative adversarial network (IG-GAN) for multimodal images. Note that some modalities are capable of collecting complete scene information irrespective of weather or illumination variations, whereas visible images possess the merit of rich texture details. Enlightened by this, detail and content streams are first cooperatively established rather than independently for cross-modal complementarity and consistency enhancement. Specifically, a details-guided interactive running-in module (GIR$_1$) and a content-guided interactive running-in module (GIR$_2$) are developed. This is conducive to ensuring that the advantages of the dominant modality can be fully utilized, while the other modality can assist in cross-modal feature alignment, enhancement, and complementary fusion. Regarding the comprehensive integration of dual-stream features, we further construct a hierarchical fusion and reconstruction

branch. In this branch, both a shallow interactive fusion (SIF) module and a multilevel interactive fusion (MIF) module are built. Furthermore, for fine decoding and fused image generation, we propose a high-level interactive fusion and reconstruction module (HRM) capable of absorbing multimodal, multigranularity local–global features. In addition, to guarantee that the fusion network can generate complete, sharpened, and high-fidelity images, we design a loss function involving the mutual game between the generator and two discriminators. Qualitative and quantitative experimental results show that compared with state-of-the-art methods, IG-GAN exhibits apparent superiority over others on four commonly used benchmarks. Our main contributions can be summarized as follows.

1) A novel unsupervised multimodal image fusion method, called IG-GAN, is proposed to fully explore the modal-specific advantageous information regarding detail richness and content completeness while enhancing cross-modal commonalities collaboratively.

2) In the generator, a guided details stream and a guided content stream are established for multilevel inter-modality alignment, cooperation, and enhancement. Specifically, a $GIR_1$ and a $GIR_2$ are developed. This means that both the content and detail streams are built with multimodal interactive promotion rather than operating independently. In each stream, both the leading role of each dominant modality and the auxiliary contribution of the weaker one are considered.

3) Concerning dual-stream feature integration, the generator also involves a SIF module, a MIF module, and a high-level fusion and reconstruction module (HRM). This promotes multigranularity, multilevel integration of dual-stream features, and fused image generation.

4) To ensure the fusion performance of our IG-GAN without supervision, a novel loss function simultaneously involving detail richness, content integrity, and high fidelity is devised for network training. It boosts the mutual game between the generator and two discriminators. Qualitative and quantitative experimental results on four benchmark datasets demonstrate the superiority of IG-GAN.

The remainder of the article is organized as follows. Section II primarily reviews typical work related to IG-GAN. Section III gives a detailed description of IG-GAN and its core modules (e.g., $GIR_1$, $GIR_2$, SIF, MIF, and HRM). In Section IV, quantitative and qualitative experimental results and discussion are provided. Ultimately, the conclusions and future research are presented in Section V.

## II. RELATED WORKS

### A. Multiclassification GAN

In the context of infrared and visible image fusion, it is commonly acknowledged that the former lacks visual information such as details and textures, but can effectively depict significant objects under low illumination. Conversely, the latter is rich in detailed texture but is highly susceptible to changes in lighting conditions. Under low illumination,

scene information in visible images may appear incomplete or even completely missing with low contrast. To address the integration of saliency and details from infrared and visible images, a GANMcC [12] is designed.

The main idea of GANMcC is to transform multimodal image fusion into a simultaneous estimation of the contribution from infrared and visible image distributions. GANMcC comprises a generator and two discriminators. The generator aims to maximize the probability that the fused image comes from both visible and infrared images. In contrast, the discriminator adopts a multiclass classifier to determine that the fused image is neither an infrared image nor a visible image. With continuous adversarial learning, the generator can estimate the probability distribution of both infrared and visible images, enabling the generation of fused images with significant contrast and rich texture details. Equation (1) depicts the loss function in favor of both the generator and discriminator

$$\mathcal{L} = \mathcal{L}_G + \mathcal{L}_D \tag{1}$$

where $\mathcal{L}_G$ is the loss of the generator and $\mathcal{L}_D$ corresponds to the loss of the discriminators.

$\mathcal{L}_G$ consists of both content loss $\mathcal{L}_{G_{con}}$ and adversarial loss $\mathcal{L}_{G_{adv}}$. Regarding $\mathcal{L}_{G_{con}}$, to preserve the details and textures in multimodal images, intensity loss, and gradient loss are provided as shown in the following equation:

$$\mathcal{L}_{G_{con}} = \beta_1 \|I_f - I_{ir}\|_F^2 + \beta_2 \|\nabla I_f - \nabla I_{vis}\|_F^2 \\ + \beta_3 \|\nabla I_f - \nabla I_{ir}\|_F^2 + \beta_4 \|I_f - I_{vis}\|_F^2 \tag{2}$$

where $I_f$, $I_{ir}$, and $I_{vis}$ represent the fused image, the infrared image, and the visible image, respectively. $\nabla$ denotes the second-order gradient operator. $\| \cdot \|_F$ is the Frobenius-norm. $\beta_1$, $\beta_2$, $\beta_3$, and $\beta_4$ are the tradeoff parameters.

The adversarial loss $\mathcal{L}_{G_{adv}}$ can be described as

$$\mathcal{L}_{G_{adv}} = \left(D(I_f)[1] - d\right)^2 + \left(D(I_f)[2] - d\right)^2 \tag{3}$$

where $D$ is the discriminator, in which $d$ is the image modal label used to determine the type of fused image. $D(\cdot)[1]$ and $D(\cdot)[2]$ represent the probability that the fused image is the visible image or the infrared image, respectively.

Concerning discriminators, they are capable of judging that the distribution of the fused image is different from both the visible and infrared images. Therefore, the corresponding discriminator loss $\mathcal{L}_D$ is defined as

$$\mathcal{L}_D = \mathcal{L}_{D_{vis}} + \mathcal{L}_{D_{ir}} + \mathcal{L}_{D_{fused}} \tag{4}$$

where $\mathcal{L}_{D_{vis}}$ aims to measure the probability that the visible image is classified as an infrared image. Analogously, $\mathcal{L}_{D_{ir}}$ depicts the probability that the infrared image is classified as a visible image. Meanwhile, $\mathcal{L}_{D_{fused}}$ represents the probability that the fused image is classified as a visible or infrared image.

The mutual game between the generator and the discriminators, based on a loss function, is beneficial for producing a high-quality fused image. Unfortunately, due to the local modeling attributes of convolutional operators, the exploration of intermodality or intramodality global semantics is lacking.

## B. SwinFusion

SwinFusion [46] is a versatile image fusion framework based on cross-domain distance learning and the Swin transformer [47]. Unlike existing transformer-based image fusion methods that mainly focus on the interaction of information within a domain, SwinFusion takes a step further by exploring the contextual relationship between multisource images. To address this limitation, SwinFusion employs the Swin transformer as the backbone to model long-range dependencies between domains and design cross-domain attention for feature fusion. Specifically, the Swin transformer is introduced to delve deeper into the semantics extracted by a CNN from shallow features. Subsequently, cross-attention is utilized for effective feature fusion. This approach ensures the integration of cross-domain context information, leading to exceptional image fusion results. The total loss, denoted as $\mathcal{L}_{\text{total}}$, is composed of SSIM loss $\mathcal{L}_{\text{ssim}}$ [49], texture loss $\mathcal{L}_{\text{tex}}$, and intensity loss $\mathcal{L}_{\text{int}}$ as described in the following equation:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{ssim}} + \mathcal{L}_{\text{text}} + \mathcal{L}_{\text{int}} \tag{5}$$

where the SSIM loss $\mathcal{L}_{\text{ssim}}$ can be expressed as

$$\mathcal{L}_{\text{ssim}} = w_1\big(1 - \text{ssim}(\boldsymbol{I}_f, \boldsymbol{I_1})\big) + w_2\big(1 - \text{ssim}(\boldsymbol{I}_f, \boldsymbol{I_2})\big) \tag{6}$$

where ssim$(\cdot)$ represents the structural similarity operation, the balancing parameters $w_1$ and $w_2$ are set to 0.5. $\boldsymbol{I_1}$ and $\boldsymbol{I_2}$ denote the bi-modal images, and $\boldsymbol{I}_f$ is the fused image.

To characterize the richness of texture in images, the following equation further provides the texture loss $\mathcal{L}_{\text{tex}}$:

$$\mathcal{L}_{\text{tex}} = \frac{1}{\text{HW}}\big\| |\nabla \boldsymbol{I}_f| - \max(|\nabla \boldsymbol{I_1}|, |\nabla \boldsymbol{I_2}|)\big\|_1 \tag{7}$$

where $|\cdot|$ denotes the absolute operation, $\|\cdot\|_1$ is $l_1-$norm, and max$(\cdot)$ refers to the element-wise maximum selection. $H$ and $W$ denote the height and width of the image.

In addition, SwinFusion adopts the intensity loss $L_{\text{int}}$ to describe the element-by-element spatial discrepancy between the fusion image and the original bimodal images, as expressed in the following equation:

$$\mathcal{L}_{\text{int}} = \frac{1}{\text{HW}}\big\| \boldsymbol{I}_f - M(\boldsymbol{I_1}, \boldsymbol{I_2})\big\|_1 \tag{8}$$

where $M(\cdot)$ is an element-wise aggregation operation.

Note that SwinFusion is an encoder–decoder-based fusion network, which adopts the SSIM loss, text loss, and intensity loss to optimize the fusion network. In comparison, many GAN-based fusion methods [12], [50] leverage both the generator loss and the discriminator loss for fusion network optimization, which contribute to generating high-quality fused images, ensuring the naturalness and realism of the fusion results.

When delving deeper into the SwinFusion model, we noticed that the model adopts a cross-attention mechanism for feature fusion in the second stage of the Swin transformer, but lacks independent fusion branches. Consequently, this may limit the flexibility and efficiency of the model when dealing with complex data. In this regard, it is sensible to introduce an independent fusion branch to enhance the overall performance and adaptability of the model.

## III. Our Method

In this section, we first introduce the framework of our proposed IG-GAN. Then, the specific modules in the generator and discriminator networks are described, respectively. Finally, to optimize the designed network and enable it to produce complete and detailed images without supervision, a corresponding loss function is provided.

### A. Framework of IG-GAN

To yield high-quality fused images, multimodal image fusion needs to explore and utilize intermodal semantic consistency for better feature alignment and enhancement. In addition, intermodality complementary information is required to compensate for the deficiencies of single-source images. To address this, we propose a dual-stream IG-GAN through the mutual game between the generator and discriminators, as depicted in Fig. 2.

As depicted in Fig. 2, the generator comprises a guided details stream, a guided content stream, and a hierarchical feature fusion and image reconstruction branch. The former primarily focuses on complementary feature mining and aligned feature enhancement, while the latter is mainly responsible for further multiview and multilevel feature fusion, as well as fused image restoration. Specifically, in each stream, the guided interactive running-in modules (i.e., GIR$_1$ and GIR$_2$) are introduced based on four consecutive guided Swin transformer blocks. Concerning hierarchical fusion, we first establish a SIF module for the generation of multiview and low-level fusion features $\boldsymbol{S}_f$. Subsequently, a MIF module is utilized to produce both low and high-level fusion features $\boldsymbol{M}_f$. Following the ResNet approach, we advocate feeding both low-level and high-level fusion features into an aggregation block that involves concatenation and a $1 \times 1$ convolution. This is followed by three consecutive transformer blocks. Finally, after the patch expansion operation in the transformer, the aggregation block, accompanied by multiple transformer modules, serves as the fusion feature decoding and fused image reconstruction.

As an unsupervised image fusion network, two discriminators, namely, Discriminator$_1$ and Discriminator$_2$, are employed to distinguish between the fused images and source images. They play a crucial role in preventing significant discrepancies between the fused and original images and avoiding artifacts. Therefore, in Section III-B, we first introduce the guided dual-stream and hierarchical fusion and image restoration parts in the generator, respectively.

Influenced by the working principle of the sensor, there are significant differences in the details (such as texture and contrast) and integrity of various modal images in bad weather and low illumination. For example, SAR images can capture complete scenes unaffected by clouds, rain, fog, and lighting changes. Unfortunately, they lack texture and contrast information. On the other hand, optical images are typically rich in detailed textures, but are susceptible to weather and illumination changes, leading to potential severe loss or pollution of scene information. Therefore, it is wise to explore and
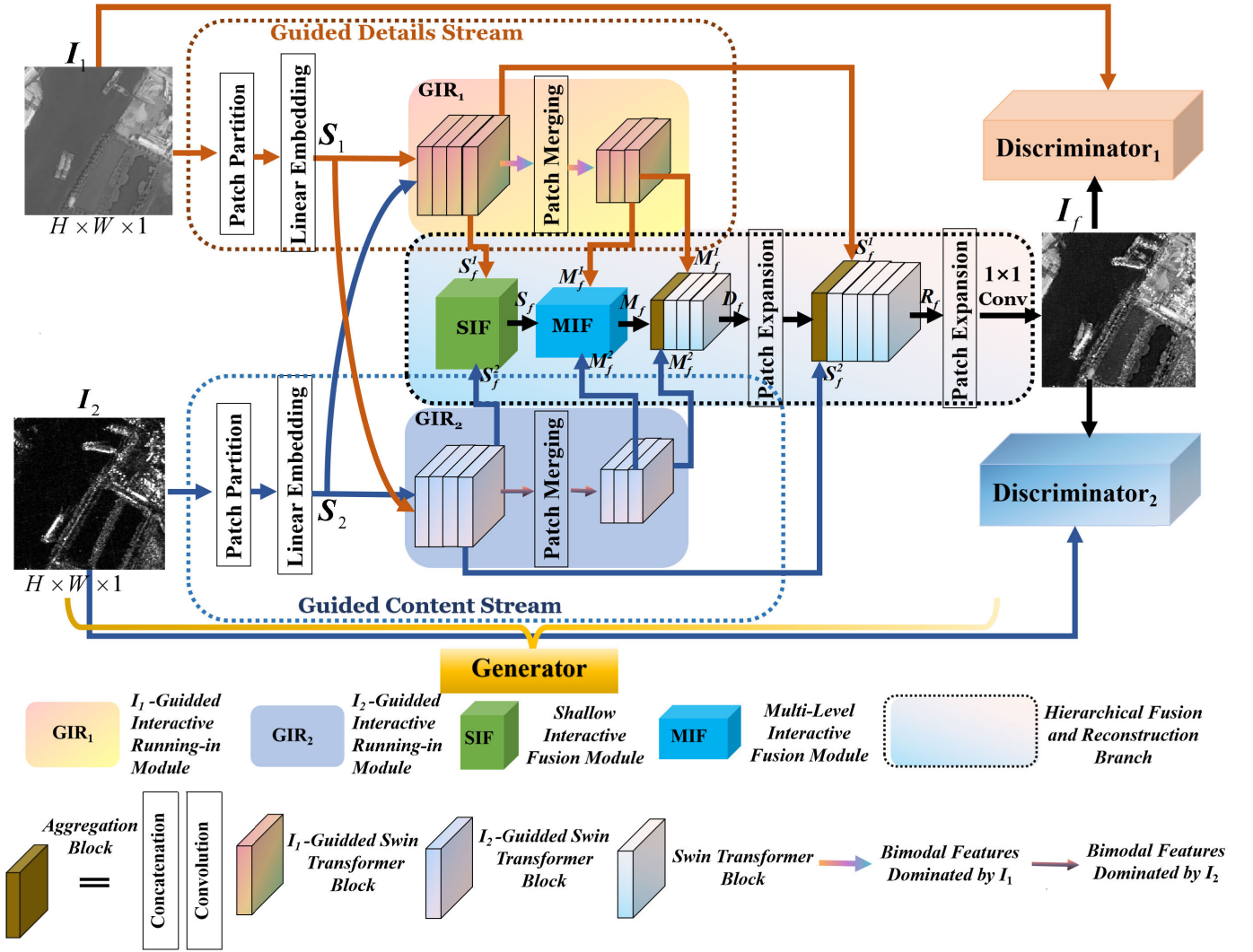
Fig. 2. Framework of IG-GAN. As a generative fusion network, IG-GAN has $I_1$-guided details stream, $I_2$-guided content stream, as well as multimodal hierarchical fusion and a reconstruction branch in the generator. This facilitates the comprehensive exploration, alignment, and enhancement of cross-modal semantic consistency. Meanwhile, modality-wise different advantages respecting detail richness and content completeness are fully considered. To enhance the fidelity of the fused images, dual discriminators are leveraged to engage in a game against the generator.

aggregate their complementary semantics in terms of rich details and content integrity, respectively. In this context, we construct a dual-stream architecture for better feature extraction and fusion, where different modalities dominate the details stream and content stream based on their performance advantages.

Let $I_1$ be the input optical image with size $H \times W \times 1$, while $I_2$ denotes the corresponding infrared or SAR image with the same size. Then, as shown in Fig. 2, $I_1$ should play a dominant role in the details stream, as it is beneficial for retaining contextual details. Regarding $I_2$, aligning spatial and semantic information with $I_1$ allows $I_1$ to assist $I_2$ in capturing more complete scene content while preserving effective detail information. For the content stream, $I_2$ should assume a critical role, providing more comprehensive scene content that remains unaffected by lighting and weather conditions. In this context, $I_1$ has the auxiliary effect of enhancing detailed textures by aligning spatially and semantically with $I_2$.

*B. Guided Dual-Stream*

Note that intermodality features have both uniqueness and relevance, involving some common semantic information. In this view, compared with common independent feature extraction mechanisms from various modalities, it is sensible to explore cross-modal information cooperatively, which helps strengthen the heterogeneous or complementary information. In this regard, we give a $I_1$-guided interactive running-in module (GIR$_1$) and a $I_2$-guided interactive running-in module (GIR$_2$) for comprehensive feature alignment, cooperation, and consistency feature enhancement.

*1) Guided Interactive Running-In Module:* The guided interactive running-in module aims to explore the commonality and uniqueness among modalities and allows multimodal images to be included jointly. As described in Fig. 2, after passing through the U-shaped network, the features from $I_1$ and $I_2$ are concatenated to form joint features in two streams. Then, they are fed into GIR$_1$ and GIR$_2$ with the first half of the channels dominant and the second half of the channels
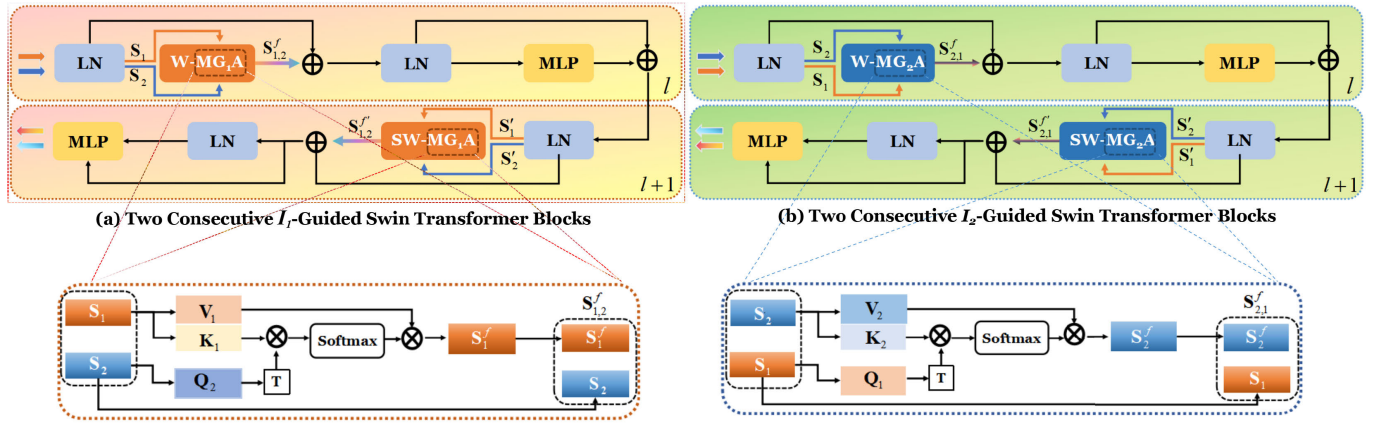
**(a) Two Consecutive $I_1$-Guided Swin Transformer Blocks**

**(b) Two Consecutive $I_2$-Guided Swin Transformer Blocks**

Fig. 3. Illustration of the proposed $I_1$- and $I_2$-guided Swin transformer blocks used in GIR$_1$ and GIR$_2$, respectively. Different from the standard Swin transformer block only involving a single modality as input, the guided ones corresponding to (a) and (b) take the features from dual-modality as input. In specific, for (a), owing to $I_1$-guided multihead attention in W-MG$_1$A and SW-MG$_1$A, the outputs of details stream are mainly dominated by features from $I_1$, whereas the ones from $I_2$ helps feature correction and enhancement via cross-modal feature interaction alignment. Regarding (b), $I_2$-guided multihead attention is given in W-MG$_2$A and SW-MG$_2$A. In this case, $I_2$ dominates the content stream, while $I_1$ helps to strengthen consistency features through feature interaction alignment.

auxiliary. This ensures that $I_1$ and $I_2$ play different roles in the two streams based on their contribution. Compared with independent feature extraction from each modality, this is conducive to aligning and enhancing the commonality between modalities. Moreover, heterogeneous or complementary information is strengthened.

Note that the Swin transformer possesses multiview perception and long-range modeling ability [51]. In addition, with the utilization of patch merging, consecutive Swin transformer modules can capture semantic features at lower scales. Motivated by the merits of the Swin transformer, we have developed the $I_1$-guided and $I_2$-guided Swin transformer blocks used in GIR$_1$ and GIR$_2$ for full interaction alignment and enhancement for $I_1$ and $I_2$. Specifically, inspired by cross-attention, Swin transformer's multihead attention is improved to $I_1$-guided and $I_2$-guided multihead attention mechanisms.

As demonstrated in Fig. 2, for the guided details stream, $4 \times I_1$-guided Swin transformer blocks are built to extract detailed semantic information dominated by optical images. Meanwhile, the features from infrared or SAR images are leveraged to assist in semantic alignment, correction, and consistency enhancement. Note that patch merging is an effective resolution reduction mechanism that does not cause information loss [39]. Therefore, to capture global semantics at larger scales, the patch merging operation is leveraged in both GIR$_1$ and GIR$_2$. Assuming that $S_1$ and $S_2$ have a spatial size of $(H/4) \times (W/4)$ after patch partition, then, after patch merging, the spatial size of the feature map is reduced to $(H/8) \times (W/8)$. After that, $3 \times I_1$-guided Swin transformer blocks are further adopted to extract semantic information.

Analogous to the details stream, for the content stream, a backbone content extraction module is constructed, which contains $4 \times I_2$-guided Swin transformer blocks dominated by SAR or infrared images. At the same time, the features from optical images are also introduced into this stream, which play an auxiliary role in semantic consistency enhancement via interactive alignment. After a patch merging operation, $3 \times$

$I_2$-guided Swin transformer blocks are capable of exploring interaction features with size $(H/8) \times (W/8)$.

For ease of understanding, before analyzing the proposed guided attention mechanisms and Swin transformer, we first recall the standard ones. Let $S$ be the input. Then, the multihead self-attention MHead($S$) in W-MSA and SW-MSA of the standard Swin transformer blocks can be described as

$$\text{MHead}(S) = \text{Concat}\big(\textbf{\textit{head}}^1, \textbf{\textit{head}}^2, \ldots, \textbf{\textit{head}}^h\big)W^O \quad (9)$$

where $\textbf{\textit{head}}^i$ represents the $i$th head of MHead($S$). $h$ denotes the number of heads. $\textbf{\textit{head}}^i$ can be obtained as defined in the following equation:

$$\textbf{\textit{head}}^i = \text{soft max}\left(\frac{K^i\big(Q^i\big)^\text{T}}{\sqrt{d_k}} + B\right)V^i \quad (10)$$

where $Q^i$, $K^i$, and $V^i$ correspond to different linear mappings of $S$, which satisfy $Q^i = S \cdot W_i^Q$, $K^i = S \cdot W_i^K$, and $V^i = S \cdot W_i^V$. $d_k$ denotes the dimension of $K^i$, $B$ is the relative position encoding that can be learned.

As shown in (10), despite the global statistics of the standard Swin transformer, the existing multihead self-attention tends to focus on a single source input $S$ [46], [52]. Consequently, the association between modalities is neglected, which is not conducive to the enhancement of intermodality consistency and the fusion of complementarity. In this connection, it is necessary to make reasonable use of the auxiliary role of other modalities while affirming the dominant role of stronger modalities. This means that we should start with the importance of each mode, and delegate our attention to the important leading mode, while other modes assist in enhancing.

For better illustration, Fig. 3 provides a visual representation of the $I_1$- and $I_2$-guided Swin transformer blocks. Specifically, the $I_1$- and $I_2$-guided attention mechanisms are displayed as well.

As described in Fig. 3, the guided Swin transformer blocks contribute to the mutual running-in of cross-modality features via a bimodal guided attention mechanism. For example,

W-MG$_1$A and SW-MG$_1$A contain the $I_1$-guided multihead attention module. In this module, the features of $I_1$ play a crucial role, whereas those of $I_2$ assist in aligning cross-modal features and enhancing consistency features. In Fig. 3(b), W-MG$_2$A and SW-MG$_2$A involve the $I_2$-guided multihead attention module. Here, the features from $I_2$ dominate, while those from $I_2$ are used for cross-modal feature alignment and enhancement. The $I_1$-guided multihead attention in W-MG$_1$A and SW-MG$_1$A can be described as follows:

$$G_1\text{MHead}(S_1, S_2)$$
$$= \text{Concat}(g_1 head^1, g_1 head^2, \ldots, g_1 head^h) \cdot W^O \quad (11)$$

where $g_1 head^i$ denotes the $i$th $I_1$-guided head in $G_1\text{MHead}(S_1, S_2)$, Concat means the concatenation operation, and $W^O$ is a linear embedding matrix.

Analogously, the $I_2$-guided multihead attention $G_2\text{MHead}(S_2, S_1)$ in W-MG$_2$A and SW-MG$_2$A can be depicted as

$$G_2 M Head(S_2, S_1)$$
$$= \text{Concat}(g_2 head^1, g_2 head^2, \ldots, g_2 head^h) \cdot W^O \quad (12)$$

where $g_2 head^i$ denotes the $i$th $I_2$-guided head in $G_2\text{MHead}(S_2, S_1)$.

Inspired by the cross-attention, we give the specific formulas of guided attention as follows:

$$\begin{cases} g_1 head^i = \text{softmax}\left( \dfrac{K_1^i \cdot (Q_2^i)^\top}{\sqrt{d_k}} + B \right) \cdot V_1^i \\[4mm] g_2 head^i = \text{softmax}\left( \dfrac{K_2^i \cdot (Q_1^i)^\top}{\sqrt{d_k}} + B \right) \cdot V_2^i \end{cases} \quad (13)$$

where $K_1^i$, $V_1^i$, and $Q_1^i$ are different projections of $S_1$. Meanwhile, $K_2^i$, $V_2^i$, and $Q_2^i$ represent different linear projections of $S_2$. If the number of heads $h$ is 1, we will abbreviate $K_1^i$, $V_1^i$, and $Q_1^i$ in (13) as $K_1$, $V_1$, and $Q_1$, respectively. Correspondingly, for $h = 1$, $K_2^i$, $V_2^i$, and $Q_2^i$ in (13) are abbreviated as $K_2$, $V_2$, and $Q_2$, respectively.

From (13) both the $I_1$- and $I_2$- guided attention mechanisms take the association between modalities into consideration. After inserting (13) into (11) and (12), we could get the output of W-MG$_1$A and W-MG$_2$A (i.e., $S_{12}^f$ and $S_{21}^f$) in Fig. 3.

Note that despite the subtle differences between SW-MG$_1$A and SW-MG$_2$A in window partitioning and cross-window attention, both adopt guided attention mechanisms to model long-range dependencies. Therefore, the outputs of SW-MG$_1$A and SW-MG$_2$A, i.e., $S_{12}^{f'}$ and $S_{21}^{f'}$ in Fig. 3, are attainable by substituting (13) into (11) and (12), respectively.

### C. Hierarchical Fusion and Reconstruction

Multimodal image fusion aims to explore and utilize the intermodal useful information to compensate for the deficiency of single-source images and yield high-quality fused images. For example, to fully leverage low-level spatial information and high-level structure information among multiscale features, a bi-directional hierarchical feature collaboration (BHFC) module is given in [53]. Tang et al. [54] gave a hierarchical multimodal fusion architecture to explore multiple bidirectional translation processes, thereby generating dual multimodal fusion embeddings. Motivated by this, we propose to build a hierarchical fusion branch involving a SIF module, a MIF module, and an HRM to fully integrate the multimodality, multigranularity, and multiview features. As depicted in Fig. 2, the former concentrates on the fusion of dual-stream low-level features. Meanwhile, the latter resorts to aggregating dual-stream local and long-range context semantics from various levels.

*1) Shallow Interactive Fusion:* Clearly, aligning and integrating dual-stream information is of great significance for effective multimodal fusion. In this regard, a SIF module and a MIF module are built and used in cascade.

SIF aims to fuse the shallow multiview features from two streams for subsequent processing. Inside the SIF, we first perform concatenation and $3 \times 3$ convolution operations on the fine resolution features from the GIR$_1$ and GIR$_2$ (i.e., $S_f^1$, $S_f^2$). After GIR$_1$ and GIR$_2$, $S_f^1$ and $S_f^2$ is obtained with the size of $(H/4) \times (W/4) \times (2 * d)$, respectively. Note that multihead self-attention embodies automatic focus statistics, which could reveal and exploit the significance of each channel. Therefore, we further employ the multihead self-attention mechanism to adaptively integrate the valuable shallow features from the detail and content streams. The specific calculation process involved in SIF is defined in the following equation:

$$S_f = \text{MHead}(\text{Conv}(\text{Concat}(S_f^1, S_f^2))) \quad (14)$$

where $\text{Concat}(\cdot, \cdot)$ and Conv stand for the concatenation and $3 \times 3$ convolution, respectively. As shown in Fig. 2, $S_f^1$ denotes the output from the front part of GIR$_1$ in the detail stream, $S_f^2$ corresponds to the output from the front part of GIR$_2$ module in the content stream. Then, through SIF, the fused low-level features $S_f$ is captured with the size of $(H/4) \times (W/4) \times (4 * d)$.

*2) Multilevel Interactive Fusion:* For the sake of further aggregating the deep semantics from two streams, but not forgetting shallow features, a MIF module is constructed. MIF is committed to integrating multiview features from SIF and deep semantics from the latter part of GIR$_1$ and GIR$_2$. Specifically, since patch merging has the advantage over pooling in terms of information preservation, MIF first utilizes patch merging to obtain lower-scale features from the SIF module (i.e., $S_f$). Then, to absorb multilevel local and long-range contextual information from two modalities (i.e., $M_f^1$, $M_f^2$), multimodal features of diverse scales are concatenated followed by $1 \times 1$ convolution for feature aggregation. Finally, similar to SIF, the multihead self-attention mechanism is applied for attentive feature fusion. The expression for MIF is given by

$$M_f = \text{MHead}\Big(\text{Conv}\Big(\text{Concat}\Big(\text{PM}(S_f), M_f^1, M_f^2\Big)\Big)\Big) \quad (15)$$

where $S_f$ is the output of SIF. PM represents the operation of patch merging. Due to the operation of patch merging, the spatial size of $S_f$ is reduced to $(H/8) \times (W/8)$, which is the same as that of $M_f^1$ and $M_f^2$. Then, based on (15), we could capture the multilevel semantic features $M_f \in \mathbb{R}^{(H/8) \times (W/8) \times (8*d)}$.

From (15), it is obvious that MIF not only fuses multiscale and multimodal features but also achieves multiview fusion via multihead attention.

*3) High-Level Interactive Fusion and Reconstruction:* To generate high-quality fused images, the core problem is to comprehensively explore and aggregate compatible and credible complementary features. To this end, we further construct a HRM branch. This part focuses mainly on hierarchical feature fusion through Swin transformer blocks. The main reason is that the Swin transformer has excellent long-range modeling capability. Particularly, multihead self-attention can obtain different perceptions from cross-modal hierarchical features with multiscales and levels. This is beneficial for comprehensive and multiview analysis of explored features.

As manifested in Fig. 2, the concatenation and $3 \times 3$ convolution are utilized for feature aggregation from two streams and the fusion branch. Specifically, the aggregation block takes multilevel $M_f$ and high-level $M_f^1$, $M_f^2$ semantic information from dual-stream as inputs. Then, the aggregated features are injected into $3\times$ Swin transformer blocks demonstrated as follows:

$$D_f = \mathrm{ST}_{3\times}\big(\mathrm{Conv}\big(\mathrm{Concat}\big(M_f, M_f^1, M_f^2\big)\big)\big) \qquad (16)$$

where $\mathrm{ST}_{3\times}$ represents three consecutive Swin transformer blocks. Here, Conv means $3 \times 3$ convolution. Through (16), we could acquire the high-level aggregated semantic features $D_f \in \mathbb{R}^{(H/8)\times(W/8)\times(4*d)}$. There are two reasons for using convolution in aggregation blocks. The first is to employ their local perception ability to perceive detailed information further. The second is to project the multimodal and multilevel features into one shared space for semantic alignment.

Regarding feature decoding, the patch expansion operation is first leveraged, which first restores the size of the feature map from $(H/8) \times (W/8)$ to $(H/4) \times (W/4)$. Then, enlightened by residual connection, the second convolution block combines low-level information $S_f^1$, $S_f^2$ from two streams with the expanded $D_f$. After that, the aggregated features are decoded based on four consecutive Swin transformer blocks, which guarantees that multimodal features of different scales, levels, and views can be considered. The specific construction process of $R_f \in \mathbb{R}^{(H/4)\times(W/4)\times(2*d)}$ can be described as

$$R_f = \mathrm{ST}_{4\times}\big(\mathrm{Conv}\big(\mathrm{Concat}\big(\mathrm{PE}(D_f), S_f^1, S_f^2\big)\big)\big) \qquad (17)$$

where PE denotes the patch expansion operation, which could increase the size of the feature maps.

Then, through size expansion and channel reduction, we can reconstruct the $H \times W \times 1$ fused image $I_f$ as

$$I_f = \mathrm{Conv}\big(\mathrm{PE}\big(R_f\big)\big) \qquad (18)$$

where Conv is the $1 \times 1$ convolution for channel reduction.

## D. Discriminators

For generator optimization in IG-GAN, it is sensible to introduce discriminators for adversarial learning. Specifically, we give two discriminators for the mutual game with the generator. They are designed to classify the generated image and misclassify the original multimodal images. For example.
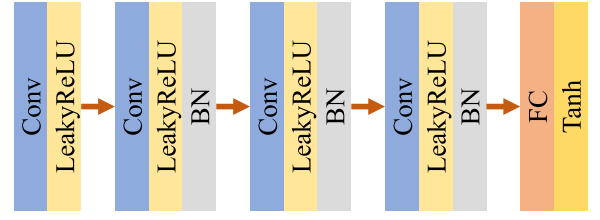


Fig. 4. Architecture of discriminator.

The first one can identify that the first mode image is true, while the fused image is false. The second discriminator is used to discriminate the fused image from the second mode image. Concretely, it can recognize the first modal image as true but the fused image as false.

Concerning structural symmetry, the two discriminators have the same structure but do not share parameters. As described in Fig. 4, each discriminator contains four convolutional modules, which involve convolution (Conv), LeakyReLU, and batch normalization (BN) layers. Finally, the fully connected (FC) layer is used, followed by the Tanh activation function.

## E. Loss Function

The loss function is of crucial importance to guide the network optimization and boost the mutual game between the generator and discriminators. To promote the training of IG-GAN without supervision, a comprehensive loss function respecting the generator and discriminators is given by

$$\mathcal{L} = \mathcal{L}_G + \mathcal{L}_{\mathrm{Dis1}} + \mathcal{L}_{\mathrm{Dis2}} \qquad (19)$$

where $\mathcal{L}_G$ reflects the generator loss, and $\mathcal{L}_{\mathrm{Dis1}}$ and $\mathcal{L}_{\mathrm{Dis2}}$ correspond to the losses of discriminator$_1$ and discriminator$_2$, respectively.

*1) Generator Loss:* The core of multimodal fusion is to use the intermodality complementary information to enrich texture details and maintain content integrity. In this regard, the generator loss is given by

$$\mathcal{L}_G = \lambda_1 \mathcal{L}_{\mathrm{con}} + \lambda_2 \mathcal{L}_{\mathrm{ei}} + \mathcal{L}_{\mathrm{com}} \qquad (20)$$

where $\lambda_1$ and $\lambda_2$ are balancing parameters and $\mathcal{L}_{\mathrm{con}}$, $L_{\mathrm{ei}}$, and $\mathcal{L}_{\mathrm{com}}$ represent the content loss, edge intensity (EI) loss, and compatibility loss, respectively.

*a) Content loss:* The content loss $\mathcal{L}_{\mathrm{con}}$ aims to retain the completeness of spatial structures. Note that the structural similarity index (SSIM) is one of the most widely used indices for image fusion, which reflects the structural similarity between images from three perspectives (i.e., light intensity, contrast, and structure). Therefore, we adopt SSIM to manifest the content completeness of the fused image. In addition, as the mean square error (mse) is an intuitive index for measuring the discrepancy between modalities pixel-by-pixel, it is employed for constructing the content loss. The specific content loss involving structural similarity loss $\mathcal{L}_{\mathrm{ssim}}$ and mse loss $\mathcal{L}_{\mathrm{mse}}$ is

$$\mathcal{L}_{\mathrm{con}} = a\mathcal{L}_{\mathrm{ssim}} + b\mathcal{L}_{\mathrm{mse}} \qquad (21)$$

where $a$ and $b$ are harmonic coefficients, which are set to 5 and 1 in the experiment, respectively.

$\mathcal{L}_{\mathrm{ssim}}$ is defined in (6). Specifically, considering the need to treat each modality equally, in (6), both $w_1$ and $w2$ are set to 1. The following equation gives the definition of $\mathcal{L}_{\mathrm{mse}}$:

$$\mathcal{L}_{\mathrm{mse}} = w_3 \| \boldsymbol{I_f} - \boldsymbol{I_1} \|_2 + w_4 \| \boldsymbol{I_f} - \boldsymbol{I_2} \|_2 \qquad (22)$$

where $\| \cdot \|_2$ represents the $l_2$-norm and $w_3$ and $w_4$ are tradeoff parameters, which are set to 1 in this article.

*b) EI loss:* It is well known that the richer the texture, the larger the total gradient of the image. Therefore, to enhance the texture details of the fused image, the EI loss is expressed as

$$\mathcal{L}_{\mathrm{ei}} = 1 - \left( \left( \nabla_x \boldsymbol{I_f} \right)^2 + \left( \nabla_y \boldsymbol{I_f} \right)^2 \right)^{\frac{1}{2}} \qquad (23)$$

where $\nabla$ denotes the gradient operator and $x$ and $y$ represent the directions of the derivative. The above three losses form our generator loss function, which aims to maximize the probability that the fused image comes from both modalities.

In contrast, the discriminator adopts a multiclass classifier to determine that the fused image is neither an infrared image nor a visible image. Then, with continuous adversarial learning, the generator can estimate the probability distribution of both the infrared and visible images, which could generate a fused image with significant contrast and rich texture details.

*c) Compatibility loss:* The fused image is the complementary fusion result of the multimodal input images. This means that the fused image should be compatible with the input images. Assume the input image is true. Then, the generator should try to maximize the probability that the fused image is still true. The following equation gives the specific formulation of compatibility loss:

$$\mathcal{L}_{\mathrm{com}} = \mathbb{E}\left[\log\left(1 - D_1(\boldsymbol{I_f})\right)\right] + \mathbb{E}\left[\log\left(1 - D_2(\boldsymbol{I_f})\right)\right] \qquad (24)$$

where $\mathbb{E}$ means expectation, $\boldsymbol{I_f}$ represents the fused image, and $D_1$ and $D_2$ stand for the discriminator$_1$ and discriminator$_2$, respectively.

*2) Discriminator Loss:* The primary task of discriminators is to improve the performance of the generator through game confrontation. Note that the compatibility loss helps the generated fused image to be compatible or in fidelity with the input multimodal images. Therefore, the discriminator should recognize the incompatibilities and judge the fused image as false. The following equation gives the discriminator loss $\mathcal{L}_{\mathrm{Dis}}$:

$$\mathcal{L}_{\mathrm{Dis}} = \mathbb{E}\left[-\log D(\boldsymbol{I_s})\right] + \mathbb{E}\left[-\log\left(1 - D(\boldsymbol{I_f})\right)\right] \qquad (25)$$

where $D$ stands for discriminator, $\boldsymbol{I_s}$ represents for the source image, and $\boldsymbol{I_f}$ denotes the fused image generated by the generator. The loss function described above is used for both source images.

## IV. EXPERIMENTS

This section first describes the experimental settings, e.g., four commonly used multimodal datasets, comparative methods, evaluation metrics, and parameters setting. Second, quantitative and qualitative comparisons with several state-of-the-art methods are provided. Finally, the effectiveness of the specific design is assessed by ablation studies.

### A. Experimental Settings

*1) Dataset:* The OS dataset consists of fine-resolution optical and SAR images [55]. The optical images were collected via the Google Earth platform, whereas SAR images were captured by the Chinese C-band sensor Gaofen-3, in spotlight mode. The dataset contains 10 692 image pairs with size $256 \times 256$. The TNO dataset is an infrared and visible image dataset provided by TNO in The Netherlands. These image pairs include various scenes with size $360 \times 270$, $505 \times 510$, and $768 \times 576$ [56]. The RGB-NIR Scene dataset includes 477 images captured in RGB and near-infrared (NIR) [57]. There are nine scenes in total: country, field, forest, indoor, mountain, old building, street, urban, and water.

*2) Comparative Methods:* To evaluate the effectiveness of our proposed method, 14 state-of-the-art fusion methods were employed for performance comparison, including DDc-GAN [37], FusionGAN [36], GAN-FM [58], GANMcC [12], PMGI [29], RFNNest [28], SDDGAN [32], STDFusion-Net [59], U2Fusion [27], SwinFusion [46], TGFuse [45], DetailGAN [50], ATFuse [60], and PSFusion [61]. Among the 14 methods, DDcGAN, FusionGAN, GAN-FM, GANMCC, SDDGAN, TGfuse, and DetailGAN all achieve multimodal image fusion through the architecture of GAN.

*3) Evaluation Metrics:* Concerning quantitative comparison, seven typical evaluation criteria were adopted. EI, spatial frequencies (SFs), and average gradient (AG) mainly attempt to reflect the EI, SFs, and AG of a fused image, respectively. Meanwhile, The sum of the correlations of differences (SCD) and the correlation coefficient (CC) are self-explanatory [62]. Both visual information fidelity (VIF) and VIF for fusion (VIFF) were introduced to characterize the VIF of the images [63], [64].

*4) Settings:* All the experiments were performed on an NVIDIA GeForce GTX 3090 GPU with batch size 24. In the first 25 epochs of the training process, the learning rate was set as 0.002 without discriminators. After that, the learning rate was set to 0.0001. Before feeding the image pairs into our IG-GAN, all images are resized to $256 \times 256$ in advance. Moreover, all comparison methods are set with reference to their authors.

### B. Results on the OS Dataset

*1) Quantitative Comparison:* To evaluate the effectiveness of IG-GAN, this section mainly focuses on a quantitative comparison of the OS dataset. Table I provides the quantitative comparison with 11 popular methods in terms of EI, SF, SCD, VIF, AG, CC, and VIFF on the OS dataset. The optimal and suboptimal results are marked in bold and underlined font, respectively.

From Table I we can observe that GANMcC and DetailGAN exhibit superior results respecting CC. This demonstrates the effectiveness of their loss function, which requires the fused image to be consistent with the original inputs. Except for CC, our IG-GAN has obvious advantages over the others in terms of EI, SF, SCD, VIF, AG, and VIFF. For example, our EI is 165.012, which is almost 12.3 higher than the next-best comparator and 116 higher than the lowest one. Regarding

SF, we achieved 1.24 advantages over the second-placed and 37.95 advantages over the worst-placed. In terms of AG, our method is greater than the third-place TGfuse by nearly 3.9. Concerning VIFF, we achieved 0.765, whereas the worst is 0.198, which is over 0.56 lower. As EI, SF, SCD, VIF, AG, and VIFF characterize the detail texture, integrity, and fidelity of fused images from different views, the obvious advantage in these criteria shows the superiority of our fusion method regarding detail richness and content integrity.

*2) Qualitative Comparison:* For illustrative purposes, Fig. 5 shows an example of the fused images generated by each method. The first two images in this figure correspond to the original optical (OPT) and SAR images, respectively. For better display, an enlarged view of the object area from the red line frame is presented.

From Fig. 5, we can see that balancing the completeness and clarity of the objects is a challenge for most fusion methods. For example, respecting DDcGAN, FusionGAN, PMGI, RFNNest, DetailGAN, and PSFusion, the integrity of the objects is acceptable. In comparison, the fused images are blurred and lack precise edge contours. By contrast, GANMcC, SwinFusion, U2Fusion, and ATFuse yield clearer outlines of the objects. Unfortunately, the completeness and contrast are not satisfactory. In view of GAN-FM, STD-fusionNet, TGFuse, and PSFusion, their fused images are superior. While in terms of the brightness and saliency of the objects, they are still some way behind the results of IG-GAN. Apparently, our IG-GAN can give complete objects, clear contours, and obvious contrast.

### C. Results on the TNO Dataset

*1) Quantitative Comparison:* This section provides the comparative experiments corresponding to different methods conducted on the TNO dataset. Table II shows the specific experimental results respecting seven criteria. Bold and underlined results represent the optimal and suboptimal results, respectively.

As illustrated in Table II, ATFuse performs excellently on the TNO dataset. It can acquire suboptimal fusion results for EI, SF, and AG. This means that the fused images are rich in structural and detailed information. Note that DDcGAN and GANMcC achieve the suboptimal VIF and CC, respectively. This demonstrates that they are good at preserving the original input information. Regarding IG-GAN, instead of performing well in some single aspects, it is superior to others on the whole. For example, in terms of EI, SF, SCD, VIF, AG, and VIFF, it always performs the best. Even for CC, it is also the third-best only after DetailGAN and GANMcC. This reveals that IG-GAN is advantageous to producing a high-quality fused image involving great texture, complete content, and high-fidelity vision.

*2) Qualitative Comparison:* To be more intuitive, Fig. 6 further depicts an example of the original images and the corresponding fused images generated by the 14 methods. The first two images in this figure correspond to the infrared (IR) and visible (VIS) images, respectively. For better display, an enlarged view of the object area from the red line frame is provided

As depicted in Fig. 6, constrained by illumination, only weak objects are challenging to recognize in the visible images. In contrast, the infrared image can present the objects with deficient details. In this scenario, image fusion should prioritize the infrared image over the visible light image. Otherwise, as observed in DDcGAN and FusionGAN, the visible light image significantly influences the fused image. Consequently, the objects in the fused image become highly blurred, as indicated in Table II. As demonstrated in Fig. 7, IG-GAN tends to outperform others for all test images. In alignment with Fig. 7, our IG-GAN can produce high-quality fusion images with complete content, a clear outline, and strong contrast.

### D. Results on the RGB-NIR Scene Dataset

*1) Quantitative Comparison:* To evaluate the effectiveness of IG-GAN, this section further provides a quantitative comparison conducted on the RGB-NIR Scene dataset. Table III shows the corresponding results respecting EI, SF, SCD, DF, AG, and VIFF. Bold and underlined represent the optimal and suboptimal results, respectively.

From Table III, we can observe that IG-GAN still has excellent performance. Despite not being the best for all criteria, it can achieve the top two for these criteria. For example, in terms of EI, SCD, DF, and VIFF, it is suboptimal and quite close to the optimal. Regarding SF, our method is superior to the others. For AG, though IG-GAN is the third-best place, it is closely different from the optimal 8.911 and suboptimal 8.856 by 0.08 and 0.025, respectively. The outstanding performance on four popular datasets demonstrates that IG-GAN achieves the preservation of intermodal consistency and the comprehensive fusion of complementary information.

*2) Qualitative Comparison:* To be more intuitive, Fig. 8 further gives an example of the original images and the corresponding fused images generated by 14 methods. In this figure, the first two images correspond to the NIR and visible (VIS) images, respectively. For better display, we also provide an enlarged view of the red object area.

As shown in Fig. 8, some methods integrate the texture and brightness from VIS and NIR well. While respecting edge sharpness and visual contrast, our IG-GAN still outperforms the others. Specifically, the tree is clearly described with branches and leaves that stand out and have the highest contrast and structural information. This reflects its unique advantages in detail and content exploration and cross-modal concordance enhancement.

### E. Computational Complexity and Efficiency

Note that model parameter quantity (Params), floating-point operations per second (FLOPs), and inference time are three important metrics for measuring the complexity and efficiency of deep learning algorithms. Therefore, to evaluate the complexity and efficiency of our IG-GAN, we further give the comparison between IG-GAN and four typical multimodal fusion methods regarding params, FLOPs, and Inference time. When the input image size is $256 \times 256$, Table IV presents a comparative analysis of model parameters and computational

TABLE I

QUANTITATIVE COMPARISON WITH STATE-OF-THE-ARTS ON THE OS DATASET

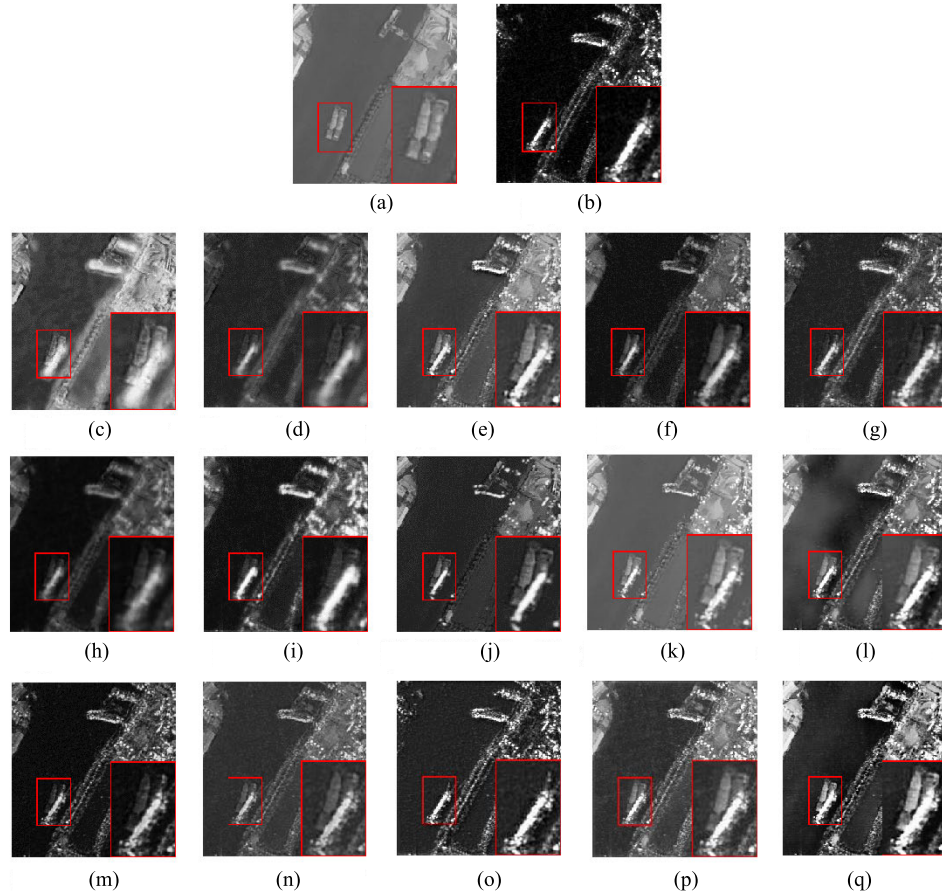| | EI | SF | SCD | VIF | AG | CC | VIFF |
|---|---|---|---|---|---|---|---|
| DDcGan [37] | 92.974 | 20.417 | 1.095 | _1.261_ | 9.041 | 0.75 | 0.511 |
| FusionGAN [36] | 48.579 | 10.3 | 0.813 | 0.459 | 4.611 | 0.747 | 0.233 |
| GANFM [58] | 126.069 | 32.792 | 1.518 | 1.014 | 12.896 | 0.78 | 0.533 |
| GANMcC [12] | 65.644 | 16.01 | 1.348 | 0.602 | 6.659 | _0.798_ | 0.405 |
| PMGI [29] | 105.072 | 26.639 | 1.379 | 0.752 | 10.621 | 0.783 | 0.485 |
| RFNNest [28] | 57.028 | 11.094 | 1.393 | 0.706 | 5.325 | 0.784 | 0.444 |
| SDDGAN [32] | 95.341 | 21.002 | 1.452 | 1.144 | 9.147 | 0.763 | _0.682_ |
| STDFusionNet [59] | 67.365 | 16.944 | 1.073 | 0.633 | 6.548 | 0.729 | 0.198 |
| SwinFusion [46] | 91.094 | 22.927 | 1.479 | 0.739 | 9.186 | 0.767 | 0.341 |
| TGFuse [45] | 130.072 | 37.135 | 1.507 | 0.953 | 13.687 | 0.762 | 0.542 |
| U2Fusion [27] | 123.799 | 31.13 | 1.496 | 0.88 | 12.515 | 0.783 | 0.512 |
| DetailGAN [50] | 78.123 | 16.766 | 1.386 | 0.536 | 7.557 | **0.8** | 0.422 |
| ATFuse [60] | _152.741_ | _47.01_ | 0.833 | 0.924 | _16.752_ | 0.729 | 0.372 |
| PSFusion [61] | 122.293 | 32.817 | _1.574_ | 0.93 | 12.918 | 0.782 | 0.557 |
| IG-GAN | **165.012** | **48.251** | **1.656** | **2.001** | **17.564** | 0.784 | **0.765** |



Fig. 5. Qualitative comparison on the OS dataset. (a) OPT. (b) SAR. (c) DDcGAN. (d) FusionGAN. (e) GAN-FM. (f) GANMcC. (g) PMGI. (h) RFNNest. (i) SDDGAN. (j) STDFusionNet. (k) SwinFusion. (l) TGFuse. (m) U2Fusion. (n) DetailGAN. (o) ATFuse. (p) PSFusion. (q) IG-GAN.

complexity between SwinFusion, TGFuse, ATFuse, PSFusion, and IG-GAN in terms of Params, FLOPs, and inference time.

There are three reasons for taking SwinFusion, TGFuse, ATFuse, and PSFusion for complexity and efficiency comparison. The first reason is based on the architectural similarities between SwinFusion and IG-GAN, both of which are hybrid fusion networks combining CNN and transformer. In this regard, it is necessary to compare the complexity and efficiency between SwinFusion and IG-GAN. The second reason is that, as given in Tables I–III, TGFuse is an effective GAN-based fusion method, which also involves the transformer and CNN. Hence, besides the effectiveness, we further provide the complexity and efficiency comparison between

TABLE II

QUANTITATIVE COMPARISON WITH SEVERAL STATE-OF-THE-ART METHODS ON THE TNO DATASET

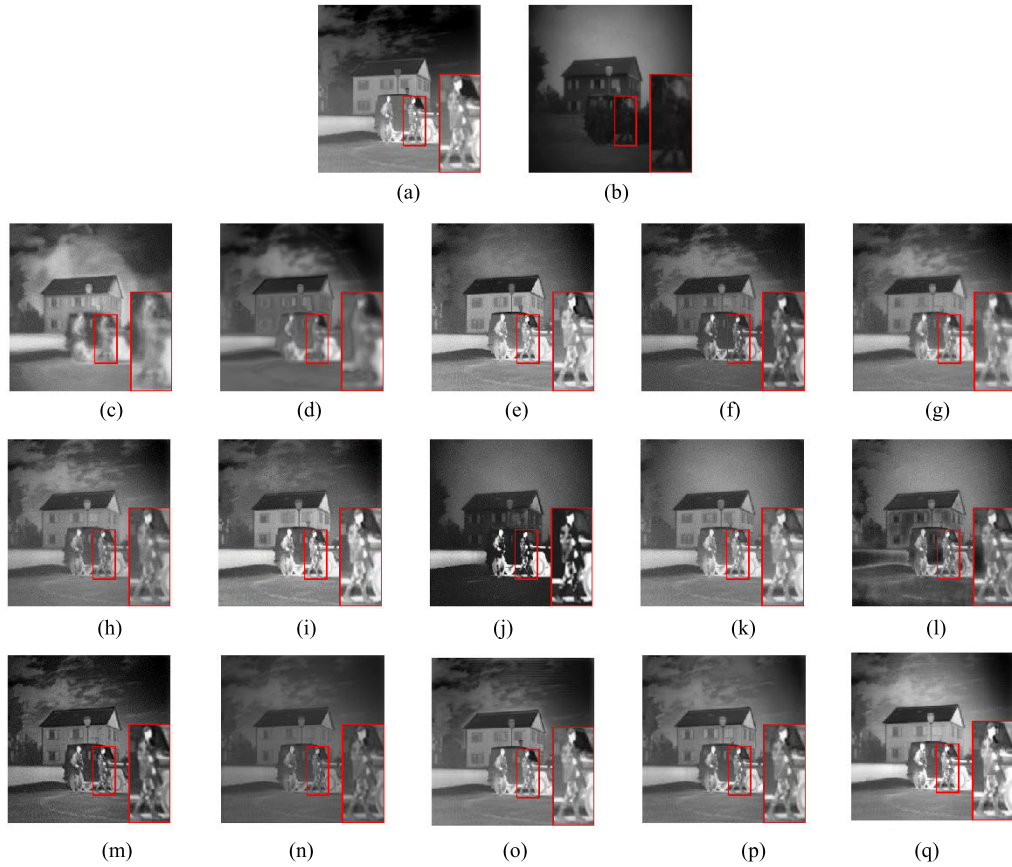| | EI | SF | SCD | VIF | AG | CC | VIFF |
|---|---|---|---|---|---|---|---|
| DDcGAN [37] | 44.486 | 11.193 | 1.479 | **1.191** | 4.550 | 0.702 | 0.612 |
| FusionGAN [36] | 24.142 | 6.240 | 1.037 | 0.524 | 2.417 | **0.727** | 0.258 |
| GAN-FM [58] | 46.763 | 12.526 | 1.545 | 1.038 | 4.846 | 0.682 | 0.494 |
| GANMcC [12] | 25.895 | 6.139 | 1.347 | 0.615 | 2.546 | 0.683 | 0.422 |
| PMGI [29] | 36.832 | 8.749 | 1.528 | 0.885 | 3.606 | 0.704 | 0.543 |
| RFNNest [28] | 28.644 | 5.873 | 1.568 | 0.729 | 2.682 | 0.701 | 0.513 |
| SDDGAN [32] | 36.799 | 8.991 | 1.556 | 1.159 | 3.608 | 0.676 | 0.658 |
| STDFusionNet [59] | 44.111 | 11.807 | 1.361 | <u>1.188</u> | 4.456 | 0.659 | 0.473 |
| SwinFusion [46] | 41.237 | 10.639 | 1.71 | 0.757 | 4.142 | 0.681 | 0.472 |
| TGFuse [45] | 41.748 | 11.044 | 1.498 | 0.857 | 4.232 | 0.670 | 0.482 |
| U2Fusion [27] | 51.445 | 11.861 | 1.607 | 1.083 | 5.060 | 0.700 | <u>0.683</u> |
| DetailGAN [50] | 22.994 | 5.15 | 1.567 | 0.386 | 2.134 | <u>0.718</u> | 0.3 |
| ATFuse [60] | 50.587 | 12.512 | 1.5516 | 0.747 | 4.885 | 0.668 | 0.36 |
| PSFusion [61] | <u>53.311</u> | <u>12.926</u> | <u>1.613</u> | 0.85 | **7.173** | 0.687 | 0.494 |
| IG-GAN | **55.225** | **13.903** | **1.697** | 1.177 | <u>5.418</u> | 0.704 | **0.717** |



Fig. 6. Quantitative comparison on the TNO dataset. (a) IR. (b) VIS. (c) DDcGAN. (d) FusionGAN. (e) GAN-FM. (f) GANMcC. (g) PMGI. (h) RFNNest. (i) SDDGAN. (j) STDFusionNet. (k) SwinFusion. (l) TGFuse. (m) U2Fusion. (n) DetailGAN. (o) ATFuse. (p) PSFusion. (q) IG-GAN.

IG-GAN and TGFuse. The third reason is that ATFuse and PSFusion are highly competitive fusion methods proposed in the past two years. Regarding this, besides the quality of the fused image, we should further compare the complexity and complexity of IG-GAN with ATFuse and PSFusion. Table IV presents a comparative analysis of model parameters and computational complexity between SwinFusion, TGFuse,

ATFuse, PSFusion, and IG-GAN in terms of Params, FLOPs, and inference time.

According to Table IV, although the number of model parameters for IG-GAN is not the least, it exhibits low model complexity and high fusion efficiency. In specific, IG-GAN has the second lowest number of model parameters after ATFuse. The PSFusion has a significantly higher model size and
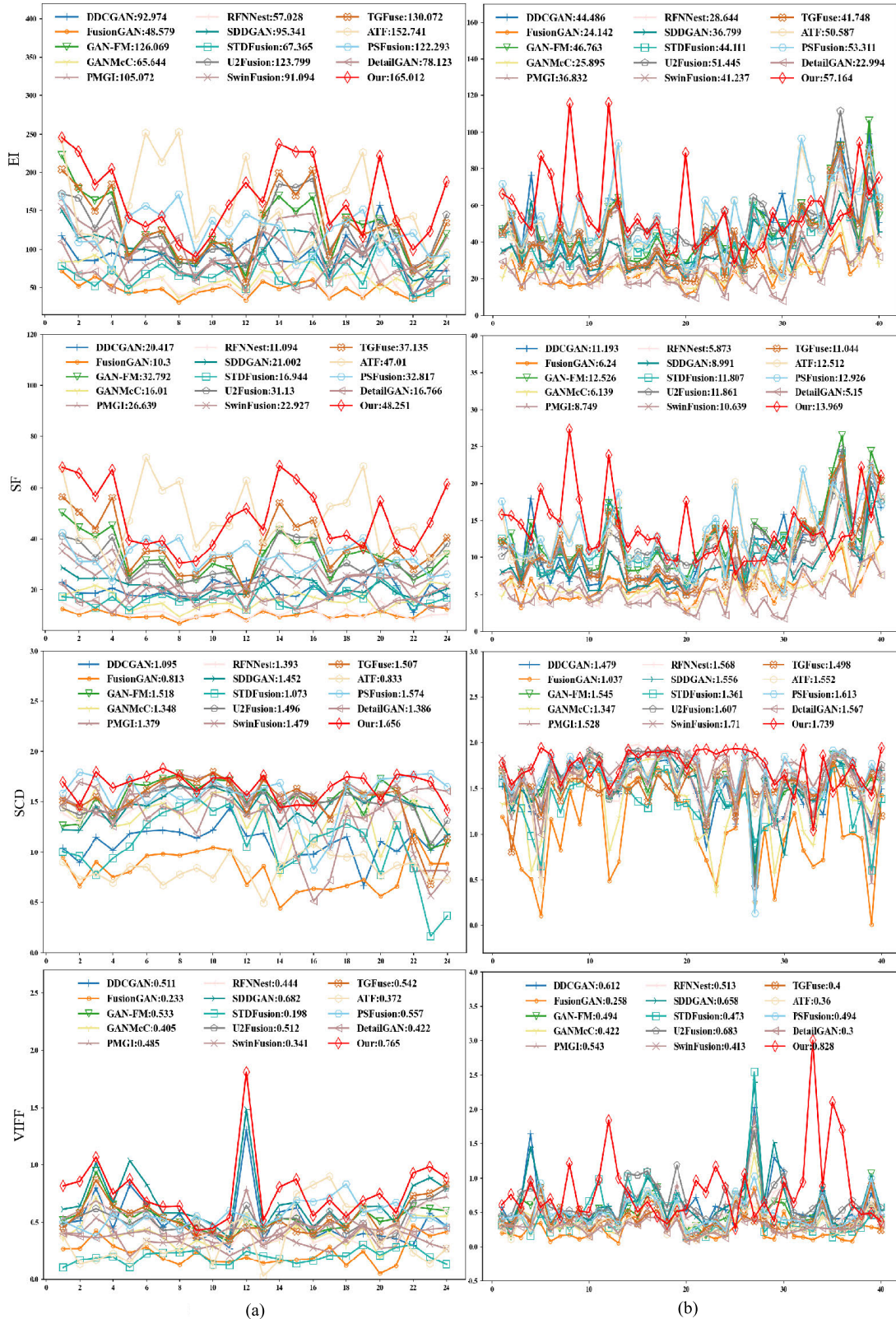
Fig. 7. Performance comparison curves with 14 deep fusion methods on the three datasets. (a) OS dataset. (b) TNO dataset.

computational complexity compared to IG-GAN. Specifically, PSFusion's model size is almost 13 times larger than IG-GAN, leading to a corresponding increase in FLOPs, reaching 1.234T compared to IG-GAN's 15.904G. This highlights the tradeoff

between fusion performance and computational complexity in PSFusion. This reveals the tradeoff between the fusion performance and computational complexity in PSFusion. Compared to SwinFusion, the model size of IG-GAN is reduced to
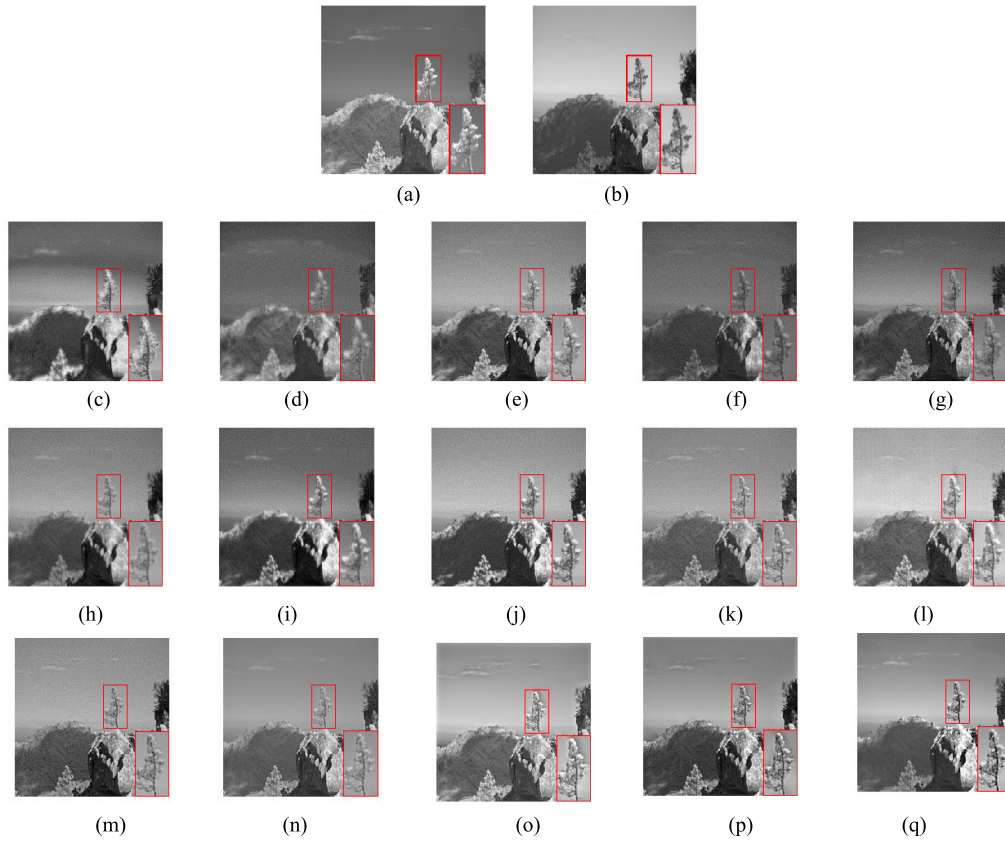
Fig. 8. Qualitative comparison on the RGB-NIR dataset. (a) NIR. (b) VIS. (c) DDcGAN. (d) FusionGAN. (e) GAN-FM. (f) GANMcC. (g) PMGI. (h) RFNNest. (i) SDDGAN. (j) STDFusionNet. (k) SwinFusion. (l) TGFuse. (m) U2Fusion. (n) DetailGAN. (o) ATFuse. (p) PSFusion. (q) IG-GAN.

TABLE III
QUANTITATIVE COMPARISON ON THE RGB-NIR SCENE DATASET

| | EI | SF | SCD | DF | AG | VIFF |
|---|---|---|---|---|---|---|
| DDcGAN [37] | 70.778 | 19.593 | 0.709 | 9.611 | 7.236 | 0.700 |
| FusionGAN [36] | 46.960 | 12.275 | 0.330 | 5.978 | 4.683 | 0.427 |
| GAN-FM [58] | 81.528 | 23.374 | 1.020 | 11.347 | 8.417 | 0.708 |
| GANMcC [12] | 52.968 | 13.981 | 0.697 | 6.773 | 5.310 | 0.519 |
| PMGI [29] | 67.805 | 17.491 | 0.828 | 8.643 | 6.736 | 0.629 |
| RFNNest [28] | 45.661 | 9.834 | 1.158 | 4.832 | 4.294 | 0.582 |
| SDDGAN [32] | 54.594 | 13.389 | 0.925 | 6.678 | 5.319 | 0.711 |
| STDFusionNet [59] | 78.242 | 19.202 | **1.358** | 9.050 | 7.575 | **0.879** |
| SwinFusion [46] | 72.477 | 19.752 | 1.140 | 9.712 | 7.433 | 0.650 |
| TGFuse [45] | 72.904 | 20.300 | 1.114 | 9.657 | 7.421 | 0.685 |
| U2Fusion [27] | **89.776** | 22.339 | 1.237 | 11.142 | **8.911** | 0.772 |
| DetailGAN [50] | 55.34 | 12.524 | 1.098 | 6.118 | 5.302 | 0.567 |
| ATFuse [60] | 83.192 | 22.385 | 1.086 | 11.283 | 8.575 | 0.591 |
| PSFusion [61] | 85.687 | <u>23.241</u> | 1.203 | **11.676** | <u>8.856</u> | 0.728 |
| IG-GAN | <u>86.445</u> | **23.941** | <u>1.318</u> | <u>11.608</u> | 8.831 | <u>0.863</u> |

over 75%, which leads to a 0.357M decrease in FLOPs. In addition, IG-GAN's inference time is less than a quarter of SwinFusion's. This demonstrates that IG-GAN has a lower model complexity and higher fusion efficiency compared to SwinFusion. Concerning TGFuse, it is a generative adversarial fusion network consisting of a spatial transformer and a channel transformer. Despite that TGFuse has less inference time than IG-GAN, it relies on high model complexity, with model parameters exceeding 155 times that of IG-GAN. Therefore, Table IV indicates that although IG-GAN is not the lightest fusion method, it has a relatively small model complexity and high algorithm efficiency.

### F. Ablation Studies

The excellent performance of IG-GAN relies on our well-designed structure and loss function. To this end,

TABLE IV
MODEL PARAMETER COMPARISON

| Model | Input | Params | FLOPs | Inference Time |
|---|---|---|---|---|
| SwinFusion | $(256 \times 256)$ | 3.895M | 63.731G | 1.035s |
| TGFuse | $(256 \times 256)$ | 549.359M | 15.945G | 0.229s |
| ATFuse | $(256 \times 256)$ | 262.939K | 5.405G | 0.213s |
| PSFusion | $(256 \times 256)$ | 45.899M | 1.234T | 1.031s |
| IG-GAN | $(256 \times 256)$ | 3.538M | 15.904G | 0.246s |

an experimental ablation study of our proposed SIF module and MIF module is first analyzed. Then, we further give an ablation analysis concerning our loss function with respect to $\mathcal{L}_{]\rangle}$ on the OS dataset.

*1) Ablation Study of Shallow and MIF Modules:* The ablation experiments were carried out under four situations as shown in Table V. For clarity, we use the blue font next to the up arrow to indicate the increment.

From Table V, we find that the fusion performance is slightly improved after employing SIF alone. When MIF alone is introduced, there is greater improvement in terms of all criteria. While for IG-GAN involving both SIF and MIF, the best performance is attained for most criteria. The main reason for this is that they are conducive to the multigrained, multilevel, and multiview integration of dual-stream features.

*2) Ablation Study of the EI Loss Function:* To assess the importance of EI loss $\mathcal{L}_{ei}$, comparative experiments without $\mathcal{L}_{ei}$ were conducted in the same settings as our original IG-GAN. Experimental results are depicted in Table VI. For clarity, we adopt the blue font next to the up arrow to indicate the increment.

As shown in Table VI, despite there being a slight and negligible decline for VIF and CC, the employment of $\mathcal{L}_{ei}$ exhibits an evident performance advantage for the other criteria.

### G. Application to Object Detection

To explore the application potential of IG-GAN to multimodal object detection, experiments on two public RGB-IR object detection datasets: 1) DroneVehicle dataset [65] and 2) FLIR dataset [66], are conducted.

*1) Dataset: a) DroneVehicle:* The DroneVehicle dataset is a large-scale vehicle detection dataset based on UAV aerial photography, containing both visible and infrared modalities. The dataset covers a full range of lighting environments and has many different occlusion information with variations in image scale and shooting angle. The dataset contains five categories, i.e., "car," "van," "bus," "truck," and "freight car." In this dataset, 17 990 pairs of images are used for training, whereas 8980 pairs of images are used for testing. In addition, all the images are resized to 640 × 640.

*b) FLIR:* The aligned FLIR dataset [66] is an autopilot dataset taken on city streets and highways and contains both daytime and nighttime lighting conditions. This dataset consists of three categories. In our experiments, we leverage 4113 pairs of images for training, whereas 515 pairs of images are used for detection. Similar to the DroneVehicle dataset, all the images are resized to 640 × 640.

*2) Metrics:* Mean average precision (mAP) denotes the mean of average precision (AP) across all categories, which is given by

$$mAP = \frac{1}{N} \sum_{i=1}^{N} AP_i \tag{26}$$

where $AP_i$ represents the AP of the $i$th category.

Concerning the positional accuracy of the detection box, mAP50, mAP75, and mAP50:95 are widely used, which reflect the mAP values under different intersection over union (IOU) thresholds. Concretely, mAP50 and mAP75 provide the mAP when the IOU threshold is 0.5 and 0.75, respectively. Respecting mAP50:95, it is the average of the mAP values at IOU thresholds ranging from 0.50 to 0.95 in steps of 0.05.

*3) Experimental Setting and Results:* To evaluate the effectiveness of IG-GAN for enhancing object detection performance, experiments are conducted on the DroneVehicle dataset and the FLIR dataset.

In specific, IG-GAN is first applied to fuse the visible and infrared image pairs and generate the fused images. Then, the fused images are fed into the detection model for training and testing. Note that YOLOv5 [67] is a popular lightweight object detection model owing to its efficiency and excellent detection performance. Regarding this, YOLOv5 (https://github.com/ultralytics/yolov5) is trained and tested based on the single modal images and the multimodal fused images through IG-GAN, respectively.

Table VII lists the detection accuracy of images before and after IG-GAN fusion respecting mAP50, mAP75, and mAP50:95. From Table VII, we can find that the IG-GAN + YOLOv5 method, significantly outperforms the YOLOv5 method on both the DroneVehicle and FLIR datasets across three metrics (mAP50, mAP75, and mAP90). For example, on the DroneVehicle dataset, the mAP50, mAP75, and mAP90 are 0.757, 0.51, and 0.467, respectively. After IG-GAN, the corresponding detection accuracies are improved by 0.71, 0.211, and 0.131, respectively. On the FLIR dataset, the integration of IG-GAN with YOLOv5 also boosts performance across multiple metrics. Fig. 9 depicts the four pairs of detection results on the DroneVehicle dataset with and without IG-GAN fusion.

TABLE V
ABLATION STUDY ON SHALLOW AND MIF MODULES

| SIF | MIF | EI | SF | SCD | VIF | AG | CC | VIFF |
|---|---|---|---|---|---|---|---|---|
| ✗ | ✗ | 147.678 | 43.010 | 1.633 | 1.682 | 15.670 | 0.784 | 0.696 |
| ✗ | ✔ | 157.242 | 45.263 | 1.655 | 1.894 | 16.684 | **0.788** (↑ 0.004) | **0.788** (↑ 0.092) |
| ✔ | ✗ | 147.781 | 43.033 | 1.610 | 1.694 | 15.599 | 0.776 | 0.696 |
| ✔ | ✔ | **165.012** (↑ **17.334**) | **48.251** (↑ **5.241**) | **1.656** (↑ **0.023**) | **2.001** (↑ **0.319**) | **17.564** (↑ **5.241**) | 0.784 | 0.765 (↑ **0.069**) |



Fig. 9. Visualization of groundtruth and detection results respecting YOLOv5, and combination of IG-GAN and YOLOv5. The red boxes indicate that the objects are missed detection.

TABLE VI
ABLATION STUDY OF THE EI LOSS

| $\mathcal{L}_{ei}$ | EI | SF | SCD | VIF | AG | CC | VIFF |
|---|---|---|---|---|---|---|---|
| ✔ | **165.012** (↑ **11.398**) | **48.251** (↑ **3.409**) | 1.656 | **2.001** (↑ **0.104**) | **17.564** (↑ **1.186**) | 0.784 | **0.765** (↑ **0.030**) |
| ✗ | 153.614 | 44.842 | **1.659** | 1.897 | 16.378 | **0.789** | 0.735 |

In Fig. 9, the first column shows the ground truth of the four scenes. Meanwhile, the second and the third column gives the detection results of YOLOv5 corresponding to images with and without IG-GAN fusion. The red boxes indicate that the objects are missing detection. From Fig. 9, we can see that IG-GAN is beneficial in reducing missed detection by

TABLE VII
COMPARISONS OF PERFORMANCES ON THE DRONEVEHICLE DATASET IN TERMS OF MAP50, MAP75, AND MAP50:95

| Method | DroneVehicle | | | FLIR | | |
|---|---|---|---|---|---|---|
| | mAP50 | mAP75 | mAP50:95 | mAP50 | mAP75 | mAP50:95 |
| YOLOv5 | 0.757 | 0.51 | 0.467 | 0.802 | 0.368 | 0.409 |
| IG-GAN + YOLOv5 | **0.828**(↑ 0.71) | **0.721**(↑ 0.211) | **0.598**(↑ 0.131) | **0.852**(↑ 0.05) | **0.454**(↑ 0.86) | **0.461**(↑ 0.53) |

integrating RGB and infrared information. For example, under low light and extremely dark conditions, it can enhance the detection performance by leveraging the infrared information. This demonstrates the potential of IG-GAN to enhance object detection tasks, particularly in challenging scenarios by leveraging multimodal complementary information.

## V. CONCLUSION

This article describes a guided dual-stream progressive IG-GAN. In this network, the details and content streams are first established with mutual collaboration rather than independently, which contributes to detail and content exploration and cross-modal concordance enhancement. Specifically, guided interactive running-in modules ($GIR_1$ and $GIR_2$) are developed within a dual stream for intermodal alignment, cooperation, and enhancement. Then, for multilevel dual-stream information fusion, a SIF module followed by a MIF module is built. Concerning fine decoding and fused image generation, an HRM is further constructed. This is beneficial to integrate multilevel local–global contextual information. In addition, for the sake of network optimization without supervision, we further provide an objected loss function facilitating the generation of complete and detailed fusion images.

Comparative experiments with 14 state-of-the-art deep fusion methods were conducted on OS, TNO, and RGB-NIR scene datasets. Quantitative experimental results show that although many methods are highly effective, IG-GAN exhibits an evident advantage over the other 14 methods in texture details. Consistent with quantitative comparison, the fused images show that IG-GAN has superiority in completeness, texture details, and contrast. In addition, an ablation study was performed concerning SIF, MIF, and $\mathcal{L}_{ei}$. The results of the ablation experiments show that these components play a crucial role in improving the fusion performance of IG-GAN.

It is worth noting that a primary objective of digital twins (DTs) is to maintain coherence across multiple datasets, such as aligning point data with image data. In this context, IG-GAN holds promise for exploring consistency within the DTs' framework in the future.

However, most existing fusion methods, including IG-GAN, rely on the registered multimodal data pairs for exploring and fusing intermodal complementary information. Therefore, for the nonpaired multimodal data, how to explore and enhancing the cross-modal consistency information and then achieve complementary information fusion remains a challenge.

In addition, to embed IG-GAN into multimodal object detection or tracking models, it is crucial to reduce the complexity of the fusion model and improve fusion efficiency.

Note that Mamba is a simpler, more efficient, and flexible architecture compared to transformer. In this regard, we will strive to introduce Mamab to build more lightweight fusion networks in the future, thereby boosting the application potential of IG-GAN for downstream tasks, e.g., object detection and tracking.

## REFERENCES

[1] D. Hong, C. Li, B. Zhang, N. Yokoya, J. A. Benediktsson, and J. Chanussot, "Multimodal artificial intelligence foundation models: Unleashing the power of remote sensing big data in Earth observation," *Innov. Geosci.*, vol. 2, no. 1, 2024, Art. no. 100055.

[2] D. Hong, J. Hu, J. Yao, J. Chanussot, and X. X. Zhu, "Multimodal remote sensing benchmark datasets for land cover classification with a shared and specific feature learning model," *ISPRS J. Photogramm. Remote Sens.*, vol. 178, pp. 68–80, Aug. 2021.

[3] S. Salcedo-Sanz et al., "Machine learning information fusion in Earth observation: A comprehensive review of methods, applications and data sources," *Inf. Fusion*, vol. 63, pp. 256–272, Nov. 2020. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S1566253520303171

[4] D. Hong et al., "Cross-city matters: A multimodal remote sensing benchmark dataset for cross-city semantic segmentation using high-resolution domain adaptation networks," *Remote Sens. Environ.*, vol. 299, Dec. 2023, Art. no. 113856.

[5] M. Chen, X. Wang, H. Wang, and S. Zhao, "A UAV-based energy-efficient and real-time object detection system with multi-source image fusion," *J. Circuits, Syst. Comput.*, vol. 31, no. 9, Jun. 2022, Art. no. 2250166.

[6] C. Li, B. Zhang, D. Hong, J. Yao, and J. Chanussot, "LRR-Net: An interpretable deep unfolding network for hyperspectral anomaly detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 61, 2023, Art. no. 5513412.

[7] X. He, Y. Chen, L. Huang, D. Hong, and Q. Du, "Foundation model-based multimodal remote sensing data classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 62, 2024.

[8] Y. Fu and X. Wu, "A dual-branch network for infrared and visible image fusion," 2021, *arXiv:2101.09643*.

[9] D. Hong, J. Yao, C. Li, D. Meng, N. Yokoya, and J. Chanussot, "Decoupled-and-coupled networks: Self-supervised hyperspectral image super-resolution with subpixel fusion," *IEEE Trans. Geosci. Remote Sens.*, vol. 61, 2023.

[10] V. Poulain, J. Inglada, M. Spigai, J.-Y. Tourneret, and P. Marthon, "High-resolution optical and SAR image fusion for building database updating," *IEEE Trans. Geosci. Remote Sens.*, vol. 49, no. 8, pp. 2900–2910, Aug. 2011.

[11] C. Li et al., "CasFormer: Cascaded transformers for fusion-aware computational hyperspectral imaging," *Inf. Fusion*, vol. 108, Aug. 2024, Art. no. 102408.

[12] J. Ma, H. Zhang, Z. Shao, P. Liang, and H. Xu, "GANMcC: A generative adversarial network with multiclassification constraints for infrared and visible image fusion," *IEEE Trans. Instrum. Meas.*, vol. 70, pp. 1–14, 2021.

[13] Z. Yu, J. Yu, C. Xiang, J. Fan, and D. Tao, "Beyond bilinear: Generalized multimodal factorized high-order pooling for visual question answering," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 29, no. 12, pp. 5947–5959, Dec. 2018.

[14] S. Li, X. Kang, L. Fang, J. Hu, and H. Yin, "Pixel-level image fusion: A survey of the state of the art," *Inf. Fusion*, vol. 33, pp. 100–112, Jan. 2017.

[15] B. K. Shreyamsha Kumar, "Image fusion based on pixel significance using cross bilateral filter," *Signal, Image Video Process.*, vol. 9, no. 5, pp. 1193–1204, Jul. 2015.

[16] Y. Liu, X. Chen, R. K. Ward, and Z. J. Wang, "Image fusion with convolutional sparse representation," *IEEE Signal Process. Lett.*, vol. 23, no. 12, pp. 1882–1886, Dec. 2016.

[17] D. P. Bavirisetti, G. Xiao, and G. Liu, "Multi-sensor image fusion based on fourth order partial differential equations," in *Proc. 20th Int. Conf. Inf. Fusion (Fusion)*, Jul. 2017, pp. 1–9.

[18] H. Li and X. Wu, "Infrared and visible image fusion using latent low-rank representation," 2018, *arXiv:1804.08992*.

[19] M. Zhou et al., "A general spatial-frequency learning framework for multimodal image fusion," *IEEE Trans. Pattern Anal. Mach. Intell.*, pp. 1–18, 2024.

[20] H. Zhang, H. Xu, X. Tian, J. Jiang, and J. Ma, "Image fusion meets deep learning: A survey and perspective," *Inf. Fusion*, vol. 76, pp. 323–336, Dec. 2021.

[21] J. Ma, C. Chen, C. Li, and J. Huang, "Infrared and visible image fusion via gradient transfer and total variation minimization," *Inf. Fusion*, vol. 31, pp. 100–109, Sep. 2016.

[22] Z. Zhao, S. Xu, C. Zhang, J. Liu, J. Zhang, and P. Li, "DIDFuse: Deep image decomposition for infrared and visible image fusion," in *Proc. Twenty-Ninth Int. Joint Conf. Artif. Intell.*, Jul. 2020, p. 976.

[23] K. R. Shahi, P. Ghamisi, B. Rasti, P. Scheunders, and R. Gloaguen, "Unsupervised data fusion with deeper perspective: A novel multisensor deep clustering algorithm," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 15, pp. 284–296, 2022.

[24] G. Masi, D. Cozzolino, L. Verdoliva, and G. Scarpa, "Pansharpening by convolutional neural networks," *Remote Sens.*, vol. 8, no. 7, p. 594, Jul. 2016.

[25] Y. Zheng, J. Li, Y. Li, J. Guo, X. Wu, and J. Chanussot, "Hyperspectral pansharpening using deep prior and dual attention residual network," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 11, pp. 8059–8076, Nov. 2020.

[26] J. Yao, B. Zhang, C. Li, D. Hong, and J. Chanussot, "Extended vision transformer (ExViT) for land use and land cover classification: A multimodal deep learning framework," *IEEE Trans. Geosci. Remote Sens.*, vol. 61, 2023, Art. no. 5514415.

[27] H. Xu, J. Ma, J. Jiang, X. Guo, and H. Ling, "U2Fusion: A unified unsupervised image fusion network," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 44, no. 1, pp. 502–518, Jan. 2022.

[28] L. A. Hui et al., "RFN-nest: An end-to-end residual fusion network for infrared and visible images," *Inf. Fusion*, vol. 73, pp. 72–86, Sep. 2021.

[29] H. Zhang, H. Xu, Y. Xiao, X. Guo, and J. Ma, "Rethinking the image fusion: A fast unified image fusion network based on proportional maintenance of gradient and intensity," in *Proc. AAAI Conf. Artif. Intell.*, 2020, pp. 12797–12804.

[30] H. Xu, J. Ma, and X.-P. Zhang, "MEF-GAN: Multi-exposure image fusion via generative adversarial networks," *IEEE Trans. Image Process.*, vol. 29, pp. 7203–7216, 2020.

[31] W. Liao, Q. Zhang, B. Yuan, G. Zhang, and J. Lu, "Heterogeneous multidomain recommender system through adversarial learning," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 34, no. 11, pp. 8965–8977, Sep. 2022.

[32] H. Zhou, W. Wu, Y. Zhang, J. Ma, and H. Ling, "Semantic-supervised infrared and visible image fusion via a dual-discriminator generative adversarial network," *IEEE Trans. Multimedia*, vol. 25, pp. 635–648, 2023.

[33] H. Zhang, Y. Song, C. Han, and L. Zhang, "Remote sensing image spatiotemporal fusion using a generative adversarial network," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 5, pp. 4273–4286, May 2021.

[34] J. Gao, Q. Yuan, J. Li, H. Zhang, and X. Su, "Cloud removal with fusion of high resolution optical and SAR images using generative adversarial networks," *Remote Sens.*, vol. 12, no. 1, p. 191, Jan. 2020.

[35] M. Zhou, J. Huang, D. Hong, F. Zhao, C. Li, and J. Chanussot, "Rethinking pan-sharpening in closed-loop regularization," *IEEE Trans. Neural Netw. Learn. Syst.*, pp. 1–15, 2023.

[36] J. Ma, W. Yu, P. Liang, C. Li, and J. Jiang, "FusionGAN: A generative adversarial network for infrared and visible image fusion," *Inf. Fusion*, vol. 48, pp. 11–26, Aug. 2019.

[37] J. Ma, H. Xu, J. Jiang, X. Mei, and X.-P. Zhang, "DDcGAN: A dual-discriminator conditional generative adversarial network for multi-resolution image fusion," *IEEE Trans. Image Process.*, vol. 29, pp. 4980–4995, 2020.

[38] D. Hong et al., "SpectralGPT: Spectral remote sensing foundation model," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 46, no. 8, pp. 5227–5244, Aug. 2024, doi: 10.1109/TPAMI.2024.3362475.

[39] V. Vibashan et al., "Image fusion transformer," 2021, *arXiv:2107.09011*.

[40] J.-F. Hu, T.-Z. Huang, L.-J. Deng, H.-X. Dou, D. Hong, and G. Vivone, "Fusformer: A transformer-based fusion network for hyperspectral image super-resolution," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022.

[41] W. G. C. Bandara and V. M. Patel, "HyperTransformer: A textural and spectral feature fusion transformer for pansharpening," 2022, *arXiv:2203.02503*.

[42] W. Liu, Y. Lin, W. Liu, Y. Yu, and J. Li, "An attention-based multiscale transformer network for remote sensing image change detection," *ISPRS J. Photogramm. Remote Sens.*, vol. 202, pp. 599–609, Aug. 2023.

[43] J. Feng, Q. Wang, G. Zhang, X. Jia, and J. Yin, "CAT: Center attention transformer with stratified spatial–spectral token for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 62, 2024.

[44] C. Li, B. Zhang, D. Hong, X. Jia, A. Plaza, and J. Chanussot, "Learning disentangled priors for hyperspectral anomaly detection: A coupling model-driven and data-driven paradigm," *IEEE Trans. Neural Netw. Learn. Syst.*, pp. 1–14, 2024, doi: 10.1109/TNNLS.2024.3401589.

[45] D. Rao, X. Wu, and T. Xu, "TGFuse: An infrared and visible image fusion approach based on transformer and generative adversarial network," *IEEE Trans. Image Process.*, 2023.

[46] J. Ma, L. Tang, F. Fan, J. Huang, X. Mei, and Y. Ma, "SwinFusion: Cross-domain long-range learning for general image fusion via Swin transformer," *IEEE/CAA J. Autom. Sinica*, vol. 9, no. 7, pp. 1200–1217, Jul. 2022.

[47] Z. Liu et al., "Swin transformer: Hierarchical vision transformer using shifted windows," 2021, *arXiv:2103.14030*.

[48] A. Prakash, K. Chitta, and A. Geiger, "Multi-modal fusion transformer for end-to-end autonomous driving," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2021, pp. 7073–7083.

[49] X. Yan, S. Z. Gilani, H. Qin, and A. Mian, "Structural similarity loss for learning to fuse multi-focus images," *Sensors*, vol. 20, no. 22, p. 6647, Nov. 2020.

[50] J. Ma et al., "Infrared and visible image fusion via detail preserving adversarial learning," *Inf. Fusion*, vol. 54, pp. 85–98, Feb. 2020.

[51] Z. Liu et al., "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, 2021, pp. 9992–10002.

[52] L. Qu, S. Liu, M. Wang, and Z. Song, "TransMEF: A transformer-based multi-exposure image fusion framework using self-supervised multi-task learning," 2021, *arXiv:2112.01030*.

[53] Z. Zhong, X. Liu, J. Jiang, D. Zhao, Z. Chen, and X. Ji, "High-resolution depth maps imaging via attention-based hierarchical multi-modal fusion," *IEEE Trans. Image Process.*, vol. 31, pp. 648–663, 2022.

[54] J. Tang, K. Li, X. Jin, A. Cichocki, Q. Zhao, and W. Kong, "CTFN: Hierarchical learning for multimodal sentiment analysis using coupled-translation fusion network," in *Proc. 59th Annu. Meeting Assoc. Comput. Linguistics 11th Int. Joint Conf. Natural Lang. Process., (ACL/IJCNLP)*, 2021, pp. 5301–5311.

[55] Y. Xiang, R. Tao, F. Wang, and H. You, "Automatic registration of optical and SAR images VIA improved phase congruency," in *Proc. IEEE Int. Geosci. Remote Sens. Symp.*, Jul. 2019, pp. 931–934.

[56] A. Toet, "The TNO multiband image data collection," *Data Brief*, vol. 15, pp. 249–251, Dec. 2017.

[57] M. Brown and S. Süsstrunk, "Multi-spectral SIFT for scene category recognition," in *Proc. CVPR*, Jun. 2011, pp. 177–184.

[58] H. Zhang, J. Yuan, X. Tian, and J. Ma, "GAN-FM: Infrared and visible image fusion using GAN with full-scale skip connection and dual Markovian discriminators," *IEEE Trans. Comput. Imag.*, vol. 7, pp. 1134–1147, 2021.

[59] J. Ma, L. Tang, M. Xu, H. Zhang, and G. Xiao, "STDFusionNet: An infrared and visible image fusion network based on salient target detection," *IEEE Trans. Instrum. Meas.*, vol. 70, pp. 1–13, 2021.

[60] L. Jian, S. Xiong, H. Yan, X. Niu, S. Wu, and D. Zhang, "Rethinking cross-attention for infrared and visible image fusion," 2024, *arXiv:2401.11675*.

[61] L. Tang, H. Zhang, H. Xu, and J. Ma, "Rethinking the necessity of image fusion in high-level vision tasks: A practical infrared and visible image fusion network based on progressive semantic injection and scene fidelity," *Inf. Fusion*, vol. 99, Nov. 2023, Art. no. 101870.

[62] V. Aslantas and E. Bendes, "A new image quality metric for image fusion: The sum of the correlations of differences," *AEU - Int. J. Electron. Commun.*, vol. 69, no. 12, pp. 1890–1896, Dec. 2015.

[63] H. R. Sheikh and A. C. Bovik, "A visual information fidelity approach to video quality assessment," *Neuroscience*, vol. 7, no. 2, pp. 2117–2128, 2005.

[64] Y. Han, Y. Cai, Y. Cao, and X. Xu, "A new image fusion performance metric based on visual information fidelity," *Inf. Fusion*, vol. 14, no. 2, pp. 127–135, Apr. 2013.

[65] Y. Sun, B. Cao, P. Zhu, and Q. Hu, "Drone-based RGB-infrared cross-modality vehicle detection via uncertainty-aware learning," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 32, no. 10, pp. 6700–6713, Oct. 2022.

[66] H. Zhang, E. Fromont, S. Lefevre, and B. Avignon, "Multispectral fusion for object detection with cyclic fuse-and-refine blocks," in *Proc. IEEE Int. Conf. Image Process. (ICIP)*, Oct. 2020, pp. 276–280.

[67] Z. Zou, K. Chen, Z. Shi, Y. Guo, and J. Ye, "Object detection in 20 years: A survey," *Proc. IEEE*, vol. 111, no. 3, pp. 257–276, Mar. 2023.

**Chenhong Sui** (Member, IEEE) received the master's and Ph.D. degrees in engineering from the Department of Telecommunications, Huazhong University of Science and Technology, Wuhan, China, in 2012 and 2015, respectively.

In 2015, she joined Yantai University, Yantai, China, where she has been serving as an Associate Professor since 2018. In 2023, she joined the Aerospace Information Research Institute, Chinese Academy of Sciences, Beijing, China, for post-doctoral research, with the main research interests include multimodal data fusion, anti-attack, and defense.

**Guobin Yang** received the master's degree from the School of Physics and Information Engineering, Yantai University, Yantai, China, in 2023.
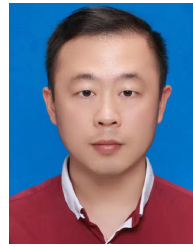
Currently, he is a Faculty Member with the Weifang Vocational College, Weifang, China. His research interests include multimodal fusion, machine learning, remote sensing, and object detection.

**Danfeng Hong** (Senior Member, IEEE) received the Dr.-Ing. degree (summa cum laude) from the Signal Processing in Earth Observation (SiPEO), Technical University of Munich (TUM), Munich, Germany, in 2019.

Since 2022, he has been a Full Professor with the Aerospace Information Research Institute, Chinese Academy of Sciences, Beijing, China. His research interests include artificial intelligence, multimodal remote sensing, foundation models, hyperspectral imaging, and Earth observation.

Dr. Hong is an Associate Editor for the IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING (TGRS) and the IEEE TRANSACTIONS ON IMAGE PROCESSING (TIP), and the Editorial Board Member of Information Fusion, *ISPRS Journal of Photogrammetry and Remote Sensing*, and *International Journal of Applied Earth Observation and Geoinformation*. He received the Jose Bioucas Dias Award for recognizing an outstanding paper at WHISPERS in 2021, the Remote Sensing Young Investigator Award in 2022, the IEEE GRSS Early Career Award in 2022, the MIT Technology Review & DeepTech "China's Intelligent Computing Innovators" award in 2024, and a Highly Cited Researcher (Clarivate Analytics/Thomson Reuters) since 2022.

**Haipeng Wang** (Member, IEEE) received the Ph.D. degree from the Institute of Information Fusion, Naval Aviation University, Yantai, China, in 2012.

In 2019, he joined the School of Computer Science and Technology, Harbin Institute of Technology, Harbin, China, to conduct post-doctoral research. Since 2020, he has been a Professor with Naval Aviation University. His main research interests include multimodal data fusion, object detection, and tracking.

**Jing Yao** (Member, IEEE) received the Ph.D. degree in mathematics from Xi'an Jiaotong University, Xi'an, China, in 2021.
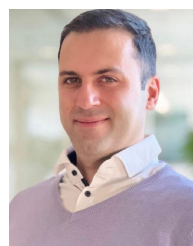
From 2019 to 2020, he was a Visiting Student at the Signal Processing in Earth Observation (SiPEO), Technical University of Munich (TUM), Munich, Germany, and at the Remote Sensing Technology Institute (IMF), German Aerospace Center (DLR), Oberpfaffenhofen, Germany. Since 2021, he has been an Assistant Professor with the Aerospace Information Research Institute, Chinese Academy of Sciences, Beijing, China. His research interests include hyperspectral and multimodal remote sensing image analysis, mainly including optimization and deep learning-based methods for image processing and interpretation applications.

Dr. Yao is serving as the Topical Associate Editor for IEEE TGRS. He was a recipient of the Best Reviewer Award of IEEE TGRS in 2024, the Best Reviewer Award of IEEE JSTARS in 2024, and the Jose Bioucas Dias Award for recognizing the outstanding paper at WHISPERS in 2021.

**Peter M. Atkinson** has been the Executive Dean of the Faculty of Science and Technology, Lancaster University, Lancaster, U.K., since 2015. Previously, he was the Executive Dean of the Faculty of Health and Medicine (a role he held simultaneously with being Executive Dean of FST from 2018 to 2019), Lancaster University. He was the Chair of the university's successful Athena SWAN institutional submission from 2018 to 2019. Previously, he was the Head of the School of Geography, University of Southampton, Southampton, U.K., from 2007 to 2012, in an executive role with full budgetary responsibility. Following this, at the University of Southampton, he was the Director of REF strategy (the academic lead for Southampton's REF2014 submission). He is currently a Distinguished Professor of spatial data science with Lancaster University and an Interdisciplinary Scientist. Specifically, his research involves the application of space–time statistics and geostatistics, machine learning and AI, and dynamic numerical modeling, to Earth observation (EO) and other spatio-temporal data, to answer a wide range of science and social science questions.

**Pedram Ghamisi** (Senior Member, IEEE) received the Ph.D. degree in electrical and computer engineering from the University of Iceland, Reykjavik, Iceland, in 2015.

He currently serves as the Head of the Machine Learning Group, Helmholtz-Zentrum Dresden-Rossendorf (HZDR), Dresden, Germany, and a Visiting Professor with Lancaster University, Lancaster, U.K. His research interests encompass deep learning, with a specific focus on remote sensing applications. For detailed information, please visit http://www.ai4rs.com.