

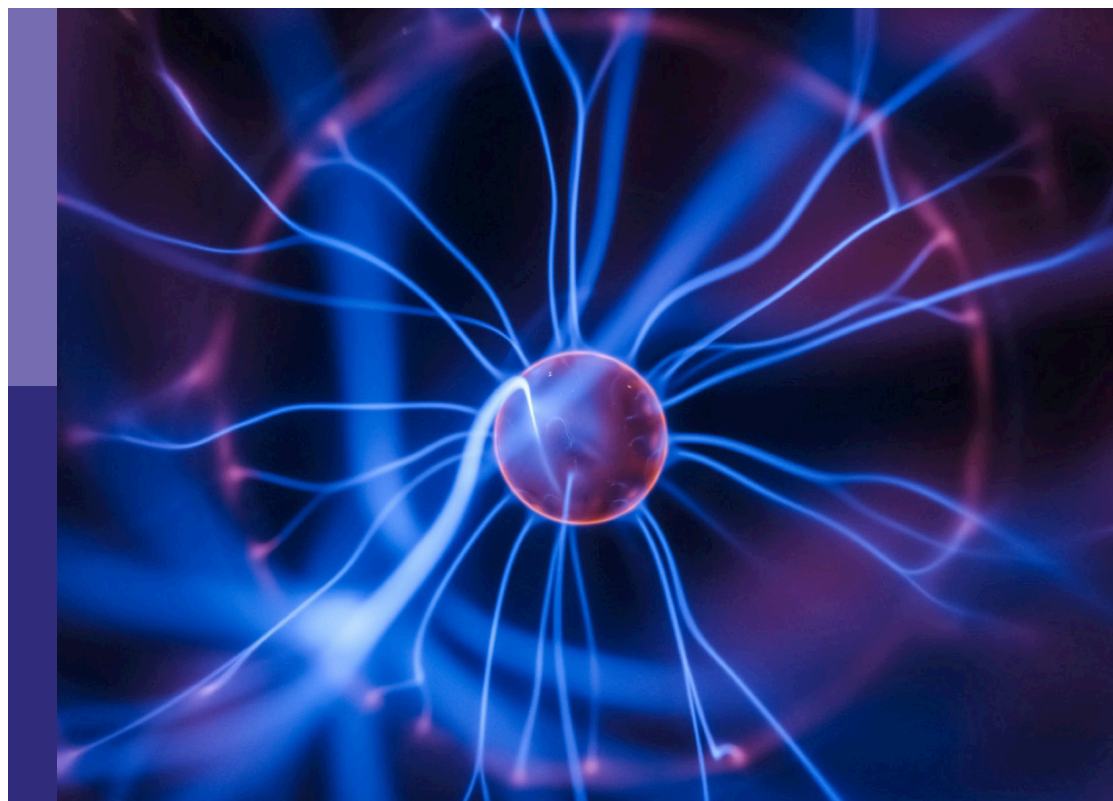
Multi-sensor imaging and fusion: methods, evaluations, and applications, volume III

Edited by

Zhiqin Zhu, Yu Liu, Huafeng Li, Jinxing Li
and Bo Xiao

Published in

Frontiers in Physics



FRONTIERS EBOOK COPYRIGHT STATEMENT

The copyright in the text of individual articles in this ebook is the property of their respective authors or their respective institutions or funders. The copyright in graphics and images within each article may be subject to copyright of other parties. In both cases this is subject to a license granted to Frontiers.

The compilation of articles constituting this ebook is the property of Frontiers.

Each article within this ebook, and the ebook itself, are published under the most recent version of the Creative Commons CC-BY licence. The version current at the date of publication of this ebook is CC-BY 4.0. If the CC-BY licence is updated, the licence granted by Frontiers is automatically updated to the new version.

When exercising any right under the CC-BY licence, Frontiers must be attributed as the original publisher of the article or ebook, as applicable.

Authors have the responsibility of ensuring that any graphics or other materials which are the property of others may be included in the CC-BY licence, but this should be checked before relying on the CC-BY licence to reproduce those materials. Any copyright notices relating to those materials must be complied with.

Copyright and source acknowledgement notices may not be removed and must be displayed in any copy, derivative work or partial copy which includes the elements in question.

All copyright, and all rights therein, are protected by national and international copyright laws. The above represents a summary only. For further information please read Frontiers' Conditions for Website Use and Copyright Statement, and the applicable CC-BY licence.

ISSN 1664-8714
ISBN 978-2-8325-7037-1
DOI 10.3389/978-2-8325-7037-1

Generative AI statement

Any alternative text (Alt text) provided alongside figures in the articles in this ebook has been generated by Frontiers with the support of artificial intelligence and reasonable efforts have been made to ensure accuracy, including review by the authors wherever possible. If you identify any issues, please contact us.

About Frontiers

Frontiers is more than just an open access publisher of scholarly articles: it is a pioneering approach to the world of academia, radically improving the way scholarly research is managed. The grand vision of Frontiers is a world where all people have an equal opportunity to seek, share and generate knowledge. Frontiers provides immediate and permanent online open access to all its publications, but this alone is not enough to realize our grand goals.

Frontiers journal series

The Frontiers journal series is a multi-tier and interdisciplinary set of open-access, online journals, promising a paradigm shift from the current review, selection and dissemination processes in academic publishing. All Frontiers journals are driven by researchers for researchers; therefore, they constitute a service to the scholarly community. At the same time, the *Frontiers journal series* operates on a revolutionary invention, the tiered publishing system, initially addressing specific communities of scholars, and gradually climbing up to broader public understanding, thus serving the interests of the lay society, too.

Dedication to quality

Each Frontiers article is a landmark of the highest quality, thanks to genuinely collaborative interactions between authors and review editors, who include some of the world's best academicians. Research must be certified by peers before entering a stream of knowledge that may eventually reach the public - and shape society; therefore, Frontiers only applies the most rigorous and unbiased reviews. Frontiers revolutionizes research publishing by freely delivering the most outstanding research, evaluated with no bias from both the academic and social point of view. By applying the most advanced information technologies, Frontiers is catapulting scholarly publishing into a new generation.

What are Frontiers Research Topics?

Frontiers Research Topics are very popular trademarks of the *Frontiers journals series*: they are collections of at least ten articles, all centered on a particular subject. With their unique mix of varied contributions from Original Research to Review Articles, Frontiers Research Topics unify the most influential researchers, the latest key findings and historical advances in a hot research area.

Find out more on how to host your own Frontiers Research Topic or contribute to one as an author by contacting the Frontiers editorial office: frontiersin.org/about/contact

Multi-sensor imaging and fusion: methods, evaluations, and applications, volume III

Topic editors

Zhiqin Zhu — Chongqing University of Posts and Telecommunications, China

Yu Liu — Hefei University of Technology, China

Huafeng Li — Kunming University of Science and Technology, China

Jinxing Li — Harbin Institute of Technology, Shenzhen, China

Bo Xiao — Imperial College London, United Kingdom

Citation

Zhu, Z., Liu, Y., Li, H., Li, J., Xiao, B., eds. (2025). *Multi-sensor imaging and fusion: methods, evaluations, and applications, volume III*. Lausanne: Frontiers Media SA. doi: 10.3389/978-2-8325-7037-1

Table of contents

04	Character-interested binary-like image learning for text image demoiréing Zhanpei Zhang, Beicheng Liang, Tingting Ren, Chengmiao Fan, Rui Li and Mu Li
18	Multi-Conv attention network for skin lesion image segmentation Zexin Li, Hanchen Wang, Haoyu Chen, Chenxin Lin and Aochen Yan
31	SGL-YOLOv9: an effective method for crucial components detection in the power distribution network Mianfang Yang, Bojian Chen, Chenxiang Lin, Wenxu Yao and Yangdi Li
46	Performance evaluation of photovoltaic scenario generation Siyu Ren, Tongxin Yang, Jun Luo, Gang Wu, Kai Mao and Bowen Liu
62	Multi-focus image fusion based on pulse coupled neural network and WSEML in DTCWT domain Yuan Jia and Tiande Ma
73	GLI-Net: A global and local interaction network for accurate classification of gastrointestinal diseases in endoscopic images Yuansen Zhang, Mengxiao Zhuang, Wenjun Chen, Xiaoqiu Wu and Qingqing Song
85	Infrared and visible image fusion driven by multimodal large language models Ke Wang, Dengshu Hu, Yuan Cheng, Yukui Che, Yuelin Li, Zhiwei Jiang, Fengxian Chen and Wenjuan Li
99	Perceptual objective evaluation for multimodal medical image fusion Chuangeng Tian, Juyuan Zhang and Lu Tang
108	Target-aware unregistered infrared and visible image fusion Dengshu Hu, Ke Wang, Cuijin Zhang, Zheng Liu, Yukui Che, Shoubing Dong and Chuirui Kong
122	Multi-sensor fusion for AI-driven behavior planning in medical applications Chang Jianming, Qin Yuanyuan, Xu Yanling, Li Li, Wu Mianhua and Wang Lulu
141	Multi-modal action recognition via advanced image fusion techniques for cyber-physical systems Zaiyong Shou and Daoyu Zhu



OPEN ACCESS

EDITED BY

Zhiqin Zhu,
Chongqing University of Posts and
Telecommunications, China

REVIEWED BY

Lingxiao Yang,
Sun Yat-sen University, China
Duo Chen,
University of Electronic Science and
Technology of China, China
Xixi Jia,
Xidian University, China

*CORRESPONDENCE

Mu Li,
✉ limu2022@hit.edu.cn

RECEIVED 11 November 2024

ACCEPTED 19 November 2024

PUBLISHED 13 December 2024

CITATION

Zhang Z, Liang B, Ren T, Fan C, Li R and Li M
(2024) Character-interested binary-like image
learning for text image demoiréing.
Front. Phys. 12:1526412.
doi: 10.3389/fphy.2024.1526412

COPYRIGHT

© 2024 Zhang, Liang, Ren, Fan, Li and Li. This
is an open-access article distributed under
the terms of the [Creative Commons
Attribution License \(CC BY\)](#). The use,
distribution or reproduction in other forums is
permitted, provided the original author(s) and
the copyright owner(s) are credited and that
the original publication in this journal is cited,
in accordance with accepted academic
practice. No use, distribution or reproduction
is permitted which does not comply with
these terms.

Character-interested binary-like image learning for text image demoiréing

Zhanpei Zhang, Beicheng Liang, Tingting Ren, Chengmiao Fan,
Rui Li and Mu Li*

Department of Computer Science and Technology, Harbin Institute of Technology, Shenzhen,
Guangdong, China

Despite the fact that the text image-based optical character recognition (OCR) methods have been applied to a wide range of applications, they do suffer from performance degradation when the image is contaminated with moiré patterns for the sake of interference between the display screen and the camera. To tackle this problem, we propose a novel network for text image demoiréing. Specifically, to encourage our study on text images, we collected a dataset including a number of pairs of images with/without moiré patterns, which is specific for text image demoiréing. In addition, due to the statistical differences among various channels on moiré patterns, a multi-channel strategy is proposed, which roughly extracts the information associated with moiré patterns and subsequently contributes to moiré removal. In addition, our purpose on the text image is to increase the OCR accuracy, while other background pixels are insignificant. Instead of restoring all pixels like those in natural images, a character attention module is conducted, allowing the network to pay more attention on the optical character-associated pixels and also achieving a consistent image style. As a result from this method, characters can be more easily detected and more accurately recognized. Dramatic experimental results on our conducted dataset demonstrate the significance of our study and the superiority of our proposed method compared with state-of-the-art image restoration approaches. Specifically, the metrics of recall and F1-measure on recognition are increased from 56.32%/70.18% to 85.34%/89.36%.

KEYWORDS

multi-sensor imaging, deep learning, text image, demoiréing, multi-channel, moiré pattern, optical character recognition

1 Introduction

Due to the huge number of text images, the automatic text recognition from a given image is quite necessary in recent years. Thanks to the techniques of optical character recognition (OCR) [1–3], image-based text detection [4, 5] and recognition [6] have been effectively improved and are widely applied to many applications, such as ID card recognition [7], table recognition [8], and license plate recognition [9, 10]. Despite the fact that these methods have achieved satisfactory performances, they are sensitively influenced by the quality of images. As displayed in Figure 1, it is a general and inevitable phenomenon that the captured image is corrupted with diverse moiré patterns due to interference between the display

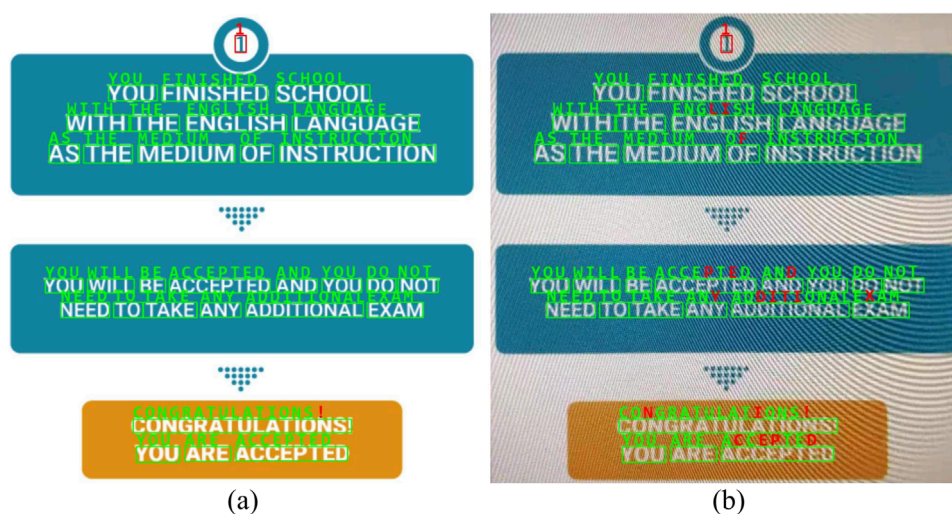


FIGURE 1
OCR on the images with (A) or without (B) moiré patterns. Characters in green denote the accurate recognition, and characters in red denote the inaccurate recognition. Making a comparison between (A) and (B), recognition accuracy on the text image is significantly influenced by the moiré patterns.

screen and the camera, resulting in significant performance degradation in both character detection and recognition. Thus, in this paper, we focus on the moiré pattern removal from the text images for OCR.

It is particularly challenging to remove moiré patterns from photographs. Different from other corruptions, such as noise [11, 12], rain [13, 14], and haze [15, 16], the moiré pattern exhibits a diverse range of characteristics. Specifically, as shown in Figure 1B, colors, thickness, and shapes (stripes or ripples) are even diverse across different areas in a photograph, and the frequency domain, as analyzed in [17], further demonstrates its complexity.

To restore the image, [18] proposed a convolutional neural network (CNN), in which a multi-resolution strategy is adopted to remove the moiré patterns from a wide range of frequencies. Inspired by this work, other studies [17, 19–22] have been proposed for image demoiréing. Despite the fact that these aforementioned works effectively obtain a moiré-free image from the input, they are only adaptive for natural images, as the structures between text images and natural images differ significantly. Compared with natural images, the key information in text images is the optical characters. In other words, the purpose of text image demoiréing is to improve the accuracy of text recognition after restoration, which encourages us to pay more attention on the optical character-associated pixels. Thus, not only the moiré patterns should be removed from the raw image but also the semantic structures of optical characters should be enhanced.

To achieve this goal, we propose the text image demoiréing network (TIDNet). Considering that the moiré pattern in the G (green) channel is statistically weaker than that in the R (red) and B (blue) channels [17], its edge patterns are roughly but adaptively extracted by our presented rough moiré pattern extraction module, regardless of whether the scales of values in the R, G, and B channels are different or similar. Furthermore, we also propose a character attention module, allowing the network to particularly pay much more attention on the optical characters for our OCR application.

In detail, as shown in Figure 2, it is obvious that under different viewpoints and capturing distances, colors of moiré-contaminated images captured from the same image differ significantly, making complete recovery more difficult. In addition, if an image is covered by watermarks (Figure 2A), it seems impossible to restore it from the contaminated images due to the missing information in image collection (Figures 2B–D). Subsequently, the inaccurate background pixel estimation may even inversely result in the degradation of performance. In fact, we need to improve the recognition accuracy. The greater the difference between the foreground and background, the easier it is to detect and recognize the text. Thus, apart from image demoiréing, we further transform diverse image styles to a consistent version, where the background pixels are white, while the foreground characters are black. Thanks to this strategy, not only the estimation for the complex background is avoided but also the difference between the characters and background pixels is enlarged, contributing to both character detection and recognition. In addition, a mask strategy and a semantic measurement are jointly introduced, allowing our model to pay much more attention on the character-associated pixels.

In order to achieve moiré pattern removal, a dataset is necessary. In addition, we create a text image dataset named HITSZ-TID, which is composed of 3,739 pairs of images. For each pair, it consists of an image contaminated with moiré patterns, as well as its associated reference image without moiré patterns. Particularly, we extract the contaminated image under multiple devices, viewpoints, and distances, ensuring the diversity and generalization of our collected dataset.

The main contributions of this paper are as follows:

- A text image demoiréing network (TIDNet) is particularly designed for text image demoiréing. Thanks to our proposed method, the recognition accuracy on text images contaminated with moiré patterns is greatly improved. Values of Recall and

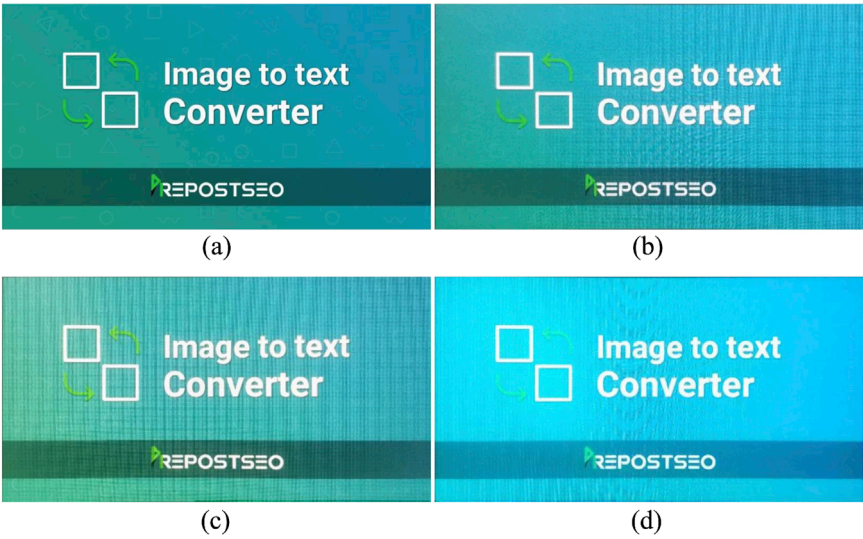


FIGURE 2
(A) Image without moiré patterns. (B–D) Images with moiré patterns, which are captured under different viewpoints and distances.

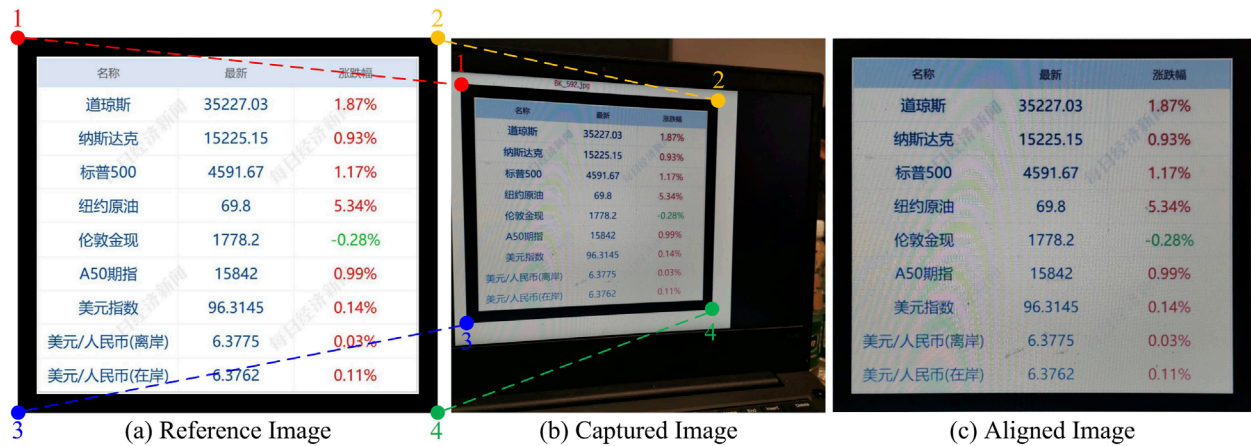


FIGURE 3
Image alignment. The reference image and the captured image contaminated with moiré patterns are aligned according to their corresponding corners: (A) reference image, (B) captured image, and (C) aligned image.

- F1-measure on recognition increased from 56.32%/70.18% to 85.34%/89.36%.
- The rough moiré pattern extraction module and character attention module are jointly introduced into our TIDNet. Due to the differences in different channels on the moiré patterns, the moiré is first removed roughly. Furthermore, the textural and semantic characters are also exploited, which are specifically adaptive for text image moiré removal.
 - A dataset HITSZ-TID which is for text image demoiréing is created. It consists of 3,739 image pairs, where each pair contains an image contaminated with moiré patterns and its corresponding reference image free from the moiré patterns. This dataset fills the gap between the OCR and image demoiréing, contributing to the research study on these two fields.

The rest of this paper is organized as follows. In [Section 2](#), some related works about image demoiréing and text image processing are briefly described. Our created dataset and proposed TIDNet are then introduced in [Section 3](#) and [Section 4](#), respectively. To demonstrate the significance of text image demoiréing for OCR and the effectiveness of our proposed method, we conducted experiments in [Section 5](#), followed by conclusion in [Section 6](#).

2 Related works

In this section, we briefly introduce the related works on image demoiréing and text image processing.

TABLE 1 Detailed information of mobile phones and display screens for capturing images.

Mobile phone	Display screen	
	Model	Resolution
Huawei Mate 30 Pro	AIDU LJ240S	1920 × 1080
Redmi Note 11 Pro	Redmi RMMNT238NF	1920 × 1080
iPhone 8 Plus	Hanpon E2206	1920 × 1080
VIVO X21S	ThinkPad E450	1366 × 768
Redmi MAX3	ThinkPad E14	1920 × 1080
iPhone 8	ThinkPad E490	1920 × 1080
Huawei Nova 5	—	—

2.1 Image demoiréing

Due to the interference of different repetitive patterns, the image contaminated with moiré patterns is an inevitable phenomenon. In recent years, various methods have been proposed for moiré pattern removal. By exploiting the prior assumption that moiré patterns are dissimilar on textures, a low-rank and sparse matrix decomposition method [23] was developed to achieve demoiréing on high-frequency textures. Different from this hand-crafted feature-based method, [18] primarily utilized the CNN for moiré image restoration. Considering that moiré patterns widely span in different resolution scales, multiple resolutions were jointly exploited in [18]. Followed by it, [21] also presented a multi-scale feature enhancement network for moiré image restoration. In addition, a coarse-to-fine strategy was presented in [24], which introduced another fine-scale network to refine the demoiréed image obtained from the coarse-scale network. In addition, instead of relying solely on real captured images like [18], [24] modeled the formation of moiré patterns and generated a large-scale synthetic dataset. Furthermore, [20] proposed a learnable bandpass filter and a two-step tone-mapping strategy for moiré pattern removal and color restoration, respectively. [25] constructed a moiré removal and brightness improvement (MRBI) database using aligned moiré-free and moiré images and proposed a CNN with additive and multiplicative modules to transfer the low light moiré image to the bright moiré-free image. Considering that the moiré patterns mainly located on the high-frequency domain, the wavelet was embedded into the network [26], in which the features represented by the wavelet transformation were then processed. To compensate for the difference in domains between the training and the testing sets, a domain adaptation mechanism was further exploited to fine-tune the output. Similarly, [27] also introduced a wavelet-based dual-branch network to separate the frequencies of moiré patterns from the image content. By exploiting progressive feature fusion and channel-wise attention, the attentive fractal network was proposed in [28]. In addition, [29] proposed another attention network named C3Net, which focuses on channel, color, and concatenation. Different from these aforementioned methods from single-image

demoiréing, the multi-frame-based image demoiréing was also studied in [19].

Despite the fact that a number of deep learning-based approaches have been proposed for moiré-free image restoration, almost all of them are designed for natural images, which are not particularly adaptive for the text images.

2.2 Text image processing

The quality of the text image has a key influence on the accuracy of ORC. According to this purpose, some works on text image processing have been studied. For instance, several artificial filters were compared on low-resolution text images [30]. Subsequently, SRCNN [31] was applied to the text image super-resolution [32]. To achieve the scene text image super-resolution, [33] designed a text-oriented network, in which the sequential information and character boundaries were enhanced. In addition, in [34], the image was decomposed into the text, foreground, and background, which were beneficial for text boundary recovery and color restoration, respectively. Considering the text-specific properties, [35] utilized the text-level layouts and character-level details for text image super-resolution. Apart from this super-resolution application, some deblurring approaches [36–42] have also been proposed for text images. Specifically, [38] introduced two-tone prior to estimate the kernel for image deblurring. The deep neural network followed by sequential highway connections was exploited to restore the blurry image to a clear image. Furthermore, by constructing a text-specific hybrid dictionary, the powerful contextual information was then extracted for blind text image deblurring [39, 43, 44]. For the text image detection and recognition, [45] proposed a mathematical model based on the Riesz fractional operator to enhance details of the edge information in license plate images, hence improving the performance. In addition, a method [46] for predicting hidden (masked) text parts was proposed to fill the gaps of non-transcribable parts in the unstructured document OCR.

Although these methods were studied for text image super-resolution or deblurring, they are not adaptive for the application of demoiréing, due to the much more complex distributions or structures of the moiré patterns. Therefore, it is quite significant to propose a specific network for text image demoiréing. A related work named MsMa-Net [47] was proposed for moiré removal in document images. However, our proposed method is quite different from that of [47]. Referring to the dataset, only 80 images were used for dataset construction, whereas 551 images were used in our dataset, resulting in 3,739 pairs. Furthermore, we further take text priors, e.g., gradient, channel, and semantic information, into account, which contribute to our performance improvement on detection and recognition. In addition, although MsMa-Net also mentioned binarization for the output, it still first enforced the output to be the same to the reference image in the color version, which was then followed by a threshold processing to achieve binarization. By contrast, our proposed dataset TIDNet directly transforms various inputs to a binary-like ground-truth without any reference estimation, making it easier to remove the influences of diverse backgrounds and contributes to image reconstruction. Thus, this work will considerably benefit future research on text image processing and OCR.

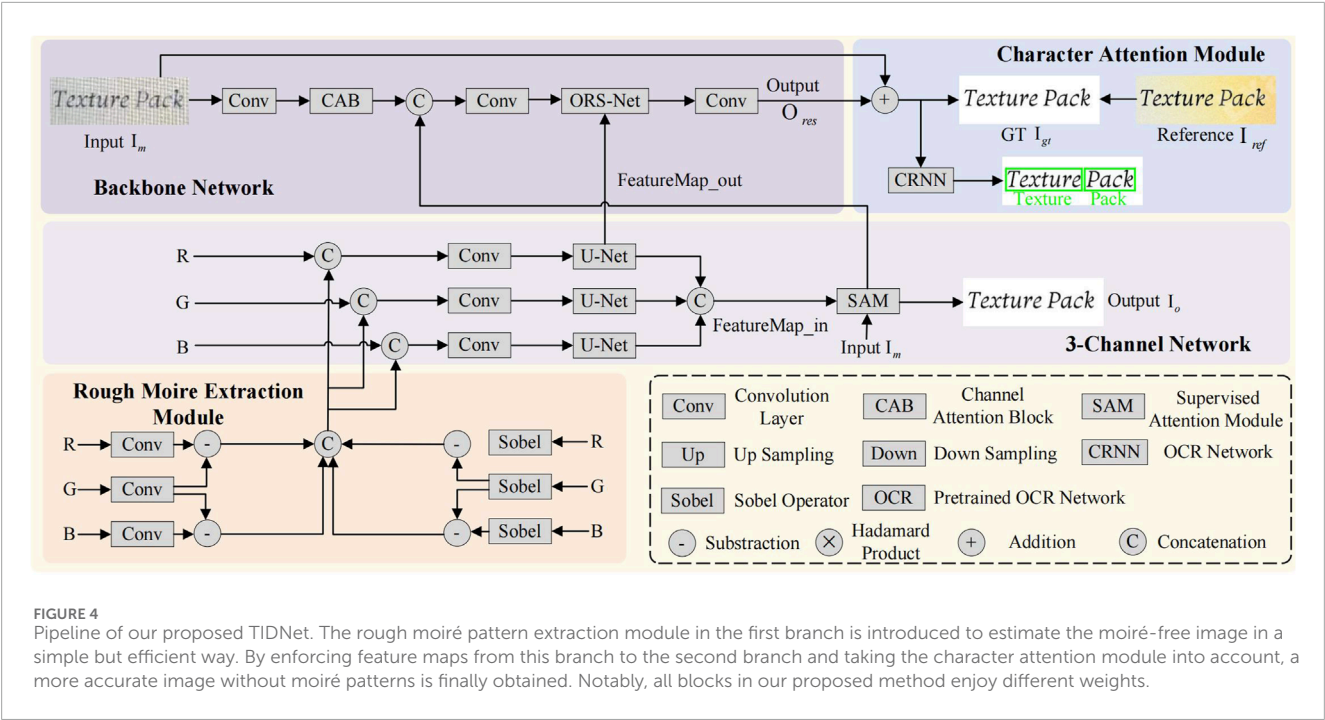


FIGURE 4 Pipeline of our proposed TIDNet. The rough moiré pattern extraction module in the first branch is introduced to estimate the moiré-free image in a simple but efficient way. By enforcing feature maps from this branch to the second branch and taking the character attention module into account, a more accurate image without moiré patterns is finally obtained. Notably, all blocks in our proposed method enjoy different weights.

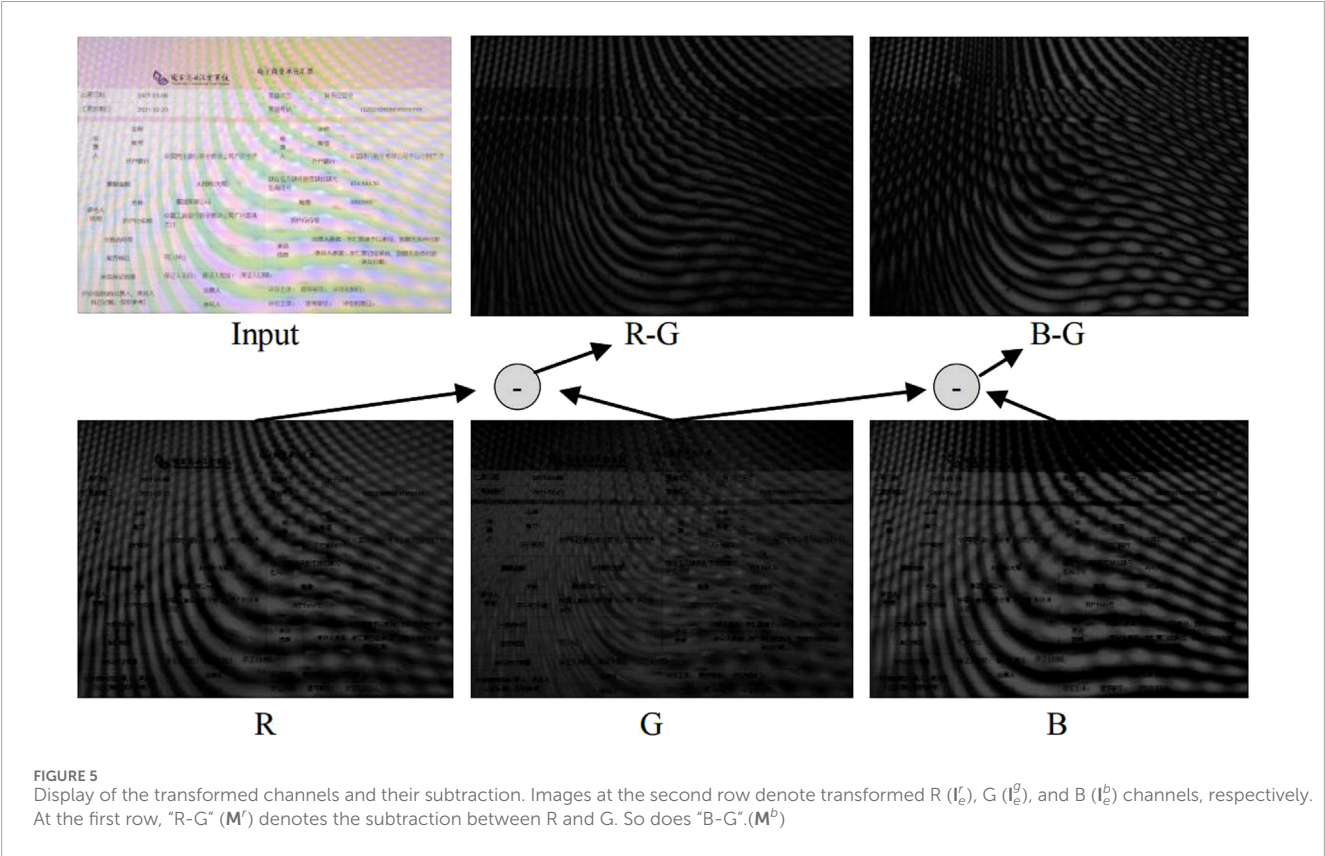
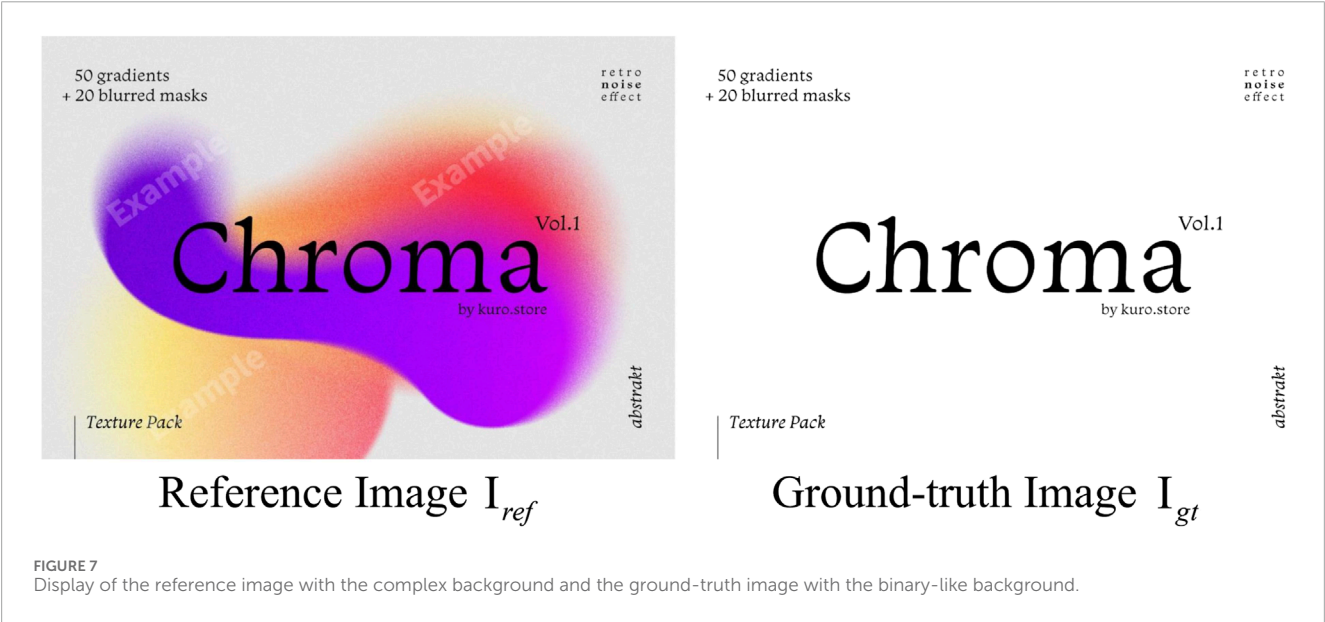
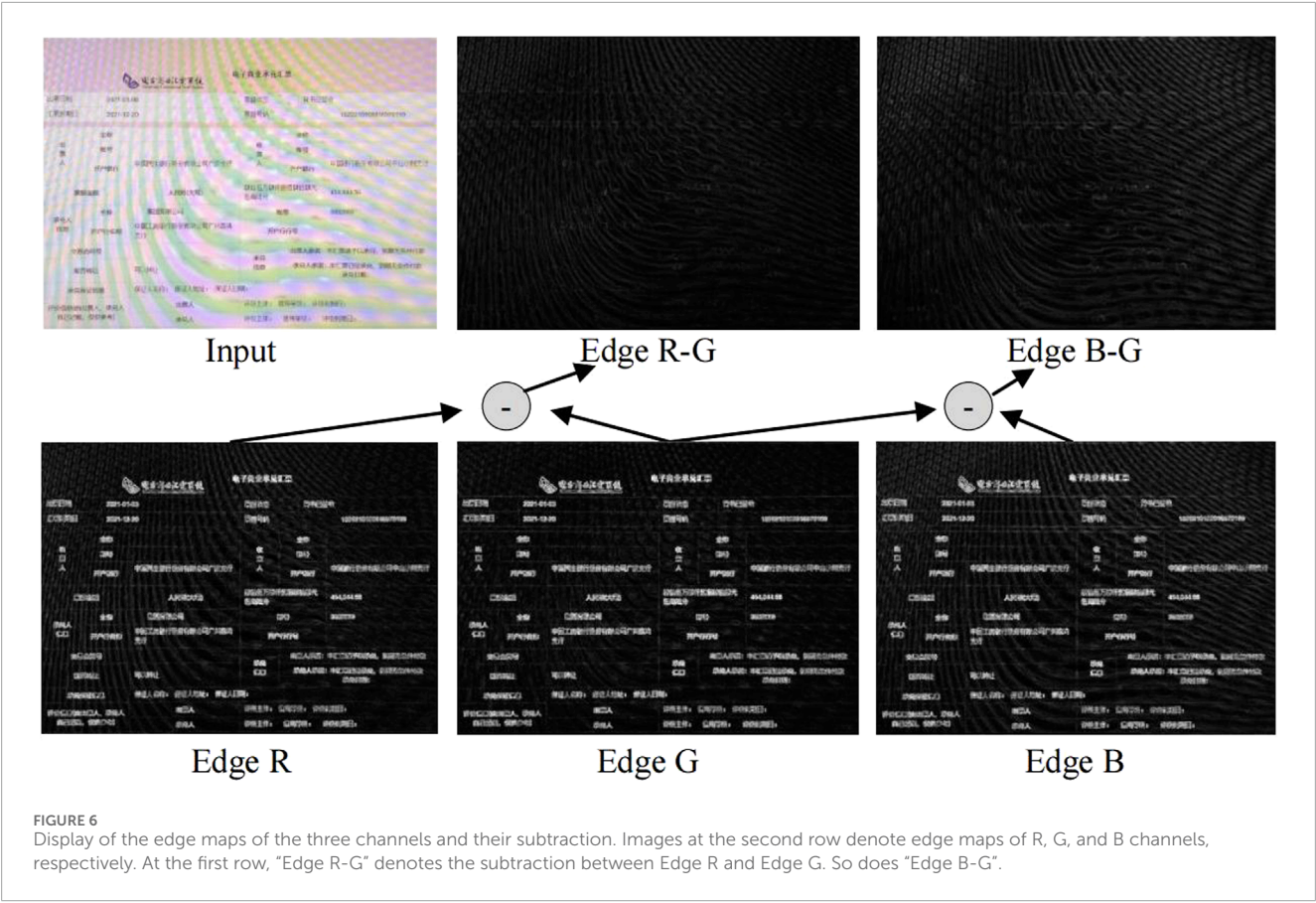


FIGURE 5 Display of the transformed channels and their subtraction. Images at the second row denote transformed $R(I_R^G)$, $G(I_G^G)$, and $B(I_B^G)$ channels, respectively. At the first row, “R-G” (M^G) denotes the subtraction between R and G. So does “B-G”. (M^B)

3 Dataset

In this study, for training and testing purposes, we collect 3,739 pairs of contaminated moiré images and uncontaminated

reference images to serve as a text image benchmark for moiré pattern removal. Specifically, we download the reference text images in Chinese or English from the internet, which are then used for capturing contaminated images.



3.1 Image capture

Similar to [18], each reference image is surrounded by a black border for alignment, which will be analyzed in Subsection 3.2. As displayed in Figure 3, the image is first located in the center of the display screen, which is then captured using a mobile phone.

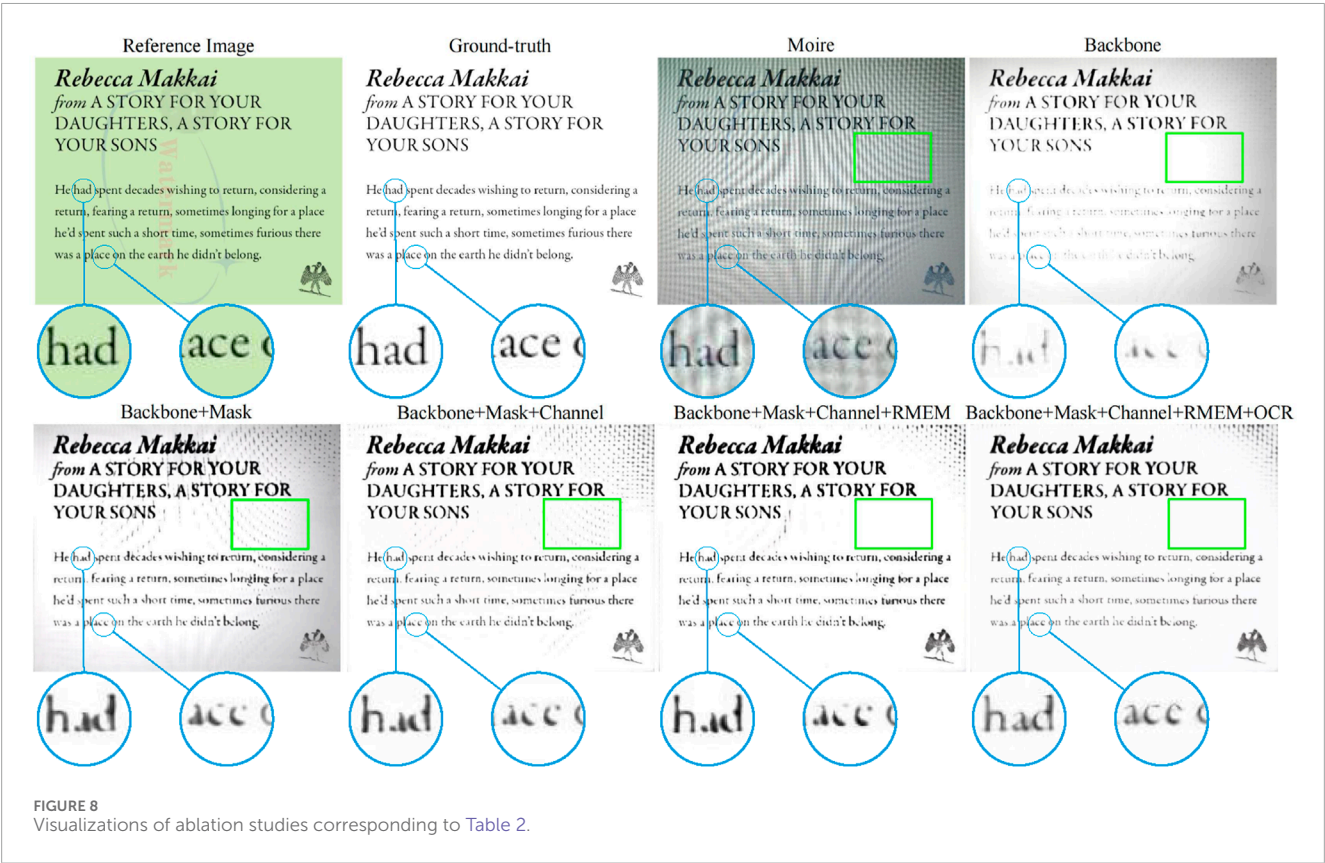
Notably, the black border is always completely captured, and each photo is taken from a random distance or viewpoint, guaranteeing the diversity of the moiré patterns.

To further enhance diversity in our created dataset, we use a variety of mobile phones and monitor screens. Table 1 lists the detailed information of our used mobile phones and display screens.

TABLE 2 Ablation studies conducted on our collected dataset.

Task	Detection			Recognition		
	Recall	Precision	F1-measure	Recall	Precision	F1-measure
Moiré	53.82%	99.15%	69.77%	56.32%	93.10%	70.18%
Backbone	24.88%	98.99%	39.77%	24.38%	80.54%	37.43%
Backbone + Mask	58.79%	99.01%	73.77%	53.85%	80.86%	64.65%
Backbone + Mask + Channel	64.13%	99.34%	77.94%	57.31%	81.05%	67.15%
Backbone + Mask + Channel + RMEM	68.93%	98.88%	81.23%	62.02%	83.02%	71.00%
Backbone + Mask + Channel + RMEM + OCR	92.92%	98.72%	95.73%	85.34%	93.78%	89.36%

“Moiré” denotes results on the raw image without any processing. “Backbone” denotes results obtained by the baseline network, which is guided by L_b . “Backbone + Mask” denotes results by adding the mask loss L_m . “Backbone + Mask + Channel” denotes results by additionally introducing the three-channel network. “Backbone + Mask + Channel + RMEM” denotes results by additionally introducing the rough moiré extraction module (RMEM). Similarly, “Backbone + Mask + Channel + OCR” denotes results by adding the OCR semantic loss L_{ocr} . Notably, the best performance is highlighted by “bold.”



Specifically, eight types of mobile phones and seven types of display screens are used for capturing images. Taking other aforementioned variables into account such as distances and viewpoints, 3,739 pairs of images are totally obtained.

3.2 Image alignment

To achieve the training phase in an end-to-end way, the contaminated image should be aligned with its corresponding

reference image at the pixel-to-pixel level. Although [18, 25] proposed the corner or patch matching algorithms for image alignment, these automatic strategies still encounter a slight misalignment. Different from the natural images, the misalignment under even several pixels would make a great influence on the text image restoration. Thus, we manually detect the corresponding corners for the text image alignment. As shown in Figures 3A, B, four corners in the reference image and contaminated image are detected, respectively, through which the geometric transformation between these two images is estimated. Finally,

TABLE 3 Results obtained by our proposed method guided by different reference images.

Task		Detection	
Metrics	Recall	Precision	F1-measure
TIDNet (color)	88.38%	99.12%	93.44%
TIDNet	92.92%	98.72%	95.73%
Task		Recognition	
TIDNet (color)	83.34%	93.52%	88.14%
TIDNet	85.34%	93.78%	89.36%

“TIDNet (color)” and “TIDNet” denote our proposed method is supervised by \mathbf{I}_{ref} with diverse backgrounds and \mathbf{I}_e with the consistent background, respectively.

we obtained the aligned image with moiré patterns, as displayed in Figure 3C.

4 Proposed method

The pipeline in our proposed method is shown in Figure 4. It is clear that there are two branches for the moiré-free image generation. From the bottom to top, our proposed rough moiré extraction module and the three-channel network are first exploited to remove the moiré pattern in a rough way. By combining the feature maps from this branch with the backbone network and introducing the character attention module, a more accurate moiré-free image is generated. Notably, we follow [48] as the backbone, in which the original resolution subnetwork (ORS-Net), channel attention block (CAB), and supervised attention module (SAM) are utilized.

4.1 Rough moiré pattern extraction module

According to [17], moiré patterns are mainly shaped in curves and stripes, which benefit from their specific properties. Obviously, extracting these properties that are different from those in the reference image would help to remove the moiré patterns. Fortunately, similar to [17], we statistically find that by decomposing the contaminated image into R, G, and B (red, green, and blue) channels, the G channel encounters much slighter moiré patterns than those in the R and B channels, as displayed in Figure 5. Of course, subtraction between the G channel and the R/B channel is a simple way to roughly obtain moiré-associated information for image restoration. However, despite the fact that different channels suffer from different moiré patterns, they also exhibit different scales of values. In other words, it is possible that one channel may have much larger or smaller values than that in the remaining one or two channels, subsequently making the aforementioned channel subtraction strategy useless. In order to tackle this problem, in this study, we introduce a learnable strategy through which the differences in value scales are adaptively alleviated.

Mathematically, let the contaminated image be $\mathbf{I}_m \in \mathbb{R}^{H \times W \times 3}$, where H and W denote the height and width of \mathbf{I}_m , respectively. By decomposing \mathbf{I}_m into the three channels, we can obtain $\mathbf{I}_m^r \in \mathbb{R}^{H \times W \times 1}$, $\mathbf{I}_m^g \in \mathbb{R}^{H \times W \times 1}$, and $\mathbf{I}_m^b \in \mathbb{R}^{H \times W \times 1}$ corresponding to the R, G, and B channels, respectively. By forwarding these three inputs into their associated convolution blocks, we can obtain

$$\mathbf{I}_e^r = \text{Conv}^r(\mathbf{I}_m^r), \mathbf{I}_e^g = \text{Conv}^g(\mathbf{I}_m^g), \mathbf{I}_e^b = \text{Conv}^b(\mathbf{I}_m^b), \quad (1)$$

where Conv^r , Conv^g , and Conv^b are the convolution blocks and $\mathbf{I}_e^r/\mathbf{I}_e^g/\mathbf{I}_e^b \in \mathbb{R}^{H \times W \times 1}$. The moiré patterns can then be roughly extracted through

$$\mathbf{M}^r = \mathbf{I}_e^r - \mathbf{I}_e^g, \mathbf{M}^b = \mathbf{I}_e^b - \mathbf{I}_e^g, \quad (2)$$

where \mathbf{M}^r and \mathbf{M}^b are both extracted features associated with the moiré patterns. In Equation 1, the scales of values for different channels are adaptively transformed to a consistent subspace, which are adaptively tuned through a task-driven strategy, so that the moiré patterns can be roughly extracted and contribute to moiré-free image generation. As shown in Figure 5, it is easy to observe that our presented technique indeed achieves the superiority.

In addition, since the edges are also an additional prior for moiré-contaminated images, we further apply the Sobel operator [49] to enhance the edge information of three channels, as shown in the second row of Figure 6. Similar to Equation 2, these edge maps associated with the “G” channel are subtracted from the other two maps via Equation 3.

$$\mathbf{M}_e^r = \mathbf{E}^r - \mathbf{E}^g, \mathbf{M}_e^b = \mathbf{E}^b - \mathbf{E}^g, \quad (3)$$

where $\mathbf{E}^r = \text{Sobel}(\mathbf{I}_m^r) \in \mathbb{R}^{H \times W \times 1}$, $\mathbf{E}^g = \text{Sobel}(\mathbf{I}_m^g) \in \mathbb{R}^{H \times W \times 1}$, and $\mathbf{E}^b = \text{Sobel}(\mathbf{I}_m^b) \in \mathbb{R}^{H \times W \times 1}$.

After obtaining \mathbf{M}^r , \mathbf{M}^b , \mathbf{M}_e^r , and \mathbf{M}_e^b , we then concatenate them as a single input, which is combined with three channel inputs. As displayed in the middle part of Figure 4, the concatenated inputs are forwarded into their corresponding convolution block and U-Net-like network. By further making a concatenation and taking the raw image \mathbf{I}_m into account again, the output $\mathbf{I}_o \in \mathbb{R}^{H \times W \times 3}$ is finally obtained through the supervised attention module [48]. By introducing the Charbonnier loss [50], \mathbf{I}_o is obtained in a supervised way, as defined in Equation 4:

$$L_0 = \sqrt{\|\mathbf{I}_o - \mathbf{I}_{gr}\|^2 + \varepsilon^2}, \quad (4)$$

where the constant ε is empirically set to 10^{-3} and \mathbf{I}_{gr} is the ground-truth image (we will analyze it in the following Subsection 4.2).

4.2 Character attention module

Different from the natural image-based restoration which focuses on all pixels equally, the purpose of our task is to increase the recognition accuracy after demoiréing. In other words, we focus on the characters rather than the surrounding background pixels. In fact, as shown in Figure 2, some images indeed include quite complex backgrounds, such as watermarking and diverse colors. Strongly enforcing the inputs to be the same to these reference images with complex backgrounds are impossible. Therefore, in this

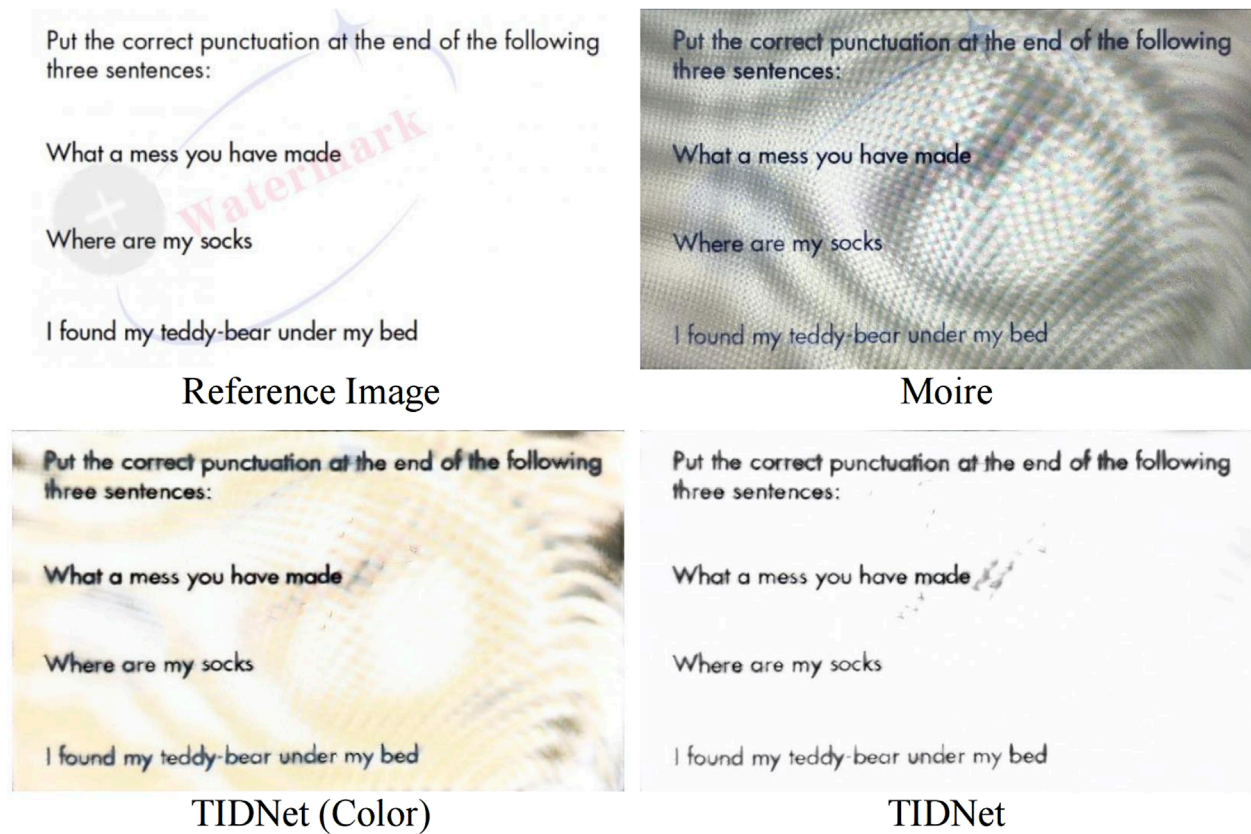


FIGURE 9 Displays of demoiréd images obtained by our TIDNet when the reference image and ground-truth image are, respectively, used as the supervised guidance.

paper, we first transform the reference image I_{ref} into a binary-like version I_{gt} , where the pixel values of characters are all close to 0, while others are all close to 255, as displayed in Figure 7. The more remarkable the characters, the greater the performance. Thanks to the generation of the image I_{gt} , we can just transform the contaminated images into a consistent style version no matter whether inputs encounter diverse backgrounds, but we can also increase the difference between the characters and the background, contributing to a more accurate text detection and recognition.

By forwarding I_m through the backbone network, we can formulate the residual output O_{res} [11] guided by I_{gt} , which is shown as follows:

$$L_b = \sqrt{\|I_m + O_{res} - I_{gt}\|^2 + \varepsilon^2}. \quad (5)$$

Equation 5 is used to encourage the reconstructed image to be similar to the ground-truth at the pixel level. Notably, feature maps obtained from SAM are also introduced into this O_{res} -related branch. For I_o , which is enforced to be similar to the ground-truth image I_{gt} , the feature maps from SAM would be beneficial for estimating O_{res} .

In addition, to further allow our model to pay much more attention on the character-associated pixels, we regard

I_{gt} as the mask for the text image enhancement, which can be formulated as Equation 6:

$$L_m = (1 - I_{gt}) \odot \sqrt{\|I_m + O_{res} - I_{gt}\|^2 + \varepsilon^2}. \quad (6)$$

Of course, images restored from the contaminated image should be easily recognized by an OCR model. Therefore, to enforce the recovered text images to exhibit their corresponding semantic priors, a text semantic loss is further introduced. Particularly, CRNN [51] followed by its pre-trained model is exploited. In this study, we use L_{ocr} to denote the semantic evaluation on the recovered image, as defined in Equation 7:

$$L_{ocr} = \text{OCR}(\text{CRNN}(I_m + O_{res}), \text{text}_{gt}), \quad (7)$$

where text_{gt} refers to the ground-truth of text information. Notably, the weights in CRNN are fixed, and the gradient would be transported to our designed network for model learning.

Taking the aforementioned analysis into account, the objective function of our proposed method is formulated as Equation 8:

$$L = \gamma L_m + \beta L_b + \lambda L_{ocr} + \eta L_o, \quad (8)$$

where γ , β , λ , and η are non-zero parameters to trade-off these four terms.

TABLE 4 Quantitative results on our collected dataset obtained by different comparison methods and TIDNet.

Task	Detection		
Metrics	Recall	Precision	F1-measure
AFN	22.88%	99.23%	37.19%
WDNet	73.82%	99.10%	84.61%
C3Net	27.97%	99.17%	43.64%
DnCNN	22.86%	99.23%	37.16%
FFDNet	43.94%	98.97%	60.86%
TIDNet	92.92%	98.72%	95.73%
Task	Recognition		
AFN	23.09%	80.13%	35.85%
WDNet	70.79%	91.80%	79.94%
C3Net	26.42%	79.11%	39.62%
DnCNN	25.14%	82.33%	38.52%
FFDNet	38.10%	79.49%	51.51%
TIDNet	85.34%	93.78%	89.36%

4.3 Implementation details

We implement our TIDNet using PyTorch [52]. The model runs on two GPUs of NVIDIA RTX 3090 with CUDA version 11.2. Except the OCR-related network CRNN, we optimize our network through the Adam optimizer with the learning rate of 2×10^{-4} . In this study, we set the maximum of epochs to 50. The learning rate is gradually reduced by following cosine annealing, and the minimum of our learning rate is 1×10^{-4} . In addition, the input image is resized to 256×256 , and the batch size is set to 12. Empirically, we first remove L_m in the first 40 epochs, after which it is exploited. Referring to the parameters γ , β , λ , and η , we empirically set them to 0.85, 0.5, 0.001, and 0.5, respectively.

5 Experiments

To demonstrate the significance of text image demoiréing and effectiveness of our proposed TIDNet, experiments are conducted on our collected dataset. In this section, the experimental settings and evaluation metrics are first described. We then conducted ablation studies to substantiate the importance of our introduced strategies. Finally, our proposed method is compared with state of the arts to further show its superiority.

5.1 Experimental settings and evaluation metrics

In this study, we divide the dataset into two subsets: one for training and another for testing. Specifically, 3,627 pairs are regarded as the training set and 112 contaminated images are used as the testing set. Notably, in testing images, there are totally 43,152 characters.

Since the final purpose of our TIDNet is to improve the OCR performance, we introduce recall, precision, and F1-measure (F1-m) scores as the quantitative evaluations for both text detection and recognition. Recall is the ratio between the number of correctly predicted characters and the number of labeled characters. It indicates how many items are correctly identified. Correspondingly, precision is the ratio between the number of correctly predicted characters and the number of all predicted characters. F1-m is a metric define by the recall score and the precision score: $\frac{Recall \times Precision}{Recall + Precision}$.

Notably, most existing methods such as natural image demoiréing and image denoising adopted the widely used quantitative evaluations: peak signal-to-noise ratio (PSNR) and structural similarity (SSIM). However, they are not suitable for our task. In our text image demoiréing, the contaminated images are enforced to be close to the binary-like ground-truths, while these guided references usually encounter imbalanced numbers of foreground and background pixels. Generally, the background pixels cover much more areas than foreground pixels. Due to this constraint, PSNR or SSIM values would be inaccurate if some character-related pixels are erased but the background is clear. In other words, the erased pixels do not make an obvious influence on PSNR or SSIM. By contrast, the detection and recognition performances of images suffered from erased characters would be remarkably influenced. Thus, in this paper, recall, precision, and F1-m are more reasonable for our task.

5.2 Ablation study

5.2.1 Is text image demoiréing necessary?

Due to the contamination of moiré patterns, it would be difficult to detect and recognize characters from the text image. As tabulated in Table 2, metrics of recall and F1-m on the contaminated images are only (53.82% and 69.77%) and (56.32% and 70.18%) for detection and recognition, respectively. However, thanks to our proposed TIDNet, these two metrics exhibit dramatic enhancement, which are (92.92% and 95.73%) and (85.34% and 89.36%). Obviously, it is quite significant for text image demoiréing.

5.2.2 Do the rough moiré extraction module and three-channel network work?

Inspired by the specific property of moiré patterns, the rough moiré extraction module is first introduced to extract the edge information related to the moiré patterns, which is then followed by our three-channel network. In detail, Table 2 shows that the 3-channel network leads to significant improvements in recall and F1-measure. By further taking the rough moiré extraction module into account, performance continues to increase.

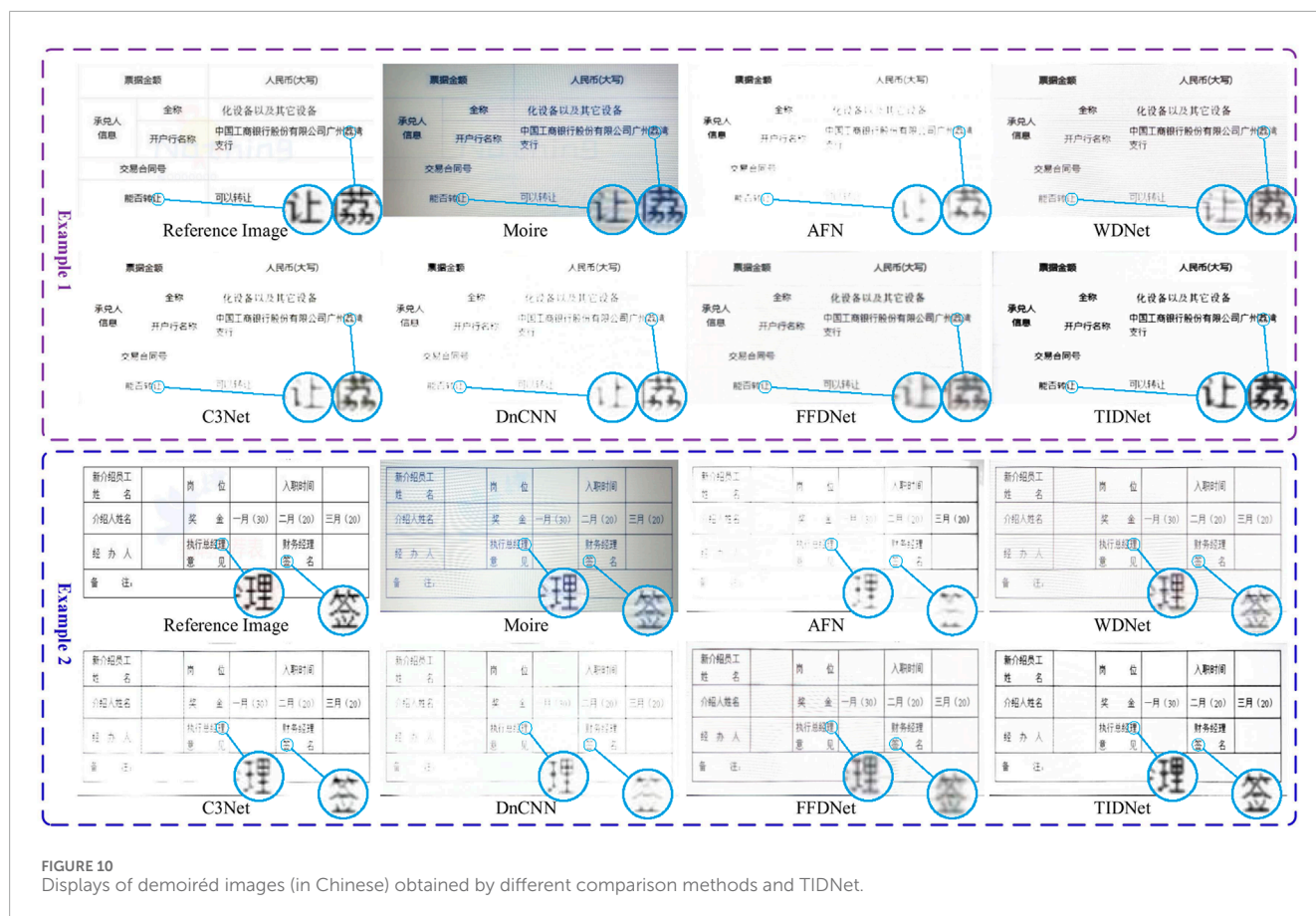


FIGURE 10 Displays of demoiré images (in Chinese) obtained by different comparison methods and TIDNet.

5.2.3 Does the character attention module work?

To enforce the network to highly focus on our interested characters, the mask loss L_m and OCR loss L_{ocr} are introduced into our proposed method. As listed in Table 2, these two losses significantly contribute to the performance improvement on both detection and recognition, exhibiting approximately 20%–35% increase. Thus, particularly focusing on the character-related pixels and exploiting their semantic information are quite significant. Notably, when only L_b is utilized, experimental results are even inferior to those obtained from the raw data. Generally, character-associated pixels cover much less areas compared with background pixels, while L_b equally pays attention on each pixel. In this case, even if the character estimation is incorrect, the influence on L_b may be slight, rendering it worthless. Fortunately, by exploiting the mask loss and OCR loss, the importance on characters are then enhanced.

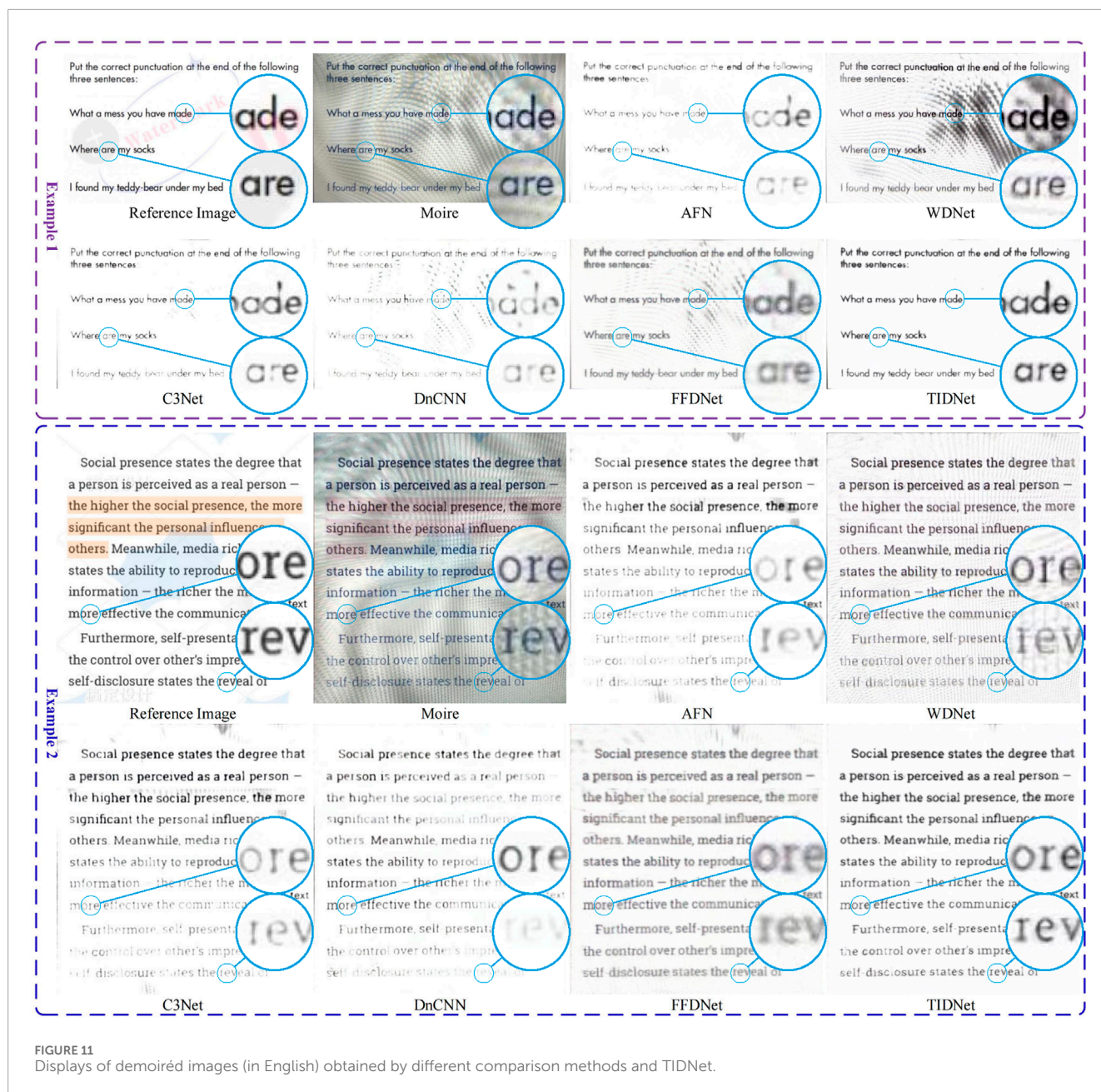
Figure 8 displays visualizations corresponding to Table 2. It is clear that when the backbone is applied, it removes the moiré patterns. However, it also regards optical characters as the moiré, which only makes the recovered image smooth and erases many character-related pixels. By contrast, thanks to our mask loss, the model highly focuses on characters, enhancing their associated pixels, as shown in the first image in the second row. Nevertheless, as character-related pixels are quite similar to some edge information, which also exists in the moiré patterns, the reconstructed image is also contaminated by moiré patterns. Thus, we further introduce the three-channel-based strategy followed our rough moiré extraction module (RMEM). Obviously, not just the moiré patterns are

alleviated, but the backgrounds are also much closer to the ground-truth compared with those obtained by “B” and “B + M.” Despite the fact that “B + M + C + RMEM” jointly enhances characters and removes moiré patterns, some recovered characters, as displayed in the enlarged details, encounter inaccurate semantic information. Fortunately, thanks to our introduced OCR loss, characters are further restored according to their semantics.

5.2.4 Does the binary-like ground-truth work?

In our proposed method, the reference image I_{ref} with diverse backgrounds is transformed to the ground-truth image I_{gt} , which is binary-like. In this way, the difference between foreground and background pixels is remarkably enlarged, allowing the network to more easily detect and recognize characters or texts. The comparison by using I_{ref} or I_{gt} as the guidance is shown in Table 3, proving the aforementioned analysis.

In addition, Figure 9 further proves the significance of using the binary-like ground-truth image I_{gt} as the guidance instead of the reference image I_{ref} . Generally, I_{ref} is corrupted with complex backgrounds such as colors and watermarking. In addition, the contaminated image may miss the information in the data collection, as shown in “Moiré” in Figure 9. Strongly enforcing the input to be identical to I_{ref} is too strict to achieve. As displayed in “TIDNet (color)” in Figure 9, the background of this recovered image is not just significantly different from I_{ref} , but it also still contains some moiré pattern-related contaminations. By contrast, due to the



consistent style of the ground-truth image, our TIDNet successfully achieve much better visualization under its guidance.

5.3 Comparison with state of the arts

To further demonstrate the effectiveness of our proposed method on the moiré pattern removal, we conducted experiments compared with state of the arts, including AFN [28], WDNNet [27], C3Net [29], DnCNN [11], and FFDNet [53]. Specifically, the first three methods are designed for image demoiréing, and the last two are designed for image restoration. To make a fair comparison, we retrain them on our collected dataset according to their released source codes.

The quantitative results on detection and recognition are tabulated in Table 4. Obviously, our presented method TIDNet dramatically outperforms these state of the arts. Compared with AFN, C3Net, and DnCNN, our achieved results are much superior to those computed by them. Specifically, they are all less than 30% and 45%, respectively, on the recall and F1-measure in text detection, whereas TIDNet achieves more than 50% improvement. Referring to FFDNet, although it is slightly better than the aforementioned method, it is still much inferior to TIDNet. In comparison to WDNNet, our proposed method also achieves noticeable performance enhancement.

The comparison visualizations in Chinese and English are, respectively, shown in Figures 10, 11. It is easy to observe that no matter whether the text images are in Chinese or English, our

presented method exhibits much better visualizations compared with existing image demoiréing and image restoration methods. Referring to AFN and C3Net, although moiré patterns are removed from the contaminated images, many character-related pixels are also erased, significantly making an inferior influence on text detection and recognition. The main reason is that these two methods regard the characters as moiré patterns since they have similar attributes. By contrast, WDNNet overcomes this problem, however its recovered images are still corrupted by more or less moiré patterns. For DnCNN, it also suffers from the similar problem compared with AFN and C3Net. Although a better visualization is obtained by FFDNet in comparison to DnCNN, its reconstructed characters are blurred. Different from these comparison approaches, our proposed method not only efficiently erases moiré patterns but also restores the characters which are quite similar to the ground-truth.

6 Conclusion

To fill the gap between the OCR and image demoiréing, in this paper, a text image-based dataset is primarily collected for text image demoiréing, allowing for supervised study. Furthermore, we propose a novel network named TIDNet, which is particularly adaptive for text image demoiréing. Inspired by the specific priors of moiré patterns, a rough moiré extraction module followed by a three-channel network is introduced so that the moiré pattern-associated information is easily extracted. Since our purpose is to improve the detection and recognition performance, a character attention module is also proposed in our TIDNet, through which the network highly pays attention on character-associated pixels and their semantic information. As a result of the aforementioned strategies, our proposed method enjoys a dramatic performance improvement on the OCR application.

Data availability statement

The raw data supporting the conclusions of this article will be made available by the authors, without undue reservation.

References

1. Mori S, Nishida H, Yamada H. *Optical character recognition*. John Wiley and Sons, Inc. (1999).
2. Kim MD, Ueda J. Dynamics-based motion deblurring improves the performance of optical character recognition during fast scanning of a robotic eye. *IEEE/ASME Trans Mechatronics* (2018) 23:491–5. doi:10.1109/tmech.2018.2791473
3. Shi X, Shen X. Oracle recognition of oracle network based on ant colony algorithm. *Front Phys* (2021) 9:768336. doi:10.3389/fphy.2021.768336
4. Guo X, Li J, Chen B, Lu G. Mask-most net: mask approximation based multi-oriented scene text detection network. In: 2019 IEEE International Conference on Multimedia and Expo (ICME); 08–12 July 2019; Shanghai, China. IEEE (2019) p. 206–11.
5. Ding H, Du Z, Wang Z, Xue J, Wei Z, Yang K, et al. Intervoxnet: a novel dual-modal audio-text fusion network for automatic and efficient depression

Author contributions

ZZ: conceptualization, data curation, formal analysis, investigation, methodology, software, validation, writing–original draft, writing–review and editing, and visualization. BL: conceptualization, formal analysis, investigation, methodology, software, and writing–original draft. TR: data curation, formal analysis, investigation, methodology, and writing–original draft. CF: data curation, investigation, software, and writing–original draft. RL: data curation, software, and writing–original draft. ML: funding acquisition, project administration, resources, supervision, writing–original draft, and writing–review and editing.

Funding

The author(s) declare that financial support was received for the research, authorship, and/or publication of this article. This project was supported in part by the National Natural Scientific Foundation of China 62472124, Shenzhen Colleges and Universities Stable Support Program GXWD20220811170130002.

Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Generative AI statement

The author(s) declare that no Generative AI was used in the creation of this manuscript.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors, and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

detection from interviews. *Front Phys* (2024) 12:1430035. doi:10.3389/fphy.2024.1430035

6. Zhan F, Lu S. Esir: end-to-end scene text recognition via iterative image rectification. In: Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition. 16–20 June 2019; Long Beach, CA, United States: IEEE (2019) p. 2059–68.

7. Satyawar W, Pratama MO, Jannati R, Muhammad G, Fajar B, Hamzah H, et al. Citizen id card detection using image processing and optical character recognition. In: Journal of physics: Conference series, 1235. Bristol, United Kingdom: IOP Publishing (2019), 012049.

8. Schreiber S, Agne S, Wolf I, Dengel A, Ahmed S. Deepdesrt: deep learning for detection and structure recognition of tables in document images. In: 2017 14th

IAPR international conference on document analysis and recognition (ICDAR); 09–15 November 2017; Kyoto, Japan, 1. IEEE (2017) p. 1162–7. doi:10.1109/icdar.2017.192

9. Zhuang J, Hou S, Wang Z, Zha ZJ. Towards human-level license plate recognition. In: Proceedings of the European Conference on computer vision. 8–14 September 2018. Munich, Germany: ECCV (2018) p. 306–21.

10. Zhu Z, Wang Z, Qi G, Mazur N, Yang P, Liu Y. Brain tumor segmentation in mri with multi-modality spatial information enhancement and boundary shape correction. *Pattern Recognition* (2024) 153:110553. doi:10.1016/j.patcog.2024.110553

11. Zhang K, Zuo W, Chen Y, Meng D, Zhang L. Beyond a Gaussian denoiser: residual learning of deep cnn for image denoising. *IEEE Trans Image Process* (2017) 26:3142–55. doi:10.1109/tip.2017.2662206

12. Jiang B, Lu Y, Wang J, Lu G, Zhang D. Deep image denoising with adaptive priors. *IEEE Trans Circuits Syst Video Tech* (2022) 32:5124–36. doi:10.1109/TCSVT.2022.3149518

13. Ren D, Zuo W, Hu Q, Zhu P, Meng D. Progressive image deraining networks: a better and simpler baseline. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition; 15–20 June 2019; Long Beach, CA, USA (2019) p. 3937–46.

14. Zhang H, Sindagi V, Patel VM. Image de-raining using a conditional generative adversarial network. *IEEE Trans Circuits Syst Video Tech* (2020) 30:3943–56. doi:10.1109/TCSVT.2019.2920407

15. Qu Y, Chen Y, Huang J, Xie Y. Enhanced pix2pix dehazing network. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2019) p. 8160–8.

16. Wang P, Zhu H, Huang H, Zhang H, Wang N. Tms-gan: a twofold multi-scale generative adversarial network for single image dehazing. *IEEE Trans Circuits Syst Video Tech* (2022) 32:2760–72. doi:10.1109/TCSVT.2021.3097713

17. He B, Wang C, Shi B, Duan LY. Mop moiré patterns using mopnet. *Proc IEEE/CVF Int Conf Comput Vis* (2019) 2424–32.

18. Sun Y, Yu Y, Wang W. Moiré photo restoration using multiresolution convolutional neural networks. *IEEE Trans Image Process* (2018) 27:4160–72. doi:10.1109/tip.2018.2834737

19. Liu S, Li C, Nan N, Zong Z, Song R. Mmdm: multi-frame and multi-scale for image demoiré. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops; 14–19 June 2020; Seattle, WA (2020) p. 434–5.

20. Zheng B, Yuan S, Slabaugh G, Leonardis A. Image demoiré with learnable bandpass filters. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition; 13–19 June 2020; Seattle, WA (2020) p. 3636–45.

21. Gao T, Guo Y, Zheng X, Wang Q, Luo X. Moiré pattern removal with multi-scale feature enhancing network. In: 2019 IEEE International Conference on Multimedia and Expo Workshops (ICMEW); 08–12 July 2019; Shanghai, China. IEEE (2019) p. 240–5.

22. Qi W, Yu X, Li X, Kang S. A moiré removal method based on peak filtering and image enhancement. *Mathematics* (2024) 12:846. doi:10.3390/math12060846

23. Liu F, Yang J, Yue H. Moiré pattern removal from texture images via low-rank and sparse matrix decomposition. In: 2015 Visual Communications and Image Processing (VCIP); 13–16 December 2015; Singapore: IEEE (2015) p. 1–4.

24. Liu B, Shu X, Wu X. Demoiré of camera-captured screen images using deep convolutional neural network (2018) arXiv preprint arXiv:1804.03809.

25. Yue H, Mao Y, Liang L, Xu H, Hou C, Yang J. Recaptured screen image demoiré. *IEEE Trans Circuits Syst Video Tech* (2020) 31:49–60. doi:10.1109/tcsvt.2020.2969984

26. Luo X, Zhang J, Hong M, Qu Y, Xie Y, Li C. Deep wavelet network with domain adaptation for single image demoiré. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops; 14–19 June 2020; Seattle, WA, USA (2020) p. 420–1.

27. Liu L, Liu J, Yuan S, Slabaugh G, Leonardis A, Zhou W, et al. Wavelet-based dual-branch network for image demoiré. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIII 16. Springer (2020) p. 86–102.

28. Xu D, Chu Y, Sun Q. Moiré pattern removal via attentive fractal network. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops; 14–19 June 2020; Seattle, WA (2020) p. 472–3.

29. Kim S, Nam H, Kim J, Jeong J. C3net: demoiré network attentive in channel, color and concatenation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops; 14–19 June 2020; Seattle, WA (2020) p. 426–7.

30. Mancas-Thillou C, Mirmehdi M. An introduction to super-resolution text. In: *Digital document processing*. Springer (2007) p. 305–27.

31. Dong C, Loy CC, He K, Tang X. Image super-resolution using deep convolutional networks. *IEEE Trans Pattern Anal Machine Intelligence* (2015) 38:295–307. doi:10.1109/tpami.2015.2439281

32. Dong C, Zhu X, Deng Y, Loy CC, Qiao Y. Boosting optical character recognition: a super-resolution approach (2015) arXiv preprint arXiv:1506.02211.

33. Wang W, Xie E, Liu X, Wang W, Liang D, Shen C, et al. Scene text image super-resolution in the wild. In: European Conference on Computer Vision. 2020: 16th European Conference. August 23–28. Glasgow, United Kingdom: Springer (2020) p. 650–66.

34. Lin K, Liu Y, Li TH, Liu S, Li G. Text image super-resolution by image matting and text label supervision. In: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW); 16–17 June 2019; Long Beach, CA, USA. IEEE (2019) p. 1722–7.

35. Chen J, Li B, Xue X. Scene text telescope: text-focused scene image super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition; 20–25 June 2021; Nashville, TN, USA (2021) p. 12026–35.

36. Mei J, Wu Z, Chen X, Qiao Y, Ding H, Jiang X. Deepdeblur: text image recovery from blur to sharp. *Multimedia Tools Appl* (2019) 78:18869–85. doi:10.1007/s11042-019-7251-y

37. Cho H, Wang J, Lee S. Text image deblurring using text-specific properties. In: European Conference on Computer Vision. October 7–13. Florence, Italy: Springer (2012) p. 524–37.

38. Jiang X, Yao H, Zhao S. Text image deblurring via two-tone prior. *Neurocomputing* (2017) 242:1–14. doi:10.1016/j.neucom.2017.01.080

39. Lee H, Jung C, Kim C. Blind deblurring of text images using a text-specific hybrid dictionary. *IEEE Trans Image Process* (2019) 29:710–23. doi:10.1109/tip.2019.2933739

40. Li J, Guo X, Lu G, Zhang B, Xu Y, Wu F, et al. Drpl: deep regression pair learning for multi-focus image fusion. *IEEE Trans Image Process* (2020) 29:4816–31. doi:10.1109/tip.2020.2976190

41. Li J, Liang B, Lu X, Li M, Lu G, Xu Y. From global to local: multi-patch and multi-scale contrastive similarity learning for unsupervised defocus blur detection. *IEEE Trans Image Process* (2023) 32:1158–69. doi:10.1109/tip.2023.3240856

42. Li J, Fan D, Yang L, Gu S, Lu G, Xu Y, et al. Layer-output guided complementary attention learning for image defocus blur detection. *IEEE Trans Image Process* (2021) 30:3748–63. doi:10.1109/tip.2021.3065171

43. Liu Y, Qi Z, Cheng J, Chen X. Rethinking the effectiveness of objective evaluation metrics in multi-focus image fusion: a statistic-based approach. *IEEE Trans Pattern Anal Machine Intelligence* (2024) 46:5806–19. doi:10.1109/tpami.2024.3367905

44. Li H, Liu J, Zhang Y, Liu Y. A deep learning framework for infrared and visible image fusion without strict registration. *Int J Comput Vis* (2024) 132:1625–44. doi:10.1007/s11263-023-01948-x

45. Raghunandan KS, Shivakumara P, Jalab HA, Ibrahim RW, Kumar GH, Pal U, et al. Riesz fractional based model for enhancing license plate detection and recognition. *IEEE Trans Circuits Syst Video Tech* (2018) 28:2276–88. doi:10.1109/TCSVT.2017.2713806

46. Karthikeyan S, de Herrera AGS, Doctor F, Mirza A. An ocr post-correction approach using deep learning for processing medical reports. *IEEE Trans Circuits Syst Video Tech* (2022) 32:2574–81. doi:10.1109/TCSVT.2021.3087641

47. Guo Y, Ji C, Zheng X, Wang Q, Luo X. Multi-scale multi-attention network for moiré document image binarization. *Signal Processing: Image Commun* (2021) 90:116046. doi:10.1016/j.image.2020.116046

48. Zamir SW, Arora A, Khan S, Hayat M, Khan FS, Yang MH, et al. Multi-stage progressive image restoration. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition; 20–25 June 2021; Nashville, TN, USA (2021) p. 14821–31.

49. Kanopoulos N, Vasanthavada N, Baker RL. Design of an image edge detection filter using the sobel operator. *IEEE J solid-state circuits* (1988) 23:358–67. doi:10.1109/4.996

50. Charbonnier P, Blanc-Feraud L, Aubert G, Barlaud M. Two deterministic half-quadratic regularization algorithms for computed imaging. In: Proceedings of 1st International Conference on Image Processing; 13–16 November 1994; Austin, TX, USA, 2. IEEE (1994) p. 168–72.

51. Shi B, Bai X, Yao C. An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. *IEEE Trans Pattern Anal Machine Intelligence* (2016) 39:2298–304. doi:10.1109/tpami.2016.2646371

52. Paszke A, Gross S, Chintala S, Chanan G, Yang E, DeVito Z, et al. *Automatic differentiation in pytorch* (2017).

53. Zhang K, Zuo W, Zhang L. Ffdnet: toward a fast and flexible solution for cnn-based image denoising. *IEEE Trans Image Process* (2018) 27:4608–22. doi:10.1109/tip.2018.2839891



OPEN ACCESS

EDITED BY

Bo Xiao,
Imperial College London, United Kingdom

REVIEWED BY

Gang Hu,
Buffalo State College, United States
Yafei Zhang,
Kunming University of Science and
Technology, China
Yimin Chen,
University of Massachusetts Lowell,
United States

*CORRESPONDENCE

Aochen Yan,
✉ aochenya@usc.edu

RECEIVED 22 November 2024

ACCEPTED 05 December 2024

PUBLISHED 20 December 2024

CITATION

Li Z, Wang H, Chen H, Lin C and Yan A (2024)
Multi-Conv attention network for skin lesion
image segmentation.
Front. Phys. 12:1532638.
doi: 10.3389/fphy.2024.1532638

COPYRIGHT

© 2024 Li, Wang, Chen, Lin and Yan. This is an
open-access article distributed under the
terms of the [Creative Commons Attribution
License \(CC BY\)](#). The use, distribution or
reproduction in other forums is permitted,
provided the original author(s) and the
copyright owner(s) are credited and that the
original publication in this journal is cited, in
accordance with accepted academic practice.
No use, distribution or reproduction is
permitted which does not comply with
these terms.

Multi-Conv attention network for skin lesion image segmentation

Zexin Li¹, Hanchen Wang¹, Haoyu Chen¹, Chenxin Lin¹ and
Aochen Yan^{2*}

¹International College, Chongqing University of Posts and Telecommunications, Chongqing, China,

²Viterbi School of Engineering, University of Southern California, Los Angeles, CA, United States

To address the trade-off between segmentation performance and model lightweighting in computer-aided skin lesion segmentation, this paper proposes a lightweight network architecture, Multi-Conv Attention Network (MCAN). The network consists of two key modules: ISDConv (Inception-Split Depth Convolution) and AEAM (Adaptive Enhanced Attention Module). ISDConv reduces computational complexity by decomposing large kernel depthwise convolutions into smaller kernel convolutions and unit mappings. The AEAM module leverages dimensional decoupling, lightweight multi-semantic guidance, and semantic discrepancy alleviation to facilitate the synergy between channel attention and spatial attention, further exploiting redundancy in the spatial and channel feature maps. With these improvements, the proposed method achieves a balance between segmentation performance and computational efficiency. Experimental results demonstrate that MCAN achieves state-of-the-art performance on mainstream skin lesion segmentation datasets, validating its effectiveness.

KEYWORDS

medical image segmentation, lightweight, melanoma, attention mechanism, Inception

Introduction

Melanoma, a highly malignant skin tumor, causes a significant number of deaths worldwide each year. Its incidence and mortality rates vary significantly depending on the region, the level of early diagnosis awareness, and the accessibility of primary care [1]. Early detection of melanoma is crucial for improving patient survival rates. However, due to the diversity and complexity of melanoma's appearance, its accurate diagnosis often relies on the experience and expertise of doctors, which somewhat limits the efficiency and accuracy of early diagnosis.

In melanoma diagnosis, image segmentation is a key step that precisely separates the lesion area from healthy skin, helping doctors identify the lesion's boundaries and assist in accurate diagnosis and treatment. Traditional segmentation methods rely heavily on complex preprocessing and manual feature extraction, making it difficult to handle the complexity of melanoma images. With the emergence of high-quality datasets, data-driven deep learning methods have rapidly gained popularity. Zhang et al. [2] proposed a novel framework that integrates multiple experts to jointly learn representations from diverse MRI modalities, effectively enhancing segmentation performance. Similarly, Li et al. [3] addressed challenges in brain tumor segmentation caused by missing modalities by utilizing a deformation-aware learning framework that reconstructs missing information, resulting in more reliable and accurate segmentation even in incomplete datasets. Among them, attention mechanisms, as an effective way to integrate local and global features,

help the model focus on the lesion areas. Dong et al. [4] enhanced the capability to capture feature information by dynamically allocating attention weights across channel and spatial dimensions, addressing the complex features, blurry boundaries, and noise interference in skin lesion segmentation. Similarly, the GL-CSAM module designed by Sun et al. [5] aims to capture global contextual information, enhancing the model's ability to perceive global features. However, they did not fully explore feature fusion between different convolutional layers. To address this issue, Qiu et al. [6] introduced a multi-level attention fusion mechanism that progressively extracts lesion boundary information using contextual information from different levels, alleviating the problem of blurry boundaries. Qi et al. [7] and Liu et al. [8] introduced single attention mechanisms to integrate contextual features, specifically designed for stroke lesion segmentation. The combination of standalone self-attention modules with convolutional layers has shown limited effectiveness in enhancing the model's non-local feature modeling capabilities. To address this limitation, Yang et al. [9] introduced a multi-attention mechanism (spatial and reverse attention). Spatial attention is used to improve the extraction of useful features, while reverse attention enhances the network's segmentation performance by applying reverse attention operations on skip connections, enabling more accurate analysis and localization of small lesion targets. Liu et al. [10] and Zhu et al. [11] enhanced the precision and detail of tumor segmentation by fusing information from multiple MRI modes such as T1, T2, and FLAIR. Zhu et al. [12] embedded a feature fusion module based on attention mechanism in the model structure to optimize the expression and integration of multi-modal features to improve segmentation accuracy. Liu et al. [13] examined the effectiveness of traditional objective evaluation indicators in the evaluation of image fusion results and proposed a statistic-based framework to compensate for the shortcomings of existing indicators. These methods have improved the segmentation task to varying degrees at different stages, achieving commendable results. However, their network designs do not fully consider how to effectively utilize spatial information, and they lack dedicated mechanisms to enhance and preserve spatial information. These shortcomings may result in suboptimal performance when handling spatial correlations.

Moreover, it is worth noting that while introducing high-quality attention mechanisms, the parameter count of the model increases, potentially compromising the real-time performance during deployment. Although high-quality attention mechanisms can enhance model performance, they are often accompanied by an increase in parameter count, which can negatively impact the real-time performance of model deployment [14]. In response to such problems, most researchers have based their efforts on the potential of deep separable convolution to improve model efficiency and effectiveness. Zhou et al. [15] constructs expansion layers using depthwise separable convolutions to efficiently extract multi-scale features with low computational overhead, enhancing the feature representation capability. Liu et al. [16], Ma et al. [17], and Feng et al. [18] adopted a similar approach by integrating depthwise separable convolution layers into the encoder. However, they often struggle to achieve precise detailed description while maintaining low computational overhead. Ruan et al. [19] combined MLP to extract global feature information, followed by feature extraction using depthwise separable convolutions (DWConv). This effectively

preserved significant features in the brain feature map while filtering out less relevant features. However, the lightweight processing of complex features remains limited. Similarly, Lei et al. [20] combined depthwise separable convolutions with bilinear interpolation to adjust the size of high-level features, making them match low-level features. However, this approach faces performance bottlenecks when further reducing the computational burden. Chen et al. [21] incorporated the advantages of asymmetric convolutions based on depthwise separable convolutions and designed an ultralight convolution module, further achieving the decoupling of spatial and channel dimensions. Existing methods still have limitations in lightweight design. Although different encoder designs effectively reduce computational load and ensure efficient feature extraction, they still lack precision in representing the blurry edges of skin lesions.

To address the contradiction between segmentation performance and lightweight design, this paper proposes a lightweight segmentation method. It aims to more accurately capture and segment the lesion area by leveraging channel and spatial redundancy, without increasing additional computational load. Specifically, the core of the segmentation framework is the Inception-Split ISDConv. Additionally, at the bridging layer stage, we introduce the AEAM, which combines the collaborative effects of spatial and channel attention with the feature calibration capabilities of the squeeze-and-excitation network. AEAM utilizes multi-scale depth-shared 1D convolutions to capture multi-semantic spatial information for each feature channel. It effectively integrates global contextual dependencies and multi-semantic spaces, while calculating channel similarity and contributions under the guidance of compressed spatial knowledge, thereby alleviating semantic differences in the spatial structure. Additionally, we introduce dynamic convolution in the encoder. Dynamic convolution dynamically aggregates multiple parallel convolution kernels based on input-relevant attention mechanisms. Assembling multiple convolution kernels is not only computationally efficient but also enhances representational capability due to the smaller size of the kernels.

The contributions of this paper can be summarized in the following three aspects:

1. In this study, a novel lightweight segmentation network named Multi-Conv Attention Network (MCAN) is proposed. It performs channel and spatial weighting on the spatial and channel redundancies in the feature map without increasing additional computational load, achieving an effect of information complementarity.
2. To address the unclear edges in skin lesions, this paper proposes ISDConv. This module performs multi-scale feature extraction using depthwise separable convolutions, multi-scale convolution kernels, and spatial and channel reconstruction convolutions. It reduces computational complexity and the number of parameters, thereby improving the model's feature representation capability while maintaining efficient feature extraction.
3. To address the insufficient utilization of redundancies in the spatial and channel feature maps, this paper proposes the Adaptive Enhanced Attention Module (AEAM). Through dimension decoupling, lightweight multi-semantic guidance,

and semantic discrepancy mitigation, AEAM achieves the collaborative effect between channel and spatial attention, enabling the model to capture and segment the lesion areas more accurately.

Related works

Attention mechanism

In the field of natural images, Li et al. [22] used a dual attention fusion module to effectively combine features from images from different sources, thereby enhancing the model's ability to focus on important regions. The attention mechanism can enhance the extraction of key features in infrared and visible images, making the fused images clearer and retaining more meaningful details [23]. In medical image segmentation, the attention mechanism is primarily used to guide the model's focus on the lesion areas in the image, assigning different weights to each pixel or feature, enhancing task-relevant features, and suppressing irrelevant background information. Huang et al. [24] prior convolutional attention mechanism that dynamically allocates attention weights across both channel and spatial dimensions. Shaker et al. [25] used a pair of mutually dependent branches based on spatial and channel attention to effectively learn discriminative features, improving the quality of segmentation masks. Fu et al. [26] used a Transformer-based spatial and channel attention module to extract global complementary information across different layers of the U-Net, which helps in learning detailed features at different scales. To address hair interference in dermoscopic images, Xiong et al. [27] proposed a multi-scale channel attention mechanism that enhances feature information and boundary awareness. Song et al. [28] argued that current popular attention mechanisms focus too much on external image features and lack research on latent features. They introduced an external-latent attention mechanism, using an entropy quantization method to summarize the distribution of latent contextual information. Similarly, Huang et al. [29] used Bi-Level Routing Attention in deep networks to discard irrelevant key-value pairs, achieving content-aware sparse attention for dispersed semantic information.

Network lightweighting

While pursuing high performance, researchers have also begun to focus on the lightweight and efficiency of medical image segmentation networks. Network structure design is one of the most popular approaches for lightweight optimization. Ma et al. [17] simplified the structure, reduced the number of parameters, and optimized the convolution operations, achieving a significant reduction in computational complexity and model size while maintaining segmentation accuracy. This enables the model to perform excellently even in resource-constrained environments, making it suitable for applications such as mobile healthcare and telemedicine. The UcUNet [30] network achieves lightweight and precise medical image segmentation by designing an efficient large-kernel U-shaped convolution module. This network leverages large-kernel convolutions to expand the receptive field while integrating

depthwise separable convolutions to reduce the computational cost, thereby maintaining high segmentation accuracy with efficient computation. Liu et al. [16] combines the lightweight characteristics of HarDNet with multi-attention mechanisms, enhancing the network's ability to capture key features and achieving more precise medical image segmentation. Sun et al. [31] introduces a contextual residual network, effectively integrating contextual information into the U-shaped network, enhancing the global understanding and stability of the segmentation. Nisa and Ismail [32] employs a dual-path structure with a ResNet encoder, combining ResNet's feature extraction capabilities with U-Net's segmentation advantages, offering an alternative effective solution for medical image segmentation. Zhao et al. [33] proposed a four-layer feature calibration branch based on an attention mechanism. The downsampling layer reduces the resolution of rectal cancer CT image feature maps to half of the original size, followed by pointwise convolution to enable interactions between channels. This method effectively expands the receptive field of subsequent convolutional layers and optimizes computational efficiency by reducing the cost of calculating spatial attention. Model compression, as another approach to simplifying network structures, removes structural redundancy while maintaining performance, making it more suitable for various applications in medical image analysis. Wang et al. [34] designed a sophisticated teacher network to learn multi-scale features, guiding a more lightweight student network to improve segmentation accuracy. Experiments showed that this method effectively acquires detailed morphological features of the brain from the teacher network. Hajabdollahi et al. [35] proposed a channel pruning algorithm for medical image segmentation tasks, which selects color channels during image processing and allows training of the target structure directly on the pre-selected key channels. However, these studies did not address how to utilize the redundancy effectively.

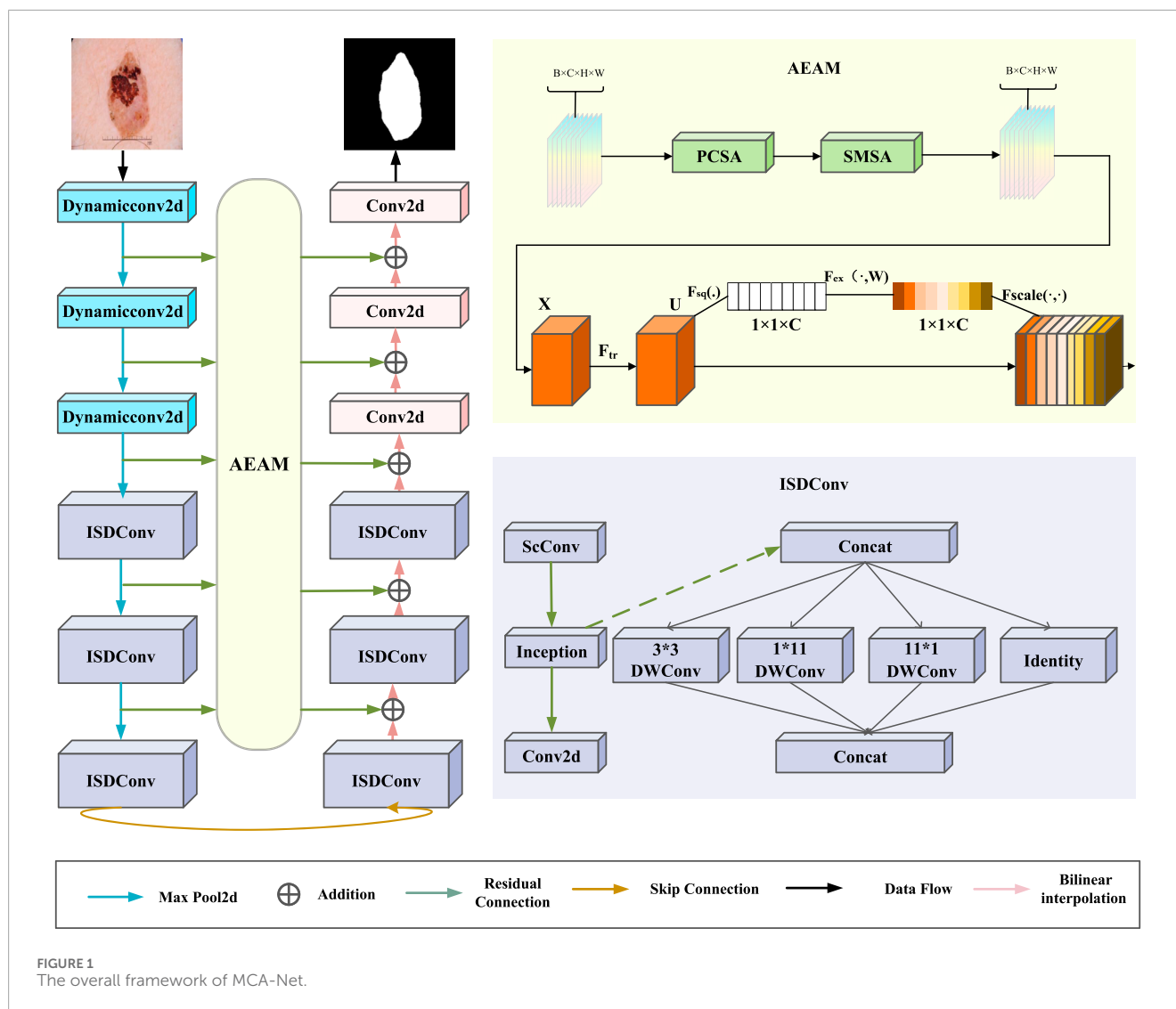
Based on the above research findings, this paper proposes a lightweight segmentation model that emphasizes spatial and channel features. This model improves segmentation accuracy and efficiency without increasing additional computational costs, providing a new and efficient solution for the medical imaging field.

Methods

The overall framework of MCA-Net

As illustrated in Figure 1, the proposed model framework consists primarily of the ISDConv module, the AEAM module, and dynamic convolution. The ISDConv module is composed of three parts: ScConv, Inception convolution, and standard convolution. By incorporating depthwise separable convolutions and group convolutions, ISDConv facilitates the model's understanding of multi-scale information within images, thereby enhancing its ability to detect and classify objects of varying sizes.

The AEAM module operates in two stages: SEattention and SCSA. SEattention enhances the network's representational capacity by explicitly modeling the interdependencies between convolutional feature channels. SCSA, in turn, is divided into two components:



SMSA and PCSA. SMSA integrates multi-semantic information and employs a progressive compression strategy to inject discriminative spatial priors into the channel self-attention mechanism of PCSA, effectively guiding channel recalibration. Within PCSA, robust feature interaction based on a self-attention mechanism further mitigates the multi-semantic information discrepancy among sub-features in SMSA.

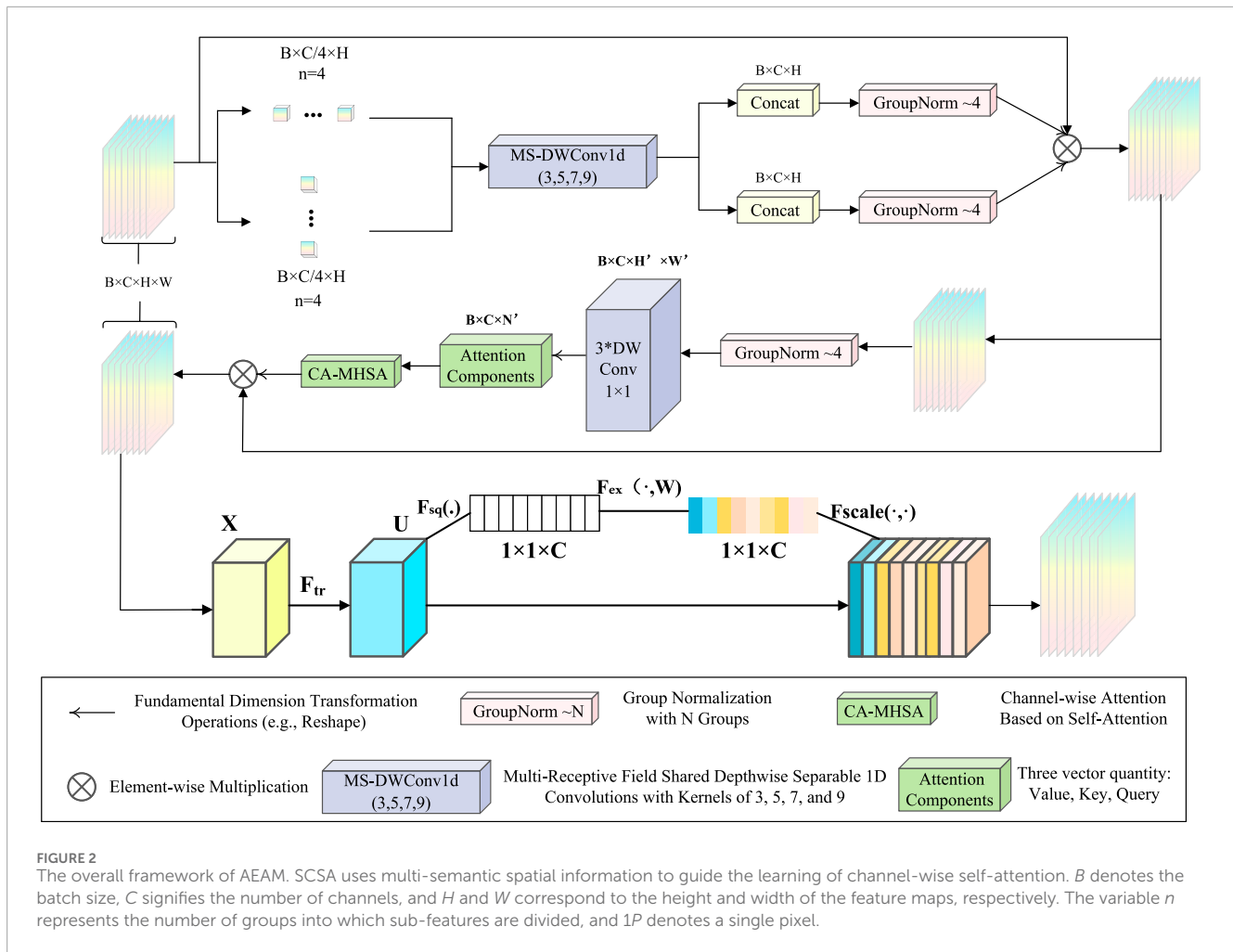
Inception-Split depth convolution

As shown in Figure 1, ISDCConv consists of a ScConv, an Inception Convolution, and a standard Conv2d layer. The Inception Convolution achieves lightweight performance by efficiently decomposing a large kernel depthwise convolution into four parallel branches along the channel dimension. These branches consist of a small square kernel, two orthogonal large kernels, and an identity mapping. The use of a small square kernel reduces computational complexity, while the orthogonal large

kernels capture different spatial information at varying scales. The identity mapping helps preserve the original input features, further enhancing the efficiency of the network. Additionally, this architecture incorporates 1×1 convolutions for dimensionality reduction before applying computationally expensive operations, minimizing the computational burden while preserving the model's ability to learn rich, multi-scale features. These four branches not only achieve higher computational efficiency than the large kernel depthwise convolution but also maintain a large receptive field, enabling the model to capture spatial context effectively for improved performance.

One of the branches employs a 3×3 kernel, which avoids the inefficiency of large square kernels. Instead, large square kernels $k_h \times k_w$ are decomposed into $1 \times k_w$ and $k_h \times 1$, significantly reducing computational complexity. Specifically, for a given input x , it is divided into four groups along the channel dimension, with the operation defined as Equation 1:

$$X_{hw}, X_w, X_h, X_{id} = \text{Split}(X) = X_{:,g}, X_{g:2g}, X_{2g:3g}, X_{3g:} \quad (1)$$



where, g represents the number of channels in each convolution branch, which is determined by the formula $g = r_g C$, where r_g is the ratio for splitting and C is the total number of input channels. The input is divided into four groups along the channel dimension based on this ratio, and the resulting split inputs are then fed into the respective parallel branches. Therefore, the following Equation 2 can be established:

$$\begin{aligned} X'_{hw} &= DWConv_{k_s \times k_s}^{g \rightarrow g}(X_{hw}) \\ X'_w &= DWConv_{1 \times k_b}^{g \rightarrow g}(X_w) \\ X'_h &= DWConv_{k_b \times 1}^{g \rightarrow g}(X_h) \\ X'_{id} &= X_{id} \end{aligned} \quad (2)$$

where k_s represents the 3×3 kernel size, k_b denotes the kernel sizes of 11×1 and 1×11 , X_{hw} represents the feature map, X_w refers to the features in the width direction, and X_h refers to the features in the height dimension of the image. After processing each input x_i through its respective branch, the outputs X' are concatenated along the channel dimension. The operation can be expressed as Equation 3.

$$X' = \text{Concat}(X'_{hw}, X'_w, X'_h, X'_{id}) \quad (3)$$

Adaptive Enhanced Attention Module

This paper introduces the AEAM attention module, designed to achieve synergy between channel attention and spatial attention through dimensional decoupling, lightweight multi-semantic guidance, and semantic discrepancy mitigation. As shown in Figure 2, the AEAM module consists of two main components: SEattention and SCSA.

The SCSA module is composed of two sequentially linked components: Shared Multi-Semantic Spatial Attention (SMSA) and Progressive Channel Self-Attention (PCSA). SMSA employs multi-scale, depth-sharing one-dimensional convolutions to extract spatial information at different semantic levels from four independent sub-features. This approach enables the efficient integration of diverse spatial semantics across sub-features. After SMSA modulates the feature maps, the resulting features are passed to PCSA. This component combines a progressive compression strategy with a channel-specific self-attention mechanism (CSA) to refine the feature representation further.

In this paper, a given input $X \in \mathbb{R}^{B \times C \times H \times W}$ is applied global average pooling along the height and width dimensions to create two unidirectional 1D sequence structures: $X_H \in \mathbb{R}^{B \times C \times W}$, $X_W \in \mathbb{R}^{B \times C \times H}$.

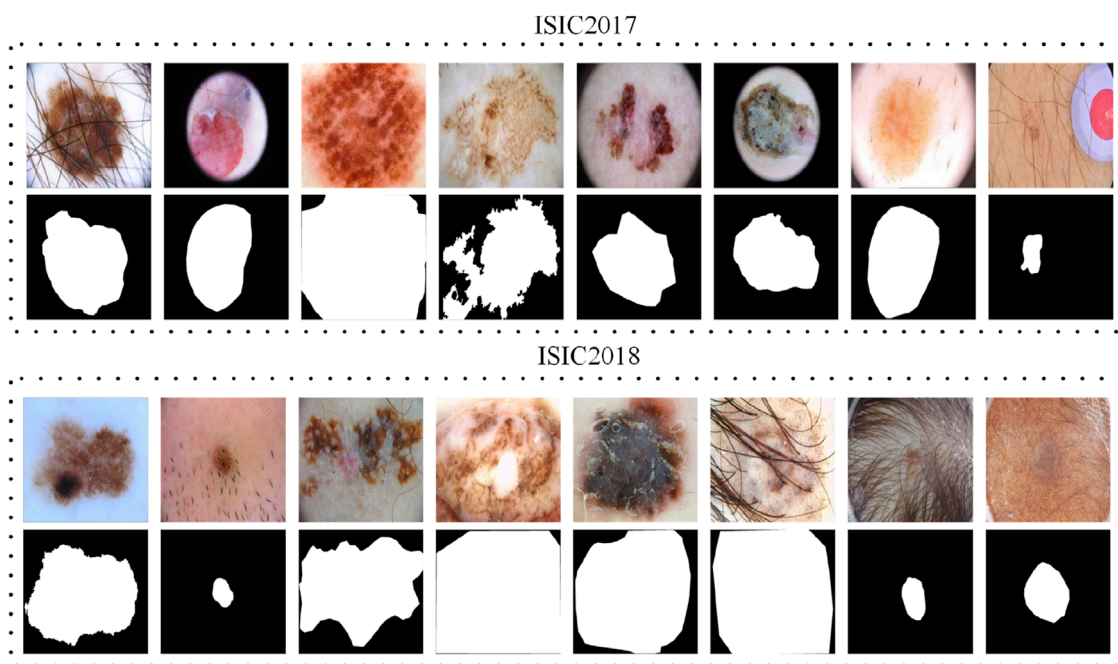


FIGURE 3
Examples of original images and their ground truth annotations from the ISIC2017 and ISIC2018 datasets.

To learn diverse spatial distributions and contextual relationships, the feature set is divided into K equally sized and independent sub-features, such that X_H^i and X_W^i , each sub-feature has a channel count of $\frac{C}{K}$, where C is the total number of channels in the original feature set. In this study, we set the default value $K = 4$, decomposing the features into H -dimensional and W -dimensional sub-features. During the decomposition process, 1D convolution is applied to each sub-feature. We employ lightweight shared convolutions for alignment, which implicitly model feature consistency across both dimensions by learning correlations.

The ablation formula is shown in Equation 4:

$$\begin{aligned}\tilde{X}_H^i &= DWConv1d_{\frac{C}{K} \rightarrow \frac{C}{K}}^{\frac{C}{K} \rightarrow \frac{C}{K}}(X_H^i) \\ \tilde{X}_W^i &= DWConv1d_{\frac{C}{K} \rightarrow \frac{C}{K}}^{\frac{C}{K} \rightarrow \frac{C}{K}}(X_W^i)\end{aligned}\quad (4)$$

Where X_H and X_W represent feature maps in height and width dimensions respectively. SEattention introduces the “Squeeze-and-Excitation” (SE) block, which enhances the network’s representational capacity by explicitly modeling the interdependencies between convolutional feature channels. The SE block employs a special mechanism that enables the network to perform feature recalibration. Through this mechanism, the block learns to selectively emphasize informative features while suppressing less useful ones by leveraging global information.

The structure of the SE block is illustrated in the lower part of Figure 2. For any given transformation F_{tr} , which maps the input X to a feature map U , which $U \in \mathbb{R}^{H \times W \times C}$, a corresponding SE block can be constructed to perform feature recalibration. The feature map U first undergoes a squeeze operation, which aggregates the feature map across the spatial dimensions to generate a channel

descriptor. The function of this descriptor is to embed the global distribution of channel feature responses, thereby enabling all layers of the network to utilize information from the global receptive field. After the aggregation, an excitation operation follows. This operation, in the form of a simple self-gating mechanism, takes the embedding as input and generates a set of modulation weights for each channel. These weights are applied to the feature map U to produce the output of the SE block, which can then be directly fed into subsequent layers of the network.

The loss function

In this study, each image in the dataset is associated with a corresponding binary mask. Skin lesion segmentation is treated as a pixel-level binary classification task, distinguishing the skin lesions from the background. The combination of Binary Cross-Entropy (BCE) loss and the Dice Similarity Coefficient (DSC) loss is used as the loss function to optimize the network parameters. This approach effectively addresses the challenge of skin lesion segmentation by balancing pixel accuracy and overlap between the predicted and ground truth masks.

The loss function, referred to as the BceDice loss, can be expressed as Equation 5:

$$\begin{aligned}L_{BCE} &= -\frac{1}{N} \sum_{i=1}^N [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)] \\ L_{Dice} &= 1 - \frac{2|X \cap Y|}{|X| + |Y|} \\ L_{BCEDice} &= \alpha_1 L_{BCE} + \alpha_2 L_{Dice}\end{aligned}\quad (5)$$

TABLE 1 Experimental comparison of MCANet with other models on the ISIC2017 dataset.

Model	Params	GFLOPs	mIoU (%)	DSC (%)
UNet (2015) [36]	7.77	13.76	76.98	86.99
TransFuse (2021) [37]	26.16	11.5	79.21	88.4
FAT-Net (2022) [38]	30	23	76.53	85
MALUNet (2022) [39]	0.175	0.083	78.78	88.13
QGD-Net (2023) [40]	0.777	—	72.58	84.1
LCAUNet (2023) [41]	13.38	18.91	76.1	86.6
SCSONet (2024) [42]	0.149	0.056	80.14	88.97
PL-Net (2024) [43]	15.03	—	77.9	85.9
UCM-Net (2024) [44]	0.499	0.047	80.71	87.66
CSAP-UNet-S (2024) [45]	27.5	8.918	81.5	88.8
ELANet (2024) [46]	0.459	8.43	82.87	90.6
MCANet (ours)	0.128	0.022	83.25	90.86

where N is the total number of samples, Y represents the ground truth label, p_i represents the predicted values, y_i denotes the true label of sample i . $|X|$ and $|Y|$ denote the ground truth and the intersection of the predicted region, respectively. α_1 and α_2 represent the weights of the two loss functions. In this study, both weights are set to 1 by default.

Experiment

Datasets

The ISIC (International Skin Imaging Collaboration) datasets are benchmark datasets widely used in medical image analysis, particularly for dermoscopic image segmentation, classification, and automated skin cancer detection. These datasets feature high-resolution dermoscopic images with comprehensive annotations, including lesion boundaries, diagnostic labels, and metadata. Covering a diverse range of skin conditions, they are designed to support tasks such as lesion segmentation, feature extraction, and disease classification. Notably, the ISIC2017 and ISIC2018 datasets have been instrumental in advancing research on melanoma detection and other skin diseases through the annual ISIC

Challenges. Our research is specifically conducted on the ISIC2017 and ISIC2018 datasets. Figure 3 are some sample images from the ISIC2017 and ISIC2018 datasets.

Experiment details

All experiments were implemented using the PyTorch framework and performed on a laptop equipped with an NVIDIA GeForce RTX 3080 Ti GPU with 8 GB of memory. Based on established practices, all images were normalized and resized to 256×256 pixels. Data augmentation techniques, including vertical flipping, horizontal flipping, and random rotations, were applied. The loss function used was the BCE-Dice loss, as defined in Equation 6.

$$L_{BCE-Dice} = \alpha \cdot \left(-\frac{1}{N} \sum_{i=1}^N (y_i \cdot \log(\hat{y}_i) + (1 - y_i) \cdot \log(1 - \hat{y}_i)) \right) + \beta \cdot \left(1 - \frac{2 \cdot \sum_{i=1}^N y_i \cdot \hat{y}_i + \epsilon}{\sum_{i=1}^N y_i + \sum_{i=1}^N \hat{y}_i + \epsilon} \right) \quad (6)$$

where y_i represents the ground truth label, \hat{y}_i denotes the predicted value, N is the total number of pixels, ϵ is a small constant which is set to 10 in this work, α and β are the weights for the BCE and Dice components. AdamW was utilized as the optimizer with an initial learning rate of 0.001, dynamically adjusted using a cosine annealing scheduler. The maximum number of iterations was set to 50, with a minimum learning rate of 0.0001. The training process was conducted over 300 epochs with a batch size of 8.

Evaluation metrics

In this study, segmentation performance is assessed using the mean Intersection over Union (mIoU), Dice Similarity Coefficient (DSC), and Accuracy (Acc), as defined in Equation 7. Additionally, the number of parameters is represented by Params, measured in millions (M), and computational complexity is quantified in GFLOPs. It is important to note that both Params and GFLOPs are calculated based on an input size of 256×256 .

$$\begin{cases} mIoU = \frac{TP}{TP + FP + FN} \\ DSC = \frac{2TP}{2TP + FP + FN} \end{cases} \quad (7)$$

Where, TP, FP, FN, and TN represent True Positives, False Positives, False Negatives, and True Negatives, respectively.

Segmentation result analysis

In this section, we conducted comparative experiments on melanoma segmentation using the ISIC2017 and ISIC2018 skin lesion segmentation datasets and evaluated the test results. The evaluation metrics include DSC, mIoU, params, and GFLOPs. The results are presented in Tables 1, 2, where we perform

TABLE 2 Experimental comparison of MCANet with other models on the ISIC2018 dataset.

Model	Params	GFLOPs	mIoU (%)	DSC (%)
UNet (2015) [36]	7.77	13.76	78.13	86.99
Unet ++ (2018) [47]	9.16	34.86	78.92	87.83
TransFuse (2021) [37]	26.16	11.5	80.63	89.27
MALUNet (2022) [39]	0.175	0.083	80.25	89.04
AMCC-Net (2023) [48]	0.845	—	80.18	89
SCSONet (2024) [42]	0.149	0.056	80.99	89.5
MCNMF-Unet (2024) [49]	0.332	0.0538	81.99	89.96
GIVTED-Net (2024) [50]	0.19	0.56	79.79	87.61
UCM-Net (2024) [44]	0.499	0.047	81.26	88.48
ELANet (2024) [46]	0.459	8.43	81.85	90.1
MCANet (ours)	0.128	0.024	83.68	91.12

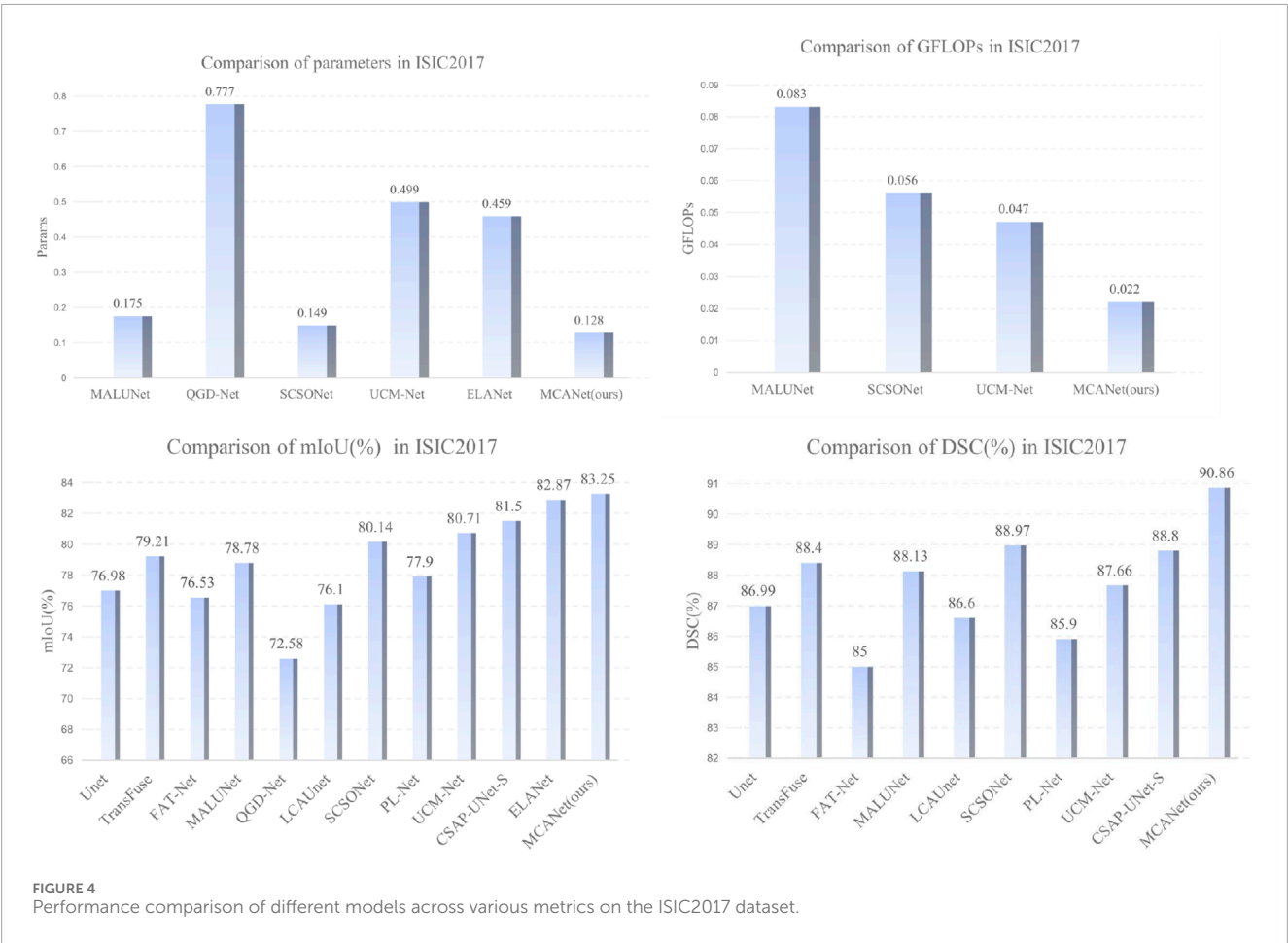
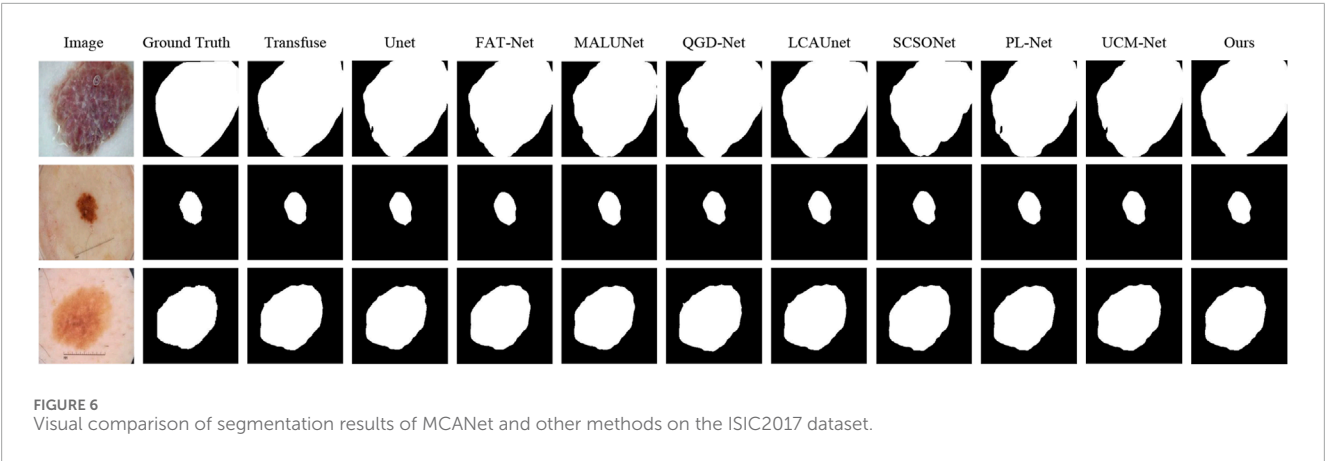
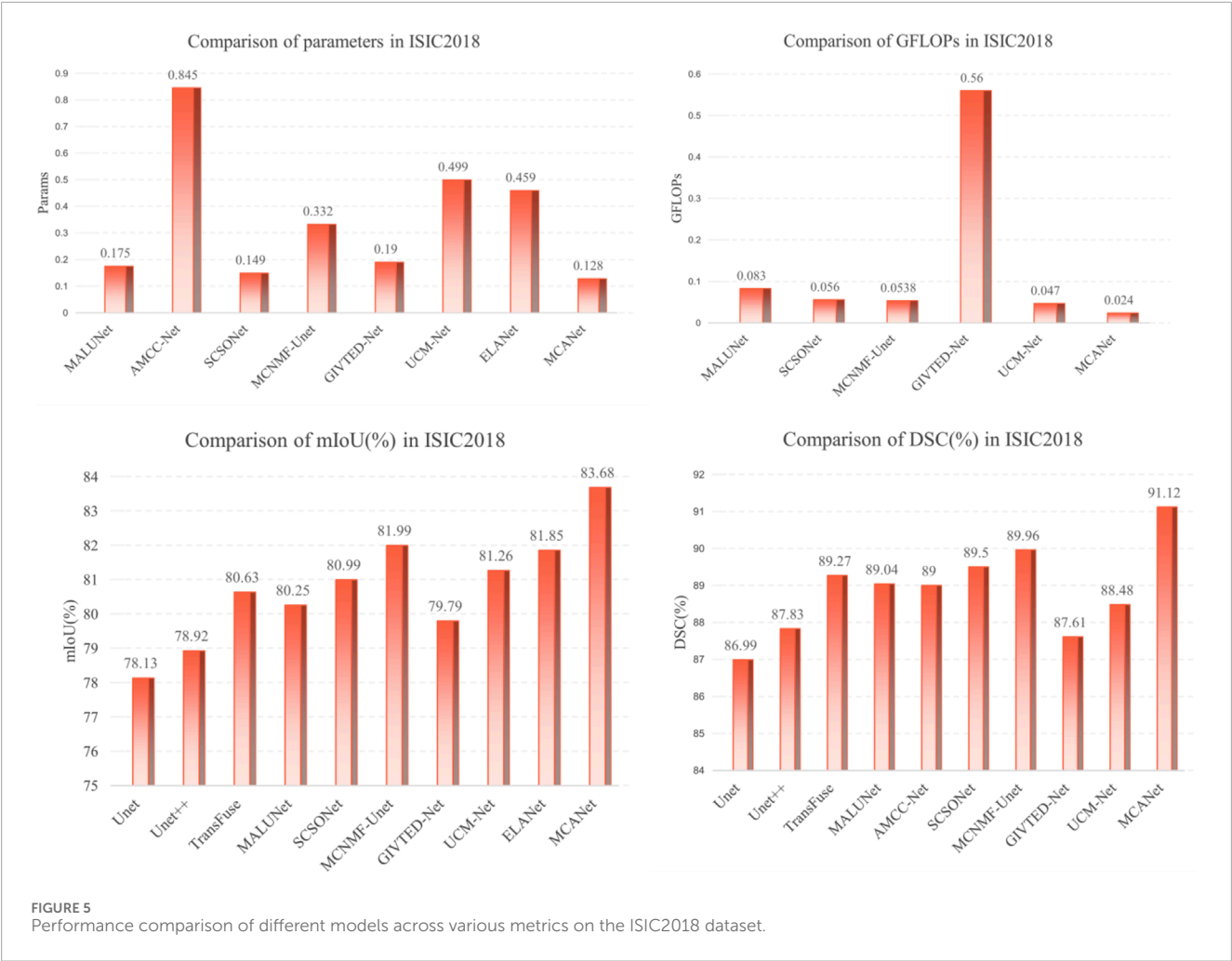


FIGURE 4 Performance comparison of different models across various metrics on the ISIC2017 dataset.



a comprehensive comparison of the proposed model with the following methods: UNet [36], Transfuse [37], FATNet [38], MALUNet [39], QGD-Net [40], LCA-UNet [41], SCSONet [42], PL-Net [43], UCM-Net [44], CSAP-UNet-S [45], and ELA-Net [46].

In addition, bar charts are utilized in this study to visually illustrate the performance of different models on various metrics, providing a clearer comparison between our method and others. Specifically, for the comparison of lightweight metrics, only models designed with lightweight objectives were selected, with the results presented in Figures 4, 5. The experimental results indicate that MCANet outperforms all other methods in both data sets in terms of DSC and mIoU metrics. Notably, MCANet achieves Dice scores

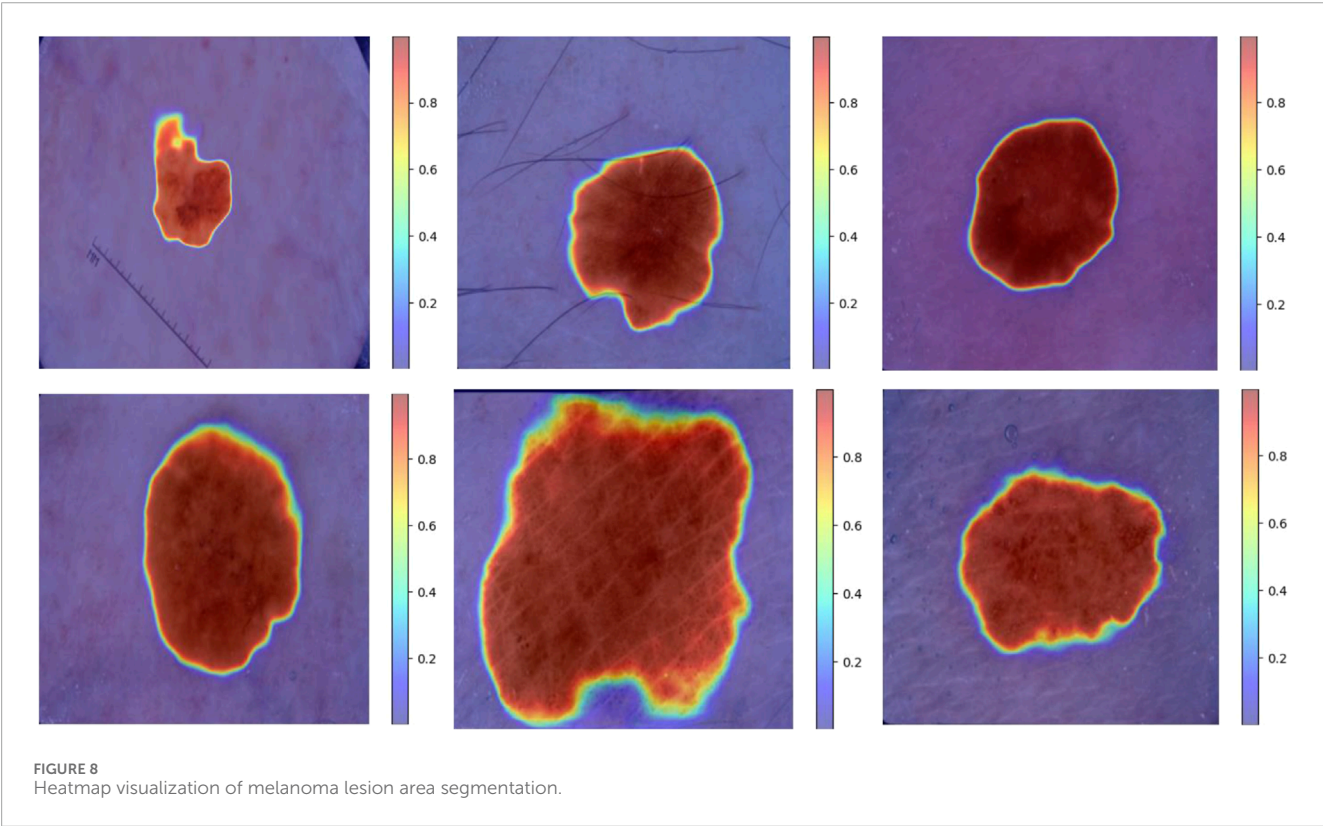
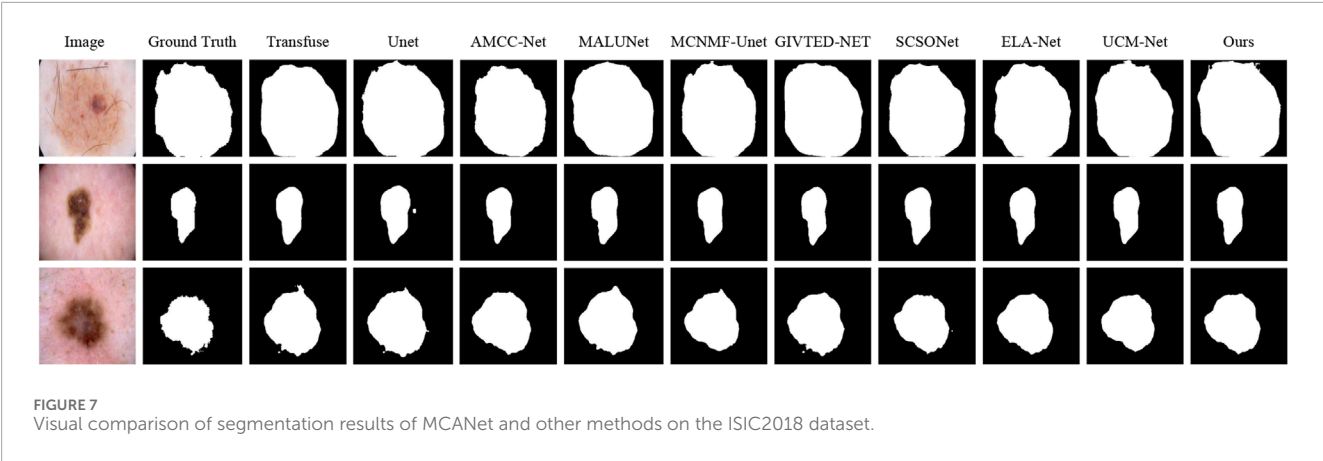


TABLE 3 Ablation experiments with different module combinations.

Model	Params	GFLOPs	mIoU (%)	DSC (%)
Base	0.112	0.021	79.01	88.27
Base + AEAM	0.127	0.022	81.39	90.34
BASE + ISDCConv	0.114	0.022	82.32	90.84
MCANet	0.128	0.024	83.25	90.86

exceeding 0.9 on the ISIC datasets, significantly outperforming all comparison models and demonstrating its superior segmentation performance.

Furthermore, to further validate the segmentation performance of the model, we present the visual segmentation results on the ISIC dataset, as shown in Figures 6, 7. Although there are some differences between the MCANet segmentation results and mask images, MCANet outperforms other models in capturing detailed information from medical images, giving it a significant advantage in accurately segmenting the areas of the injury. Specifically, Figure 8 shows that MCANet can more accurately capture the target location in segmentation tasks involving smaller lesions, with finer and more precise segmentation of the lesion boundaries. However, our study also has some limitations. First, although MCANet demonstrates impressive performance on the ISIC datasets, its generalizability to other medical imaging datasets remains to be fully explored. In addition, while the model is lightweight in design, further optimization is required to meet the strict deployment constraints

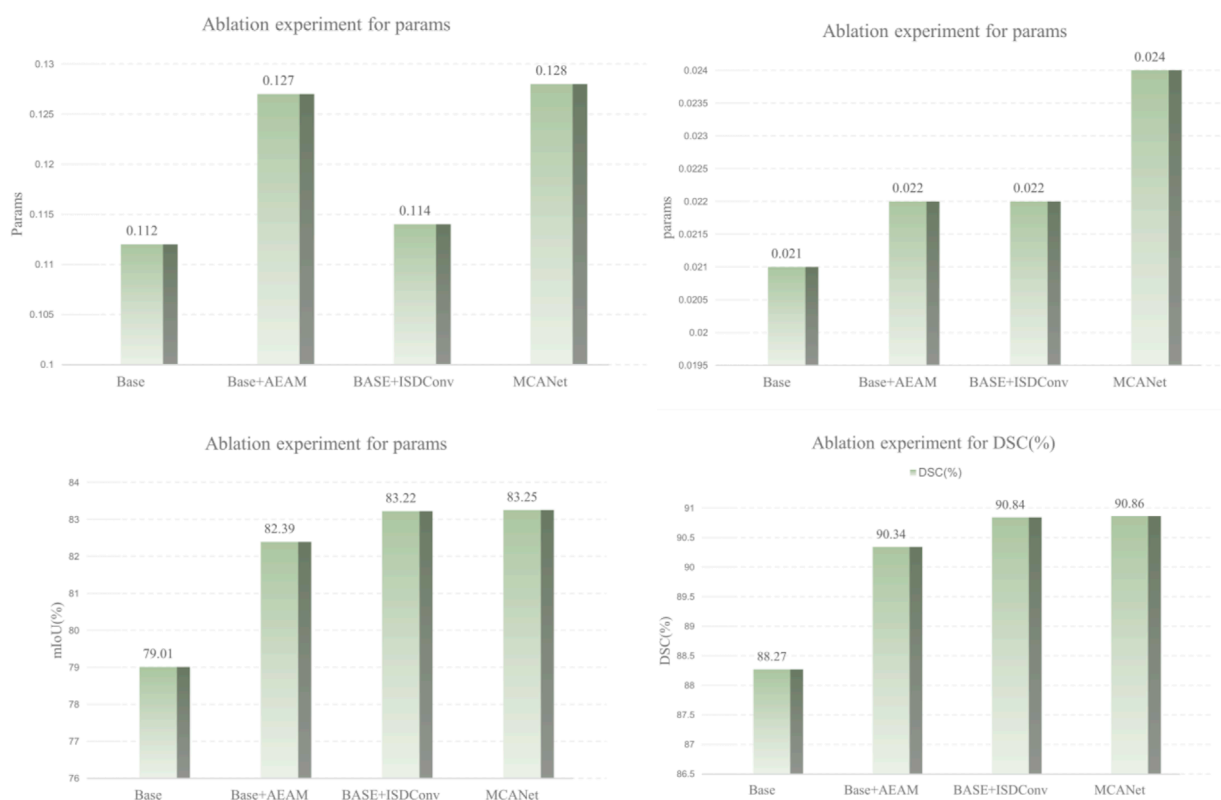


FIGURE 9
Visual representation of the impact of different modules on model performance.

of resource-constrained devices, such as smartphones or embedded systems. Another limitation lies in the annotation quality of the datasets used, as potential noise in the segmentation masks may influence the model's learning process. Finally, despite MCANet's ability to capture detailed features, there are still some challenges in handling highly irregular or extremely small lesions, which may require more advanced attention mechanisms. To address these issues, future research will focus on several directions. First, extending the evaluation to additional datasets with diverse imaging modalities can help assess the robustness and versatility of MCANet. Second, incorporating techniques such as knowledge distillation or pruning could further improve the model's efficiency for deployment in real-time scenarios. Third, exploring semi-supervised or unsupervised learning methods may reduce dependency on high-quality annotations, enabling better performance even with noisy labels. Finally, integrating advanced multi-scale feature extraction modules could enhance the model's ability to handle challenging segmentation tasks involving complex lesion patterns.

Ablation study on module effectiveness

To evaluate the contribution of each module in MCANet, we designed and conducted a series of ablation studies, with the results summarized in Table 3. Using the SCSONet baseline model as a reference, we performed comparative experiments with different combinations of the proposed modules on the ISIC dataset. Furthermore, to provide a clearer visualization of the impact of

each module on segmentation performance, we used bar charts to illustrate variations in key metrics, such as DSC and mIoU, as shown in Figure 9. In the ablation study, "Base + AEAM" represents the integration of the proposed AEAM module into the baseline model, "Base + ISDConv" denotes the addition of the ISDConv module to the baseline, and "MCANet" refers to the complete network architecture proposed in this study. From Table 3 and the bar chart, it can be observed that integrating the proposed modules into the baseline model not only results in negligible increases in parameter count and computational complexity but also leads to significant improvements in segmentation performance. Specifically, as the modules are progressively added, the segmentation performance steadily improves, with the key metrics DSC and mIoU ultimately reaching 0.9086 and 0.8325, representing increases of 2.93% and 5.37%, respectively, compared to the baseline. The bar chart further illustrates this performance improvement trend, visually highlighting the contribution of each module.

Moreover, the experimental results demonstrate that the proposed modules collaborate effectively, with the addition of individual modules not causing any degradation in overall performance but instead continuously improving segmentation accuracy. Additionally, our module design is highly adaptable, allowing for seamless integration into other network architectures without requiring significant modifications to the original structure. For instance, incorporating the AEAM or ISDConv modules into other networks results in varying degrees of performance improvement, validating the generalizability and practicality of the proposed modules.

In summary, the results of the ablation studies and their visual analysis demonstrate the significant contributions of the proposed modules to the model's performance. These improvements not only enhance the segmentation capability of MCANet but also highlight the academic significance and practical applicability of our work in the field of medical image segmentation.

Conclusion

Medical image analysis typically requires significant computational resources, which directly impact diagnostic speed and accuracy. Advanced methods like deep learning are resource-intensive, making them difficult to implement in resource-constrained environments. To address this, we propose MCAN, a novel lightweight network architecture featuring ISDConv, AEAM, and dynamic convolution. Our model reduces computational costs while maintaining performance, achieving competitive segmentation with 0.128M parameters and 0.022 GFLOPs. However, due to the limited dataset, the model's generalization ability requires further investigation.

Future research can focus on several key areas. Firstly, further optimization of lightweight techniques and attention mechanisms is needed, especially for specific types of medical images. For example, improving the prediction accuracy and robustness of melanoma images across different skin types is an important direction. Additionally, due to the limited dataset size in this study, further validation of the model's generalization ability is required. Future work should aim to expand the dataset with more representative clinical data to assess the model's performance in real-world clinical environments, particularly in resource-constrained settings such as mobile medical devices or low-resource hospitals. Finally, our method could be extended to multi-modal tasks, such as integrated diagnosis using CT and MRI, with a focus on improving the model's fusion capability while maintaining computational efficiency.

Data availability statement

The original contributions presented in the study are included in the article/supplementary material, further inquiries can be directed to the corresponding author.

Author contributions

ZL: Conceptualization, Data curation, Investigation, Methodology, Project administration, Software, Supervision,

Validation, Visualization, Writing—original draft, Writing—review and editing. HW: Conceptualization, Formal Analysis, Resources, Software, Supervision, Validation, Visualization, Writing—original draft, Writing—review and editing. HC: Data curation, Investigation, Methodology, Project administration, Resources, Supervision, Writing—review and editing. CL: Conceptualization, Formal Analysis, Resources, Software, Writing—review and editing. AY: Data curation, Formal Analysis, Funding acquisition, Supervision, Writing—review and editing.

Funding

The author(s) declare that financial support was received for the research, authorship, and/or publication of this article. This research was sponsored by the Special key project of Chongqing technology innovation and application development (CSTB2024TIAD-STX0023, CSTB2024TIAD-STX0030, CSTB2024TIAD-STX0037), Science and Technology Research Program of Chongqing Municipal Education (KJQN202400618) and “Unveiling and Leading” Project by the Chongqing Municipal Bureau of Industry and Information Technology (2022-37).

Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Generative AI statement

The author(s) declare that no Generative AI was used in the creation of this manuscript.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

References

- Schadendorf D, Van Akkooi AC, Berking C, Griewank KG, Gutzmer R, Hauschild A, et al. Melanoma. *The Lancet* (2018) 392:971–84. doi:10.1016/s0140-6736(18)31559-9
- Zhang Y, Li Z, Li H, Tao D. Prototype-driven and multi-expert integrated multi-modal mr brain tumor image segmentation. *IEEE Trans Instrumentation Meas* (2024) 74:1–14. doi:10.1109/tim.2024.3500067
- Li Z, Zhang Y, Li H, Chai Y, Yang Y. Deformation-aware and reconstruction-driven multimodal representation learning for brain tumor segmentation with missing modalities. *Biomed Signal Process Control* (2024) 91:106012. doi:10.1016/j.bspc.2024.106012
- Dong Z, Li J, Hua Z. Transformer-based multi-attention hybrid networks for skin lesion segmentation. *Expert Syst Appl* (2024) 244:123016. doi:10.1016/j.eswa.2023.123016
- Sun Y, Dai D, Zhang Q, Wang Y, Xu S, Lian C. Msca-net: multi-scale contextual attention network for skin lesion segmentation. *Pattern Recognition* (2023) 139:109524. doi:10.1016/j.patcog.2023.109524

6. Qiu S, Li C, Feng Y, Zuo S, Liang H, Xu A. Gfanet: gated fusion attention network for skin lesion segmentation. *Comput Biol Med* (2023) 155:106462. doi:10.1016/j.combiomed.2022.106462
7. Qi K, Yang H, Li C, Liu Z, Wang M, Liu Q, et al. X-net: brain stroke lesion segmentation based on depthwise separable convolution and long-range dependencies. In: Proceedings, Part III 22nd International Conference Medical Image Computing and Computer Assisted Intervention–MICCAI 2019; October 13–17, 2019; Shenzhen, China. Springer (2019) p. 247–55.
8. Liu X, Yang H, Qi K, Dong P, Liu Q, Liu X, et al. Msdf-net: multi-scale deep fusion network for stroke lesion segmentation. *IEEE Access* (2019) 7:178486–95. doi:10.1109/access.2019.2958384
9. Yang L, Fan C, Lin H, Qiu Y. Rema-net: an efficient multi-attention convolutional neural network for rapid skin lesion segmentation. *Comput Biol Med* (2023) 159:106952. doi:10.1016/j.combiomed.2023.106952
10. Liu Y, Shi Y, Mu F, Cheng J, Chen X. Glioma segmentation-oriented multi-modal mr image fusion with adversarial learning. *IEEE/CAA J Automatica Sinica* (2022) 9:1528–31. doi:10.1109/jas.2022.105770
11. Zhu Z, Wang Z, Qi G, Mazur N, Yang P, Liu Y. Brain tumor segmentation in mri with multi-modality spatial information enhancement and boundary shape correction. *Pattern Recognition* (2024) 153:110553. doi:10.1016/j.patcog.2024.110553
12. Zhu Z, He X, Qi G, Li Y, Cong B, Liu Y. Brain tumor segmentation based on the fusion of deep semantics and edge information in multimodal mri. *Inf Fusion* (2023) 91:376–87. doi:10.1016/j.inffus.2022.10.022
13. Liu Y, Qi Z, Cheng J, Chen X. Rethinking the effectiveness of objective evaluation metrics in multi-focus image fusion: a statistic-based approach. *IEEE Trans Pattern Anal Machine Intelligence* (2024) 46:5806–19. doi:10.1109/tpami.2024.3367905
14. Khan TM, Naqvi SS, Meijering E. Esdmr-net: a lightweight network with expand-squeeze and dual multiscale residual connections for medical image segmentation. *Eng Appl Artif Intelligence* (2024) 133:107995. doi:10.1016/j.engappai.2024.107995
15. Zhou Y, Kang X, Ren F, Lu H, Nakagawa S, Shan X. A multi-attention and depthwise separable convolution network for medical image segmentation. *Neurocomputing* (2024) 564:126970. doi:10.1016/j.neucom.2023.126970
16. Liu T, Liu H, Yang B, Zhang Z. LDCNet: limb direction cues-aware network for flexible HPE in industrial behavioral biometrics systems. *IEEE Trans Ind Inform* (2023) 20:8068–78. doi:10.1109/tii.2023.3266366
17. Ma T, Wang K, Hu F. Lmu-net: lightweight u-shaped network for medical image segmentation. *Med and Biol Eng and Comput* (2024) 62:61–70. doi:10.1007/s11517-023-02908-w
18. Feng L, Wu K, Pei Z, Weng T, Han Q, Meng L, et al. Mlu-net: a multi-level lightweight u-net for medical image segmentation integrating frequency representation and mlp-based methods. *IEEE Access* (2024) 12:20734–51. doi:10.1109/access.2024.3360889
19. Ruan J, Xie M, Gao J, Liu T, Fu Y. Ege-unet: an efficient group enhanced unet for skin lesion segmentation. In: International conference on medical image computing and computer-assisted intervention. Springer (2023) p. 481–90.
20. Lei T, Sun R, Du X, Fu H, Zhang C, Nandi AK. Sgu-net: shape-guided ultralight network for abdominal image segmentation. *IEEE J Biomed Health Inform* (2023) 27:1431–42. doi:10.1109/jbhi.2023.3238183
21. Chen Y, Dai X, Liu M, Chen D, Yuan L, Liu Z. Dynamic convolution: attention over convolution kernels. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (2020) p. 11030–9.
22. Li H, Cen Y, Liu Y, Chen X, Yu Z. Different input resolutions and arbitrary output resolution: a meta learning-based deep framework for infrared and visible image fusion. *IEEE Trans Image Process* (2021) 30:4070–83. doi:10.1109/tip.2021.3069339
23. Li H, Liu J, Zhang Y, Liu Y. A deep learning framework for infrared and visible image fusion without strict registration. *Int J Comp Vis* (2024) 132:1625–44. doi:10.1007/s11263-023-01948-x
24. Huang H, Chen Z, Zou Y, Lu M, Chen C, Song Y, et al. Channel prior convolutional attention for medical image segmentation. *Comput Biol Med* (2024) 178:108784. doi:10.1016/j.combiomed.2024.108784
25. Shaker AM, Maaz M, Rasheed H, Khan S, Yang MH, Khan FS. Unetr++: delving into efficient and accurate 3d medical image segmentation. *IEEE Trans Med Imaging* (2024). doi:10.1109/TMI.2024.3398728
26. Fu Y, Liu J, Shi J. Tsca-net: Transformer based spatial-channel attention segmentation network for medical images. *Comput Biol Med* (2024) 170:107938. doi:10.1016/j.combiomed.2024.107938
27. Xiong J, Tang M, Zong L, Li L, Hu J, Bian D, et al. Ina-net: an integrated noise-adaptive attention neural network for enhanced medical image segmentation. *Expert Syst Appl* (2024) 258:125078. doi:10.1016/j.eswa.2024.125078
28. Song E, Zhan B, Liu H. Combining external-latent attention for medical image segmentation. *Neural Networks* (2024) 170:468–77. doi:10.1016/j.neunet.2023.10.046
29. Huang Z, Cheng S, Wang L. Medical image segmentation based on dynamic positioning and region-aware attention. *Pattern Recognition* (2024) 151:110375. doi:10.1016/j.patcog.2024.110375
30. Yang S, Zhang X, Chen Y, Jiang Y, Feng Q, Pu L, et al. Ucnnet: a lightweight and precise medical image segmentation network based on efficient large kernel u-shaped convolutional module design. *Knowledge-Based Syst* (2023) 278:110868. doi:10.1016/j.knsys.2023.110868
31. Sun Q, Dai M, Lan Z, Cai F, Wei L, Yang C, et al. Ucr-net: U-shaped context residual network for medical image segmentation. *Comput Biol Med* (2022) 151:106203. doi:10.1016/j.combiomed.2022.106203
32. Nisa SQ, Ismail AR. Dual u-net with resnet encoder for segmentation of medical images. *Int J Adv Comp Sci Appl* (2022) 13. doi:10.14569/ijacsa.2022.0131265
33. Zhao Q, Zhong L, Xiao J, Zhang J, Chen Y, Liao W, et al. Efficient multi-organ segmentation from 3d abdominal ct images with lightweight network and knowledge distillation. *IEEE Trans Med Imaging* (2023) 42:2513–23. doi:10.1109/tmi.2023.3262680
34. Wang H, Zhang D, Song Y, Liu S, Wang Y, Feng D, et al. Segmenting neuronal structure in 3d optical microscope images via knowledge distillation with teacher-student network. In: 2019 IEEE 16th international symposium on biomedical imaging (ISBI 2019). IEEE (2019) p. 228–31.
35. Hajabdollahi M, Esfandiarpour R, Khadivi P, Soroushmehr SMR, Karimi N, Samavi S. Simplification of neural networks for skin lesion image segmentation using color channel pruning. *Comput Med Imaging Graphics* (2020) 82:101729. doi:10.1016/j.compmimag.2020.101729
36. Ronneberger O, Fischer P, Brox T. U-net: convolutional networks for biomedical image segmentation. In: proceedings, part III 18th international conference Medical image computing and computer-assisted intervention–MICCAI 2015; October 5–9, 2015; Munich, Germany. Springer (2015) p. 234–41.
37. Zhang Y, Liu H, Hu Q. Transfuse: fusing transformers and cnns for medical image segmentation. In: proceedings, Part I 24th international conference Medical image computing and computer assisted intervention–MICCAI 2021; September 27–October 1, 2021; Strasbourg, France. Springer (2021) p. 14–24.
38. Wu H, Chen S, Chen G, Wang W, Lei B, Wen Z. Fat-net: feature adaptive transformers for automated skin lesion segmentation. *Med image Anal* (2022) 76:102327. doi:10.1016/j.media.2021.102327
39. Ruan J, Xiang S, Xie M, Liu T, Fu Y. Malunet: a multi-attention and lightweight unet for skin lesion segmentation. In: 2022 IEEE International Conference on Bioinformatics and Biomedicine (BIBM). IEEE (2022) p. 1150–6.
40. Wang J, Huang G, Zhong G, Yuan X, Pun CM, Deng J. Qgd-net: a lightweight model utilizing pixels of affinity in feature layer for dermoscopic lesion segmentation. *IEEE J Biomed Health Inform* (2023) 27:5982–93. doi:10.1109/jbhi.2023.3320953
41. Zhang Q, Bai R, Peng B, Wang Z, Liu Y. Fft pattern recognition of crystal hrtem image with deep learning. *Micron* (2023) 166:103402. doi:10.1016/j.micron.2022.103402
42. Chen H, Li Z, Huang X, Peng Z, Deng Y, Tang L, et al. Scsonet: spatial-channel synergistic optimization net for skin lesion segmentation. *Front Phys* (2024) 12:1388364. doi:10.3389/fphy.2024.1388364
43. Cheng J, Gao C, Lu H, Ming Z, Yang Y, Zhu M. Pl-net: progressive learning network for medical image segmentation (2021) arXiv preprint arXiv:2110.14484.
44. Weng S, Zhu T, Zhang T, Zhang C. Ucm-net: a u-net-like tampered-region-related framework for copy-move forgery detection. *IEEE Trans Multimedia* (2023) 26:750–63. doi:10.1109/tmm.2023.3270629
45. Fan X, Zhou J, Jiang X, Xin M, Hou L. Csap-unet: convolution and self-attention paralleling network for medical image segmentation with edge enhancement. *Comput Biol Med* (2024) 172:108265. doi:10.1016/j.combiomed.2024.108265
46. Nie T, Zhao Y, Yao S. Ela-net: an efficient lightweight attention network for skin lesion segmentation. *Sensors* (2024) 24:4302. doi:10.3390/s24134302
47. Zhou Z, Rahman Siddiquee MM, Tajbakhsh N, Liang J. Unet++: a nested u-net architecture for medical image segmentation. In: Proceedings 4 Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: 4th International Workshop, DLMIA 2018, and 8th International Workshop, ML-CDS 2018, Held in Conjunction with MICCAI 2018; September 20, 2018; Granada, Spain. Springer (2018) p. 3–11.
48. Dayananda C, Yamanakkanavar N, Nguyen T, Lee B. Amcc-net: an asymmetric multi-cross convolution for skin lesion segmentation on dermoscopic images. *Eng Appl Artif Intelligence* (2023) 122:106154. doi:10.1016/j.engappai.2023.106154
49. Yuan L, Song J, Fan Y. Mcnmf-unet: a mixture conv-mlp network with multi-scale features fusion unet for medical image segmentation. *PeerJ Comp Sci* (2024) 10:e1798. doi:10.7717/peerj-cs.1798
50. Al-Fahsi RDH, Prawirosoenoto ANF, Nugroho HA, Ardiyanto I. Givted-net: ghostnet-mobile inversion vit encoder-decoder network for lightweight medical image segmentation. *IEEE Access* (2024) 12:81281–92. doi:10.1109/ACCESS.2024.3411870



OPEN ACCESS

EDITED BY

Zhiqin Zhu,
Chongqing University of Posts and
Telecommunications, China

REVIEWED BY

Guanqiu Qi,
Buffalo State College, United States
Yimin Chen,
University of Massachusetts Lowell,
United States

*CORRESPONDENCE

Bojian Chen,
✉ cbj.android@gmail.com

RECEIVED 25 October 2024

ACCEPTED 21 November 2024

PUBLISHED 24 December 2024

CITATION

Yang M, Chen B, Lin C, Yao W and Li Y (2024)
SGI-YOLOv9: an effective method for crucial
components detection in the power
distribution network.
Front. Phys. 12:1517177.
doi: 10.3389/fphy.2024.1517177

COPYRIGHT

© 2024 Yang, Chen, Lin, Yao and Li. This is an
open-access article distributed under the
terms of the [Creative Commons Attribution
License \(CC BY\)](#). The use, distribution or
reproduction in other forums is permitted,
provided the original author(s) and the
copyright owner(s) are credited and that the
original publication in this journal is cited, in
accordance with accepted academic practice.
No use, distribution or reproduction is
permitted which does not comply with
these terms.

SGI-YOLOv9: an effective method for crucial components detection in the power distribution network

Mianfang Yang, Bojian Chen*, Chenxiang Lin, Wenxu Yao and
Yangdi Li

State Grid Fujian Electric Power Research Institute, FuZhou, China

The detection of crucial components in the power distribution network is of great significance for ensuring the safe operation of the power grid. However, the challenges posed by complex environmental backgrounds and the difficulty of detecting small objects remain key obstacles for current technologies. Therefore, this paper proposes a detection method for crucial components in the power distribution network based on an improved YOLOv9 model, referred to as SGI-YOLOv9. This method effectively reduces the loss of fine-grained features and improves the accuracy of small objects detection by introducing the SPDConv++ downsampling module. Additionally, a global context fusion module is designed to model global information using a self-attention mechanism in both spatial and channel dimensions, significantly enhancing the detection robustness in complex backgrounds. Furthermore, this paper proposes the Inner-PIoU loss function, which combines the advantages of Powerful-IoU and Inner-IoU to improve the convergence speed and regression accuracy of bounding boxes. To verify the effectiveness of SGI-YOLOv9, extensive experiments are conducted on the CPDN dataset and the PASCAL VOC 2007 dataset. The experimental results demonstrate that SGI-YOLOv9 achieves a significant improvement in accuracy for small object detection tasks, with an mAP@50 of 79.1% on the CPDN dataset, representing an increase of 3.9% compared to the original YOLOv9. Furthermore, it achieves an mAP@50 of 63.3% on the PASCAL VOC 2007 dataset, outperforming the original YOLOv9 by 1.6%.

KEYWORDS

crucial component, smart grid, attention mechanism, YOLOv9, deep learning

1 Introduction

With the continuous growth in electricity demand and the ongoing expansion of the power grid, the stability and reliability of the power distribution network, as a critical hub in the power system, have become increasingly important. The primary function of the power distribution network is to transmit electrical energy from high-voltage transmission networks to low-voltage consumer networks, and its reliability directly impacts the quality and safety of electricity supply to users. Crucial components of the power distribution network include insulators, arresters, transformers, and Cut-out Switches (COS), which must withstand harsh weather conditions, high mechanical stress, and extreme voltage, making them prone

to damage [1]. Therefore, the detection and monitoring of these crucial components have become a central focus in the maintenance and management of the power distribution network.

The power distribution network cover vast areas, with numerous and complexly distributed equipment, making traditional manual inspection methods insufficient to meet the operational and maintenance demands of modern power grids. Manual inspections are not only labor-intensive and inefficient, but they are also susceptible to geographical constraints, resulting in risks of omission and false detections. With the rapid advancement of computer vision technology, image detection has gradually replaced traditional manual inspections as a non-contact detection method [2]. This technology enables comprehensive, multi-angle, and high-precision inspection of crucial components in the power distribution network, significantly enhancing the intelligence and automation of component monitoring.

In the early stages of image detection, traditional methods primarily relied on handcraft feature extraction, including characteristics such as shape, color, and texture, combined with machine learning algorithms for recognition. Murthy V S et al. utilized a combination of Support Vector Machine (SVM) and Multiresolution Analysis (MRA) to detect defects in transmission line insulators, where MRA was used to capture insulator images, and SVM was applied to detect their condition. Hao J et al. applied Canny edge detection and directional angle selection to process insulator images, followed by the Hough transform to extract linear features of the damaged sections of the insulator. Zhang K et al. [3] proposed a method based on k-means clustering and morphological techniques to segment insulator images. Yu Y et al. [4] introduced a model that uses iterative curve evolution based on texture features and shape priors to detect insulators, though this method requires pre-acquisition of shape priors, limiting its applicability and resulting in slow detection speed. Zhao Z et al. [5] proposed a method that uses orientation angle detection and binary shape priors to locate insulators at different angles. However, traditional methods generally depend on feature extraction and shallow learning classification, and some even require the support of prior knowledge. These limitations make it difficult for such methods to cope with significantly varying complex scenes and render them vulnerable to noise and background interference, leading to weak generalization capabilities. As a result, traditional methods are often suitable only for images with simple backgrounds or large objects.

Deep learning-based object detection techniques, on the other hand, offer promising new possibilities for identifying key components. Architectures like Convolutional Neural Networks (CNNs) are capable of automatically extracting image features through multiple layers, which greatly enhances detection accuracy and efficiency [6–8]. By leveraging training on large-scale datasets, these models can perform consistently across a range of complex scenarios, minimizing the need for manual intervention and reducing the risk of misjudgment. This improvement bolsters the reliability and safety of power systems, providing robust technical support for the advancement of smart grid technologies.

Deep learning-based object detection research can be generally categorized into two main approaches. The first approach includes two-stage detection models like R-CNN [9], Faster R-CNN [10], and Mask R-CNN [11], which use a region proposal network (RPN)

to generate candidate object regions, followed by classification and regression to enhance detection accuracy. Such models are typically characterized by complex architectures and high detection accuracy but relatively slow processing speed. Zhao Z et al. [12] improved the anchor generation method of the Faster R-CNN model and optimized the non-maximum suppression (NMS) in the RPN, achieving improved insulator detection, particularly for insulators with varying aspect ratios, scales, and occlusions. However, the dataset utilized by this network contains almost no images of vertically oriented insulator strings. As a result, this method is incapable of detecting missing faults in images that include such types of insulator strings. Odo A et al. [13] utilized Mask R-CNN and RetinaNet to detect insulators and U-bolts on each tower. Dong C et al. [14] introduced an enhanced Cascade R-CNN that integrates Swin-v2 with a balanced feature pyramid to strengthen feature representation, while also incorporating side-aware boundary localization for greater precision in detecting small components in power transmission lines.

Another prominent category of algorithms comprises single-stage object detection models, such as the YOLO (You Only Look Once) series [15–21] and SSD [22]. These models bypass the need for region proposal networks, allowing them to directly execute classification and regression tasks following feature extraction by the backbone network [23]. This approach significantly reduces both training and inference time, enhancing efficiency. In practical engineering applications, due to the limitations of computational resources on devices, single-stage object detection algorithms are often preferred. Qi C et al. [24] enhanced the SSD model by using the lightweight SqueezeNet architecture and adding multiple convolutional layers and connection branches, thus improving feature extraction and enabling the detection of five types of electrical equipment in substations. Siddiqui et al. [25] developed an automated real-time system for detecting electrical equipment and analyzing faults, employing a CNN-based framework to identify insulators, arresters, and COS across different materials in complex settings. However, this method operates in a simplified environment with a single detection background and lacks interference from complex backgrounds. Liu Z et al. [26] created a large-scale dataset for transmission line component detection and optimized YOLOv4 by adding a prediction layer and refining the selection of positive and negative samples during training, thereby enhancing small object detection. Qiu Z et al. [27] preprocessed insulator images using the Laplacian sharpening method and improved the YOLOv4 model structure by incorporating the lightweight MobileNet convolutional neural network. However, its detection performance on blurry and small objects was suboptimal. Liu M et al. [28] improved YOLOv5 by incorporating diversified branch blocks (DBB), efficient channel attention (ECA), and an upgraded spatial pyramid pooling (SPP) module, with TensorRT utilized for accelerated edge detection of critical components. Liu C et al. [29] integrated a CBAM mixed attention module and Swin Transformer self-attention into YOLOv7, along with adding a dedicated small object detection layer to better identify small transmission line components. Chen B et al. [30] introduced innovative methods, including the Edge Detailed Shape Data Augmentation (EDSDA) and the Cross-Channel and Spatial Multi-Scale Attention (CCSMA) module, which enhanced the detection capability of insulator edge

shapes and defect features. Additionally, the design of the Re-BiC module and the MPDIoU localization loss function optimized feature fusion and computational efficiency, leading to significant improvements in detection accuracy and speed. He M et al. [31] introduced an improved YOLOv8 model for detecting insulators and fault areas, using GhostNet and an asymmetric convolution-based feature extraction module to enhance recognition in complex environments, while the ResPANet module fused high-resolution feature maps with residual skip connections to mitigate information loss in small feature layers. However, this method fails to effectively extract the features of subtle defects, resulting in poor detection performance for small target defects.

In practical applications, the small size of most key components in the power distribution network, along with the cluttered backgrounds, makes their detection particularly challenging. This poses significant difficulties for traditional detection models, driving researchers to focus on small object detection techniques to improve both accuracy and reliability. Developing more robust and effective methods for identifying these components in complex environments remains a critical research challenge in the field. Zhu Z et al. [32] proposed a small object detection network with a multi-level perception parallel structure. This network addressed the issues of lacking global representation information and the dense distribution of small objects through a global multi-level perception module and a dynamic region aggregation module, respectively. Qi G et al. [33] introduced an improved YOLOv5 algorithm, which utilized an Adaptive Spatial Parallel Convolution module (ASPCov) to extract multi-scale local context information of small objects. Additionally, to enhance the detection performance of small objects, it employed nearest-neighbor interpolation and sub-pixel convolution algorithms to construct high-resolution feature maps with rich semantic features. Li Y et al. [34] presented a feature fusion module (CGAL) based on both global and local attention mechanisms and designed a decoupled detection framework featuring a four-head structure, thereby enabling efficient detection of small objects. Zhang T et al. [35] optimized the backbone of YOLOv5 by incorporating a Convolutional Block Attention Module (CBAM) to focus on key information for insulator and defect detection while suppressing non-essential information. Additionally, small object detection anchors and layers were added to improve the detection of small defects.

Although the aforementioned studies have made significant progress in object detection, most of the research has primarily focused on detecting high-voltage transmission lines using UAV aerial images, where the targets are relatively large and the backgrounds are comparatively simple. However, compared to high-voltage transmission lines, the detection of key components in the power distribution network presents more complex challenges. Power distribution networks are typically deployed in areas with dense human activity and diverse geographical and environmental conditions, making them prone to obstructions from trees, buildings, and other structures. Moreover, the components within the power distribution network are generally smaller, more densely distributed, and often have similar appearances, further complicating the detection task. Existing algorithms still struggle with handling the complex backgrounds typical of distribution network scenarios, and they fail to effectively address the issue of information loss for small components during the process of deep

feature extraction, which significantly impairs detection accuracy. Therefore, there is an urgent need for more advanced methods that can overcome these challenges and improve detection performance in such complex environments.

To address the challenges of detecting crucial components in the power distribution network, we propose an innovative algorithm, SGI-YOLOv9. The main contributions of this paper are as follows.

- We propose the SPD++Conv downsampling module to replace the original downsampling module in the YOLOv9 backbone, effectively reducing the loss of fine-grained features. This allows the output feature maps of the backbone to retain more detailed information, significantly improving the detection accuracy of small objects.
- A Global Context Fusion module is proposed, leveraging the ability of the self-attention mechanism to capture global information. It models global context from both spatial and channel dimensions of the feature maps. This module effectively integrates global contextual features, enabling our method to perform more robustly in challenging scenarios such as complex backgrounds and occlusions.
- We propose the Inner-PIoU loss function, which combines the advantages of Powerful-IoU and Inner-IoU. By introducing scalable auxiliary bounding boxes, this method effectively addresses the slow convergence and limited generalization capabilities of traditional IoU loss function in small object detection.

2 Materials and methods

2.1 Dataset preparation and analysis

The dataset used in this study is provided by a private user on the Roboflow platform and has been named the Components of Power Distribution Network (CPDN) [36]. It contains 3,383 images and 25,185 instances, with each image having a resolution of 640×640 . The dataset includes common crucial components in the power distribution network, such as arresters, COS, insulators, and transformers, as shown in Figure 1. It can be observed that, except for transformers, the other components contain repeating circular structures called sheds, which vary in material, number, and size. The similarity in shed structures among these components increases the difficulty of classification.

Figure 2 presents image samples from the CPDN dataset in various environments, with each crucial component marked with different colored boxes, illustrating their distribution and position within the power distribution network. It is evident that the backgrounds in the power distribution network images are highly complex, covering diverse scenes such as urban streets, residential areas, and green spaces. Due to the influence of different angles in capturing images, components in these scenes are often obscured by various objects, and there is significant overlap of targets. Additionally, it is clear from the images that the components occupy a relatively small portion of the overall frame, with targets often blending into the background or multiple components being closely arranged. These factors pose considerable challenges for detection algorithms. The small visual differences between similar components further increase the risk of misclassification.

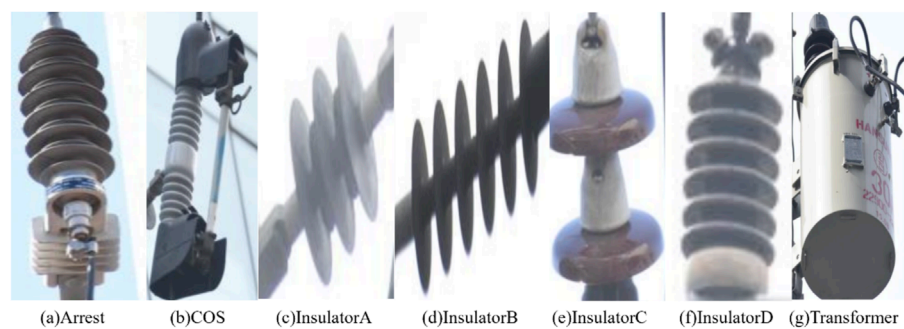


FIGURE 1 Illustration of crucial components in the power distribution network. (A) Arrest; (B) COS: Cut-out Switches; (C) Insulator (A) short polymer insulator; (D) Insulator (B) long polymer insulator; (E) Insulator (C) short porcelain insulator; (F) Insulator (D) long porcelain insulator; (G) Transformer.



FIGURE 2 Annotated examples of crucial components in the CPDN dataset.

In the CPDN dataset, the Insulator C and Transformer categories account for 7.7% and 7.2% of all instances, respectively, posing challenges related to class imbalance and small object detection. During training process, the model often assigns more weight to categories with a larger number of samples, which can lead to overfitting and reduce its ability to generalize to new datasets. To address this issue, we apply data augmentation techniques to mitigate the problem of class imbalance. Specifically, we use methods such as affine transformations, random noise, color jittering, and brightness adjustments [26] to generate diverse training samples. The augmented dataset is split into training, validation, and test sets with a 7:2:1 ratio.

2.2 Proposed method

In this study, we select YOLOv9 as the baseline model due to its various advantages. YOLOv9 introduces Programmable Gradient Information (PGI) in its architecture, which generates reliable gradient information through auxiliary reversible branches, solving the information bottleneck problem in deep network training and allowing the network to update weights more effectively. Meanwhile, a Generalized Efficient Layer Aggregation Network (GELAN) is proposed, which is based on gradient path planning and balances accuracy and inference speed.

However, in real-world transmission line applications, detecting crucial components presents multiple challenges. First, crucial components such as insulators are typically small objects, which places high demands on the model's ability to extract fine-grained features. Second, the background of transmission lines is highly complex, with many interfering factors, and the components are often occluded by other objects. As a result, YOLOv9 tends to have a higher rate of missed and false detections in these complex scenarios, particularly when detecting small objects and occluded objects. To address these issues, this study will improve YOLOv9 by enhancing feature extraction, contextual information utilization, and model training to improve the detection accuracy of small targets and enhance their robustness in occluded scenes, in order to achieve high-precision detection of crucial components.

2.2.1 Overview of SGI-YOLOv9 network

In this study, we propose an improved YOLOv9 method by optimizing two core modules in the original YOLOv9s model architecture. First, in the deep downsampling part of the backbone network, we design an SPDConv++ module to replace the original convolutional module. SPDConv++ spatially decomposes and reconstructs the input features, significantly reducing the loss of fine-grained feature information during downsampling and thereby improving the accuracy of small object detection. Second, in the neck part, we introduce a Global Context Fusion Module (GCFM), which combines spatial and channel self-attention mechanisms to model global contextual information. The GCF module effectively captures long-range contextual dependencies, enhancing robustness and detection accuracy in complex backgrounds and occluded scenarios. Additionally, during the training phase, we propose the Inner-PIoU loss function to improve convergence. The rest of the network structure and strategies remain consistent with the original YOLOv9s.

Based on the aforementioned improvements, we developed the final SGI-YOLOv9 algorithm, with the overall architecture shown in Figure 3. The following sections will provide a detailed explanation of the SGI-YOLOv9 method proposed in this paper.

2.2.2 SPDConv++ module

Small objects inherently possess limited feature information, making it essential to minimize information loss during feature extraction to maintain detection accuracy. In the original YOLOv9 architecture, a convolutional module with a stride of 2 is employed for downsampling, which inevitably results in the loss of fine-grained features, thereby impairing small object detection. To address this limitation and improve the model's small object detection capability, inspired by SPD-Conv (space-to-depth convolution) [37], we propose the SPD++ convolutional module, as illustrated in Figure 4. Specifically, for the input feature X with a size of $M \times M \times C$, we first sample and split it into four sub-features: X_1 , X_2 , X_3 and X_4 , defined as shown in Equations (1)–(4).

$$X_1 = X[0:M:2, 0:M:2] \quad (1)$$

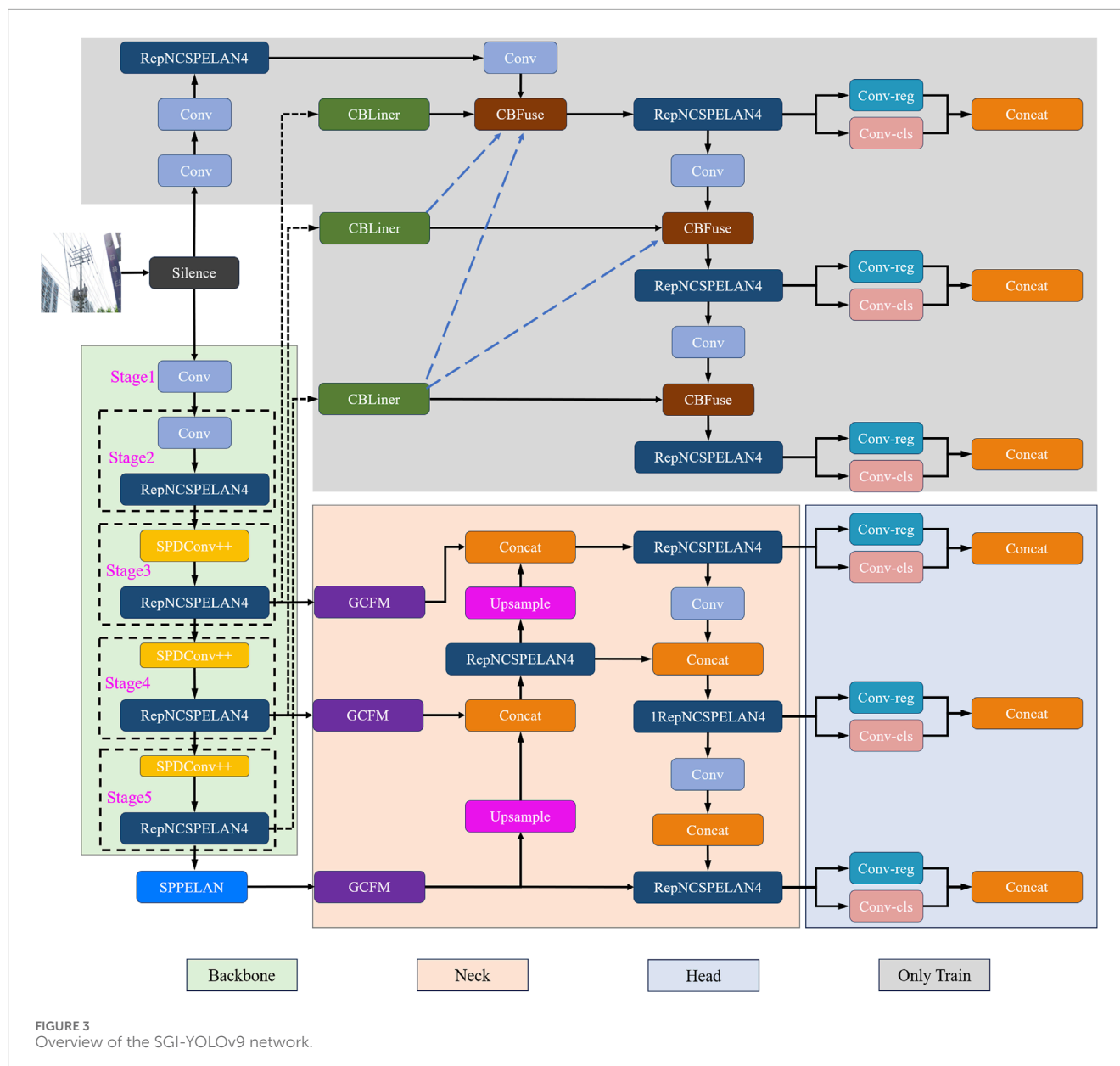
$$X_2 = X[1:M:2, 0:M:2] \quad (2)$$

$$X_3 = X[0:M:2, 1:M:2] \quad (3)$$

$$X_4 = X[1:M:2, 1:M:2] \quad (4)$$

The sub-features X_1 , X_2 , X_3 and X_4 are concatenated along the channel dimension to form X' , with dimensions $\frac{M}{2} \times \frac{M}{2} \times 4C$. At this stage, the spatial resolution of the features is half that of the input, and the number of channels is four times the input. As shown in Figure 4, the feature vectors sampled into the same sub-feature map in the input X are labeled with the same color to provide a more intuitive visualization. This demonstrates that the process of transforming X into four sub-features does not result in any feature loss, while the sub-features effectively preserve the spatial structural relationships of the original input, enabling the successful downsampling of input features without compromising information integrity. However, the concatenated sub-features have a channel count four times greater than that of the original input, inevitably introducing channel redundancy. The original SPD-Conv module employs a 1×1 convolutional layer to compress the channel dimensions to match the input, but directly applying 1×1 convolutions significantly impacts the output due to the presence of redundant information, resulting in feature loss. To address this limitation, we propose the SPD++ convolutional module, which incorporates a channel attention mechanism to emphasize important channels and suppress redundant ones. Following the channel attention module, a 1×1 convolution is applied to adjust the number of channels to match the input, effectively mitigating the adverse effects of channel redundancy.

The channel attention module begins by performing global max pooling and global average pooling on the input feature map. By using a three-layer fully connected feedforward network to interact with different channels, a set of attention weights can be learned that can suppress redundant channels and highlight important channels. After the three-layer network, the resulting pooled vectors are then



processed by a three-layer fully connected feedforward network. The outputs from both pooling operations are combined through element-wise addition, followed by the application of a sigmoid activation function to generate the channel attention weights. These weights are subsequently used to reweight the input features along the channel dimension, enhancing the model's ability to focus on the most informative channels.

In the entire SPDConv++ module, we do not use convolutions with a stride greater than 1, ensuring that downsampling is performed with minimal loss of fine-grained feature information. Compared to the original SPD convolution, we introduce a channel attention mechanism to the concatenated sub-features to highlight the more discriminative channels. Since the number of channels in the concatenated sub-features is significantly higher than that of the original input features, some redundant features are inevitable. Therefore, incorporating a channel attention mechanism

is essential to effectively reduce redundancy and enhance the model's discriminative capability.

2.2.3 Global context fusion

Contextual information is important for detecting small and occluded objects; however, traditional convolutional network architectures lack the ability to effectively integrate global contextual features. In recent years, self-attention mechanisms, due to their ability to establish long-range dependencies, have been widely used in visual tasks to fuse global contextual information [38, 39]. However, traditional visual self-attention mechanisms only perform computations in the spatial dimension, neglecting the modeling of information in the channel dimension. To more fully integrate global contextual information and further improve detection accuracy, this paper proposes a Global Context Fusion module. This module includes both spatial self-attention and channel self-attention

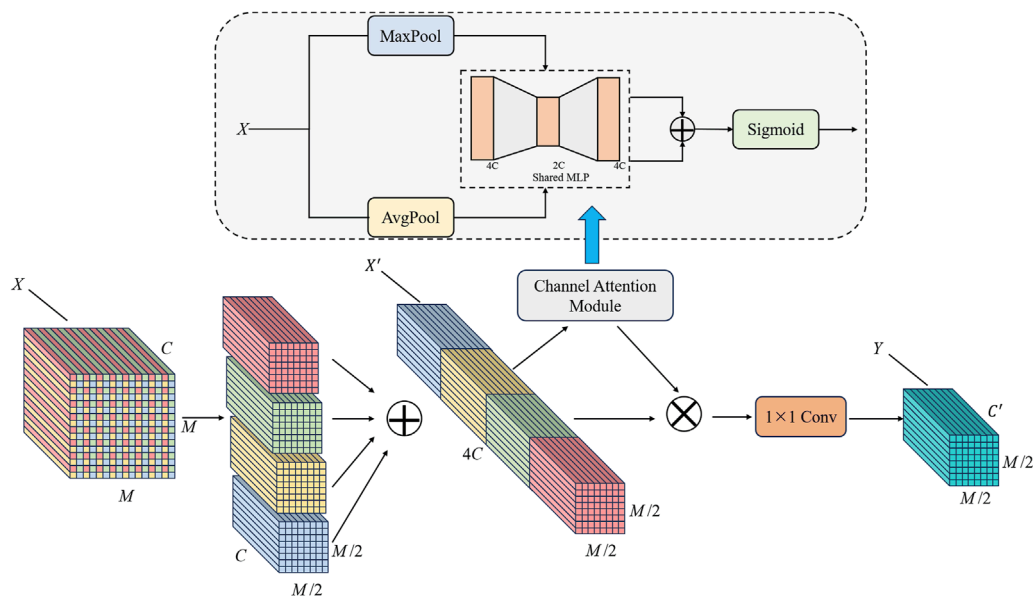


FIGURE 4
The SPDConv++ module.

mechanisms, which model global information from the spatial and channel dimensions, respectively, as shown in Figure 5A. The outputs of the spatial self-attention module and the channel self-attention module are concatenated along the channel dimension, and finally passed through a 1×1 convolutional layer to ensure that the number of output channels is the same as the input.

Specifically, as shown in Figure 5B, in the spatial self-attention module, for the input feature X , three parallel 1×1 convolutions are first applied to generate the query matrix Q , key matrix K , and value matrix V . Then, a Reshape operation is used to adjust their dimensions, such that $Q \in R^{HW \times C}$, $K \in R^{HW \times C}$ and $V \in R^{HW \times C}$. Subsequently, calculations are performed according to Equation (5), where d denotes the length of each feature vector in Q and K . Each feature vector in the Q , K , and V matrices corresponds to a patch of the input image. By computing the product of the Q matrix and K^T , the relationships between each patch and all other patches in the image are captured. These relationships are quantified as attention weights, ranging from 0 to 1, using the softmax function. The resulting weight matrix is then multiplied with the V matrix to generate a weighted output matrix. In this process, each feature vector in the output matrix is computed based on the connections of all patches in the image, thereby capturing global contextual information. After the computation, another Reshape operation is applied to adjust the result to match the dimensions of the input features. Finally, a 1×1 convolution is applied, and the result is element-wise added to the original input X .

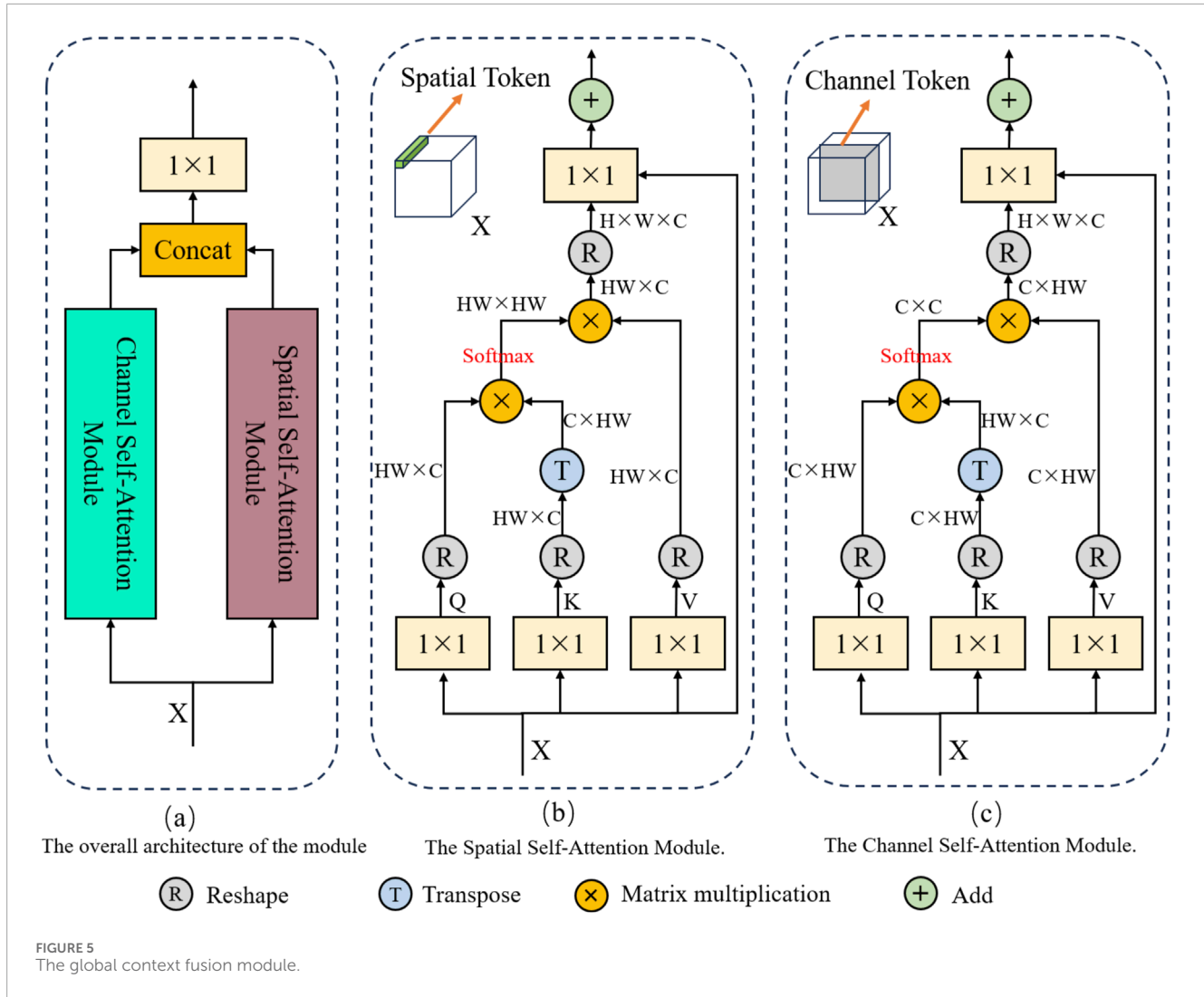
$$A(Q, K, V) = \text{SoftMax} \left[\frac{QK^T}{\sqrt{d}} \right] V \quad (5)$$

Additionally, The spatial self-attention module uses tokens corresponding to different spatial positions in the feature map as computing units to obtain contextual information from the spatial

dimension, but ignores information modeling from the channel dimension. In deep networks, the feature maps of different channels focus on expressing different feature information, so it is equally important to fuse global contextual information from the channel dimension. Therefore, in the channel self-attention module designed in this paper, we treat each channel as an independent token for self-attention mechanism calculation. As shown in Figure 5C, in the channel self-attention module, each channel of the input feature X is treated as an independent token. Therefore, in this module, after applying the Reshape operation to Q , K , and V , $Q \in R^{C \times HW}$, $K \in R^{C \times HW}$ and $V \in R^{C \times HW}$ are obtained. The subsequent computation process is the same as in the spatial self-attention module.

2.2.4 Inner-PloU

Intersection over Union (IoU) is a fundamental metric for assessing the performance of object detection systems. In these tasks, IoU quantifies the overlap between the predicted bounding box and the ground truth box, specifically calculating the ratio of the intersection area to the union area of these boxes. An effectively designed IoU-based loss function promotes better alignment of the predicted bounding box with the ground truth, thereby enhancing model convergence speed. In YOLOv9, the Complete Intersection over Union (CIoU) metric is utilized, which considers not only the overlapping area but also the distance between the center points and the aspect ratio of the boxes [40]. However, CIoU has limitations; it does not fully account for shape differences and variations between anchor boxes and ground truth boxes, potentially leading to undesirable convergence behavior [41]. Furthermore, in scenarios where the anchor box and the ground truth box do not overlap, merely increasing the size of the anchor box can lead to a reduction in CIoU loss, which is an unreasonable outcome. Consequently, during model training, CIoU may fail to adequately represent the differences between bounding boxes, resulting in



decreased model generalization and slower convergence rates. To address this limitation and improve detection accuracy, this study introduces Powerful-IoU (PIoU) for optimization [42]. The loss function for PIoU is defined as shown in Equations 6, 7.

$$P = \left(\frac{w_p^{gt}}{w_{gt}} + \frac{w_p}{w_{gt}} + \frac{h_p^{gt}}{h_{gt}} + \frac{h_p}{h_{gt}} \right) / 4 \quad (6)$$

$$L_{PIoU} = L_{IoU} + 1 - e^{-P^2} \quad (7)$$

In this equation, w_p^{gt} , w_p , h_p^{gt} , h_p represent the absolute distances between the edges of the anchor box and the target box, while w_{gt} and h_{gt} note the width and height of the target box, as shown in Figure 6. PIoU incorporates a penalty factor that utilizes the size of the target box as the denominator, along with a function that adjusts based on the quality of the anchor box. This approach effectively directs the anchor box to regress along a more efficient trajectory, leading to accelerated model convergence and enhanced detection accuracy.

Although the new loss term in PIoU contributes to accelerating model convergence, it has inherent limitations in adapting to different types of detectors and detection tasks. To address these issues, we introduce Inner-IoU to mitigate the common problems

of weak and slow convergence in various detection tasks. Inner-IoU, by utilizing additional scalable bounding boxes, effectively overcomes the shortcomings in generalization ability of existing methods, thereby enhancing the overall model performance [43]. The parameters and operational mechanism are shown in Figure 9. The calculation method for Inner-IoU is as shown in Equation 8.

$$IoU^{Inner} = \frac{inter}{union} \quad (8)$$

The calculation methods for *inter* and *union* are as shown in Equations 9, 10.

$$inter = (\min(b_r^{gt}, b_r) - \max(b_l^{gt}, b_l)) * (\min(b_b^{gt}, b_b) - \max(b_t^{gt}, b_t)) \quad (9)$$

$$union = (w^{gt} * h^{gt}) * (ratio)^2 + (w * h) * (ratio)^2 - inter \quad (10)$$

The definitions of b_r^{gt} , b_r , b_l^{gt} , b_l , b_b^{gt} , b_b , b_t^{gt} and b_t as shown in Equations 11–14.

$$b_l^{gt} = x_c^{gt} - \frac{w^{gt} * ratio}{2}, b_r^{gt} = x_c^{gt} + \frac{w^{gt} * ratio}{2} \quad (11)$$

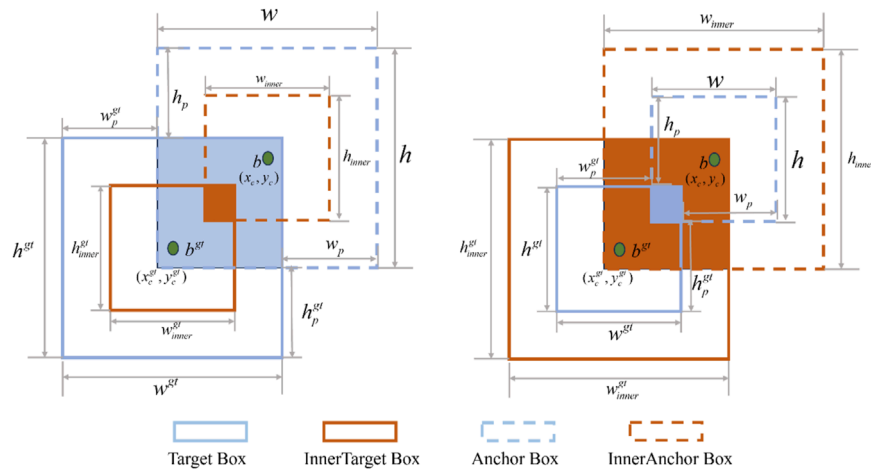


FIGURE 6
Factors of Inner-PIoU.

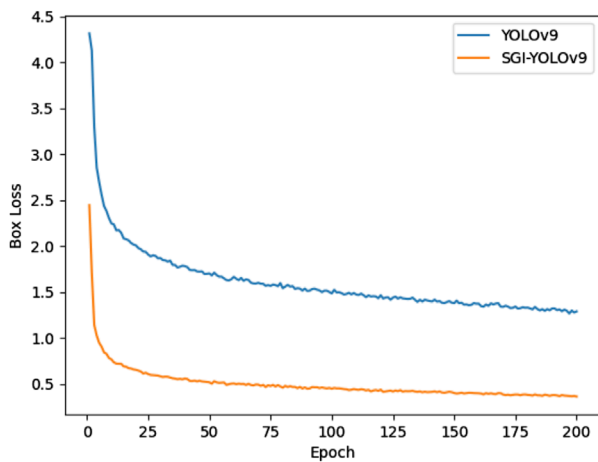


FIGURE 7
The box loss curves of SGI-YOLOv9 and original YOLOv9 models.

range, benefiting regression in cases of low IoU. Finally, we propose a novel computation method called Inner-PIoU, which combines the advantages of Powerful-IoU and Inner-IoU, fully accounting for the differences between bounding boxes. This method not only enhances the model's generalization ability and improves detection accuracy for small objects, but also reduces unexpected convergence behaviors. The formula for Inner-PIoU is shown in Equation 15.

$$L_{\text{Inner-PIoU}} = L_{\text{PIoU}} + \text{IoU} - \text{IoU}^{\text{Inner}} \quad (15)$$

3 Experimental results

To evaluate the efficacy of the proposed SGI-YOLOv9 method, we performed training and testing using the CPDN dataset as well as the PASCAL VOC 2007 dataset, followed by a comparative analysis against other state-of-the-art object detection models. This chapter offers a comprehensive overview of the experimental procedures and implementation details.

3.1 Implementation details

3.1.1 Experimental environment

All experiments were conducted under a consistent computational environment. The system specifications used in our experiments are as follows: a 15-core Intel(R) Xeon(R) Platinum 8358P CPU operating at 2.60 GHz, and an NVIDIA GeForce RTX 3090 GPU. The system ran on Ubuntu 20.04 with PyTorch 1.11.0 and CUDA 11.3. The memory capacity was 24 GB, and Python version 3.8 was employed throughout the experiments.

3.1.2 Training and evaluation metric

3.1.2.1 Training

During the model training phase, we configured the momentum parameter to 0.9 and set the weight decay coefficient to $5e-4$,

$$b_t^{gt} = y_c^{gt} - \frac{h^{gt} * \text{ratio}}{2}, b_b^{gt} = y_c^{gt} + \frac{h^{gt} * \text{ratio}}{2} \quad (12)$$

$$b_l = x_c - \frac{w * \text{ratio}}{2}, b_r = x_c + \frac{w * \text{ratio}}{2} \quad (13)$$

$$b_t = y_c - \frac{h * \text{ratio}}{2}, b_b = y_c + \frac{h * \text{ratio}}{2} \quad (14)$$

The center point of the anchor box is (x_c, y_c) , with its width and height denoted as w , and h , respectively. The center point of the target box is (x_c^{gt}, y_c^{gt}) , with its width and height represented by w^{gt} and h^{gt} . The *ratio* is the scaling factor, typically ranging from 0.5 to 1.5. When the *ratio* is less than 1, the auxiliary bounding box is smaller than the actual bounding box, narrowing the effective regression range, but the absolute value of its gradient is larger than that obtained from IoU loss. Conversely, when the *ratio* is greater than 1, the enlarged auxiliary bounding box expands the effective regression

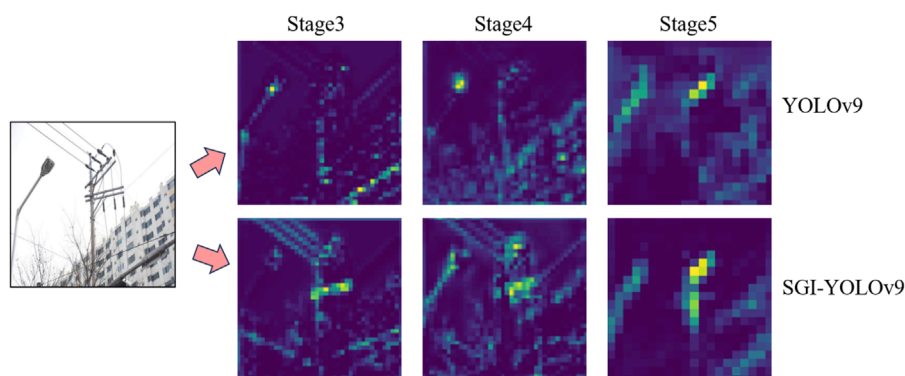


FIGURE 8
Visualization results of the feature maps.

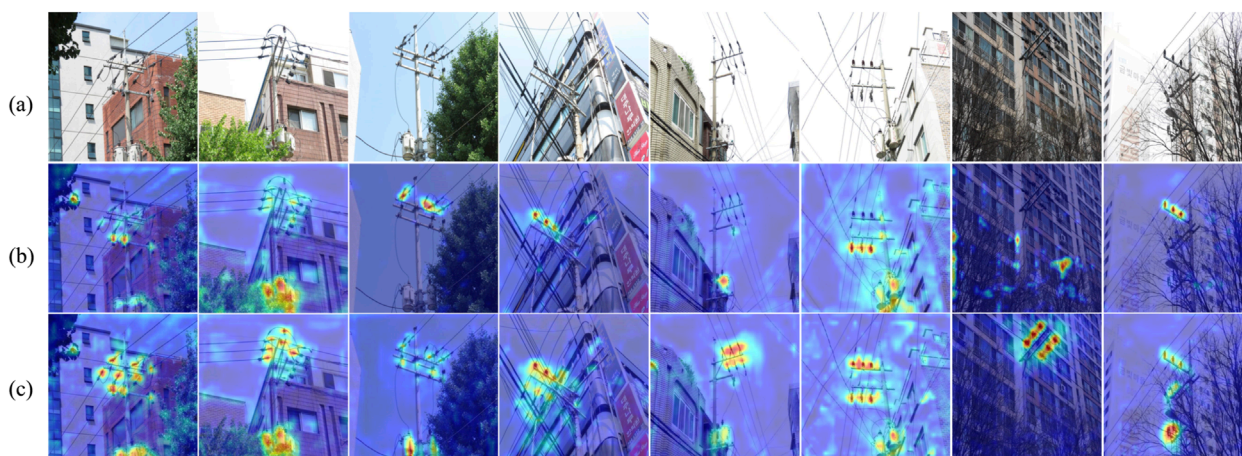


FIGURE 9
XGrad-CAM heatmap of YOLOv9 and SGI-YOLOv9.

employing stochastic gradient descent (SGD) as the optimization algorithm. The batch size was consistently maintained at 32, with a total of 200 training epochs and an initial learning rate of 0.01. Additionally, auxiliary training strategies were implemented during training; however, these strategies were not applied during the inference phase.

3.1.2.2 Evaluation Metric

This paper employs commonly used evaluation metrics in the field of object detection, including precision (P), recall (R), and mean average precision (mAP). These metrics are used to assess the effectiveness and accuracy of component detection in the power distribution network. Higher values indicate better model performance. The calculation of these metrics involves the following parameters: TP (true positives, where the prediction is positive and the actual label is also positive), FP (false positives, where the prediction is positive but the actual label is negative), and FN (false negatives, where the prediction is negative but the actual label is positive).

In object detection tasks, precision measures the degree of false positives produced by the algorithm. A higher precision indicates fewer false detections. The calculation formula is shown in Equation (16):

$$Precision = \frac{TP}{TP + FP} \quad (16)$$

In object detection tasks, recall measures the degree of missed detections by the algorithm. A higher recall indicates fewer missed detections. The calculation formula is shown in Equation (17):

$$Recall = \frac{TP}{TP + FN} \quad (17)$$

The evaluation of an object detection algorithm's performance should encompass both precision and recall metrics. By varying the confidence thresholds, corresponding precision and recall values can be derived, which are subsequently plotted to create a Precision-Recall (PR) curve, with precision represented on the vertical axis and recall on the horizontal axis. The area enclosed by the PR curve and the coordinate axes indicates the Average Precision (AP). If we

TABLE 1 Experimental results for different ratios in Inner-PIoU.

Method	Precision (%)	Recall (%)	mAP@50 (%)	mAP50-95 (%)
CIoU	80.6	70.1	75.2	45.0
Inner-PIoU (ratio = 0.9)	81.5	69.5	75.9	45.5
Inner-PIoU (ratio = 1.0)	81.1	69.6	75.8	45.5
Inner-PIoU (ratio = 1.1)	82.4	70.5	76.4	46.5
Inner-PIoU (ratio = 1.2)	82.2	69.7	76.1	45.9

denote the function associated with this curve as $p(r)$, the formula for AP is presented in Equation (18):

$$AP = \int_0^1 p(r) dr \quad (18)$$

Mean Average Precision (mAP) is calculated by determining the AP values for all target categories and then computing their average. The formula for mAP is provided in Equation (19):

$$mAP = \sum_{n=1}^N AP(n) / N \quad (19)$$

3.2 Ablation study

To determine the optimal ratio parameter for Inner-PIoU in detecting crucial components in the CPDN dataset, we conduct a series of experiments and compare the results with the CIoU used in original YOLOv9, as shown in Table 1. When the ratio is set to 1.0, indicating no auxiliary bounding box and only using Powerful-IoU, the results show a 0.6% increase in mAP@50, validating the effectiveness of Powerful-IoU in detecting crucial components in the power distribution network. When the ratio is set to 0.9, which introduces a smaller auxiliary bounding box, there is no significant improvement in mAP@50 compared to the ratio of 1.0. However, when the ratio exceeds 1.0, indicating the use of a larger auxiliary bounding box, the performance improves. Specifically, with ratio values of 1.1 and 1.2, mAP@50 increases by 0.6% and 0.3%, respectively, compared to a ratio of 1.0. Since most components in the CPDN dataset are considered small objects, further experiments demonstrate that when the ratio exceeds 1.0, the convergence of the model training for small object detection improves significantly. Consequently, this leads to a notable enhancement in detection accuracy. Therefore, we select a ratio of 1.1 for Inner-PIoU as the optimal parameter and use it as the default in subsequent experiments.

We further analyze and compare the box loss of the improved YOLOv9 with the original YOLOv9, as shown in Figure 7. The curve shows that the loss for the improved YOLOv9 is significantly lower than that of the original YOLOv9 during the initial training phase, indicating that the improved YOLOv9 adapts to the data more quickly. As the training epochs progress, both models exhibit a rapid decline in loss, but the improved YOLOv9 demonstrates a much faster decrease. This indicates that the

improved YOLOv9 learns the positions of bounding boxes more efficiently and reduces the deviation between the predicted and actual boxes more effectively. After both models converge, the loss for the SGI-YOLOv9 consistently remains lower than that of the original YOLOv9. These findings confirm that the SGI-YOLOv9, with the incorporation of Inner-PIoU, adapts to the dataset faster, achieves lower loss values during training, and converges more quickly.

Next, we conduct ablation experiments on each of the proposed modules, as shown in Table 2. The results demonstrate that each module contributes to improving the accuracy of crucial components recognition in the power distribution network. Specifically, when the ratio is set to 1.1, Inner-PIoU improves accuracy by 1.2% on the CPDN test set, while SPDConv++ and GCFM contribute improvements of 1.6% and 1.1%, respectively. These findings further validate that the proposed methods enhance the accuracy of crucial components recognition in the power distribution network effectively.

As shown in Figure 8, we compare the feature maps extracted at various stages of the backbone network between the original YOLOv9 model and the SGI-YOLOv9 model. Through feature map visualization, it is evident that after introducing the SPDConv++ method, the improved model exhibits a stronger response to edge information of crucial components in the power distribution network. Particularly in the Stage 3, Stage 4, and Stage 5 phases of the backbone network, the SGI-YOLOv9 model significantly reduces the loss of fine-grained features, preserving more detailed information. These results indicate that the SPDConv++ method effectively enhances the richness of fine-grained features in the backbone network output feature maps, further validating the effectiveness and robustness of this method in object detection from the perspective of feature visualization.

Additionally, we utilize XGrad-CAM [44] to perform a visual analysis of the attention heatmaps for both the original YOLOv9 and the SGI-YOLOv9 models, as shown in Figure 9. In this figure, (a) represents the input image, (b) shows the attention heatmap from the original YOLOv9 model, and (c) displays the attention heatmap from the SGI-YOLOv9 model. The visualization results clearly indicate that the SGI-YOLOv9 model significantly improves its focus on key components of the power transmission lines in complex backgrounds. This highlights the notable advantage of SGI-YOLOv9 in enhancing detection accuracy, particularly in complex scenes involving crucial components of the power distribution

TABLE 2 Ablation experiments for the SGI-YOLOv9 method.

Method	Inner-PloU	SPDConv++	GCFM	mAP@50 (%)	mAP50-95 (%)
YOLOv9	—	—	—	75.2	45.0
SGI-YOLOv9	✓	—	—	75.2	46.5
	✓	✓	—	78.0	47.7
	✓	✓	✓	79.1	48.5

TABLE 3 Comparison results of different models.

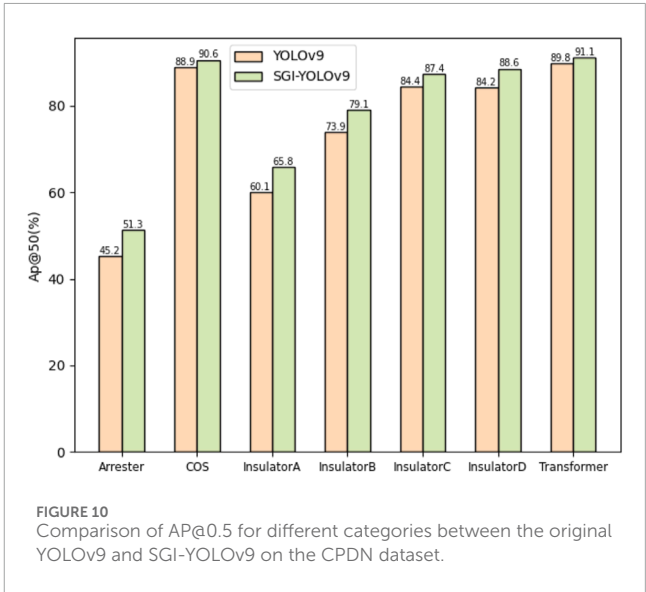
Method	Precision (%)	Recall (%)	mAP@50 (%)	map50-95 (%)
YOLOv5 [18]	83.7	67.9	72.3	40.0
YOLOv7 [19]	79.4	66.1	71.8	37.7
YOLOv8 [20]	82.1	68.3	74.8	44.0
YOLOv9 [21]	80.6	70.0	75.2	45.0
SGI-YOLOv9	85.2	72.3	79.1	48.5

network, further validating its effectiveness and reliability in practical applications.

3.3 Compare with state-of-arts on CPDN dataset

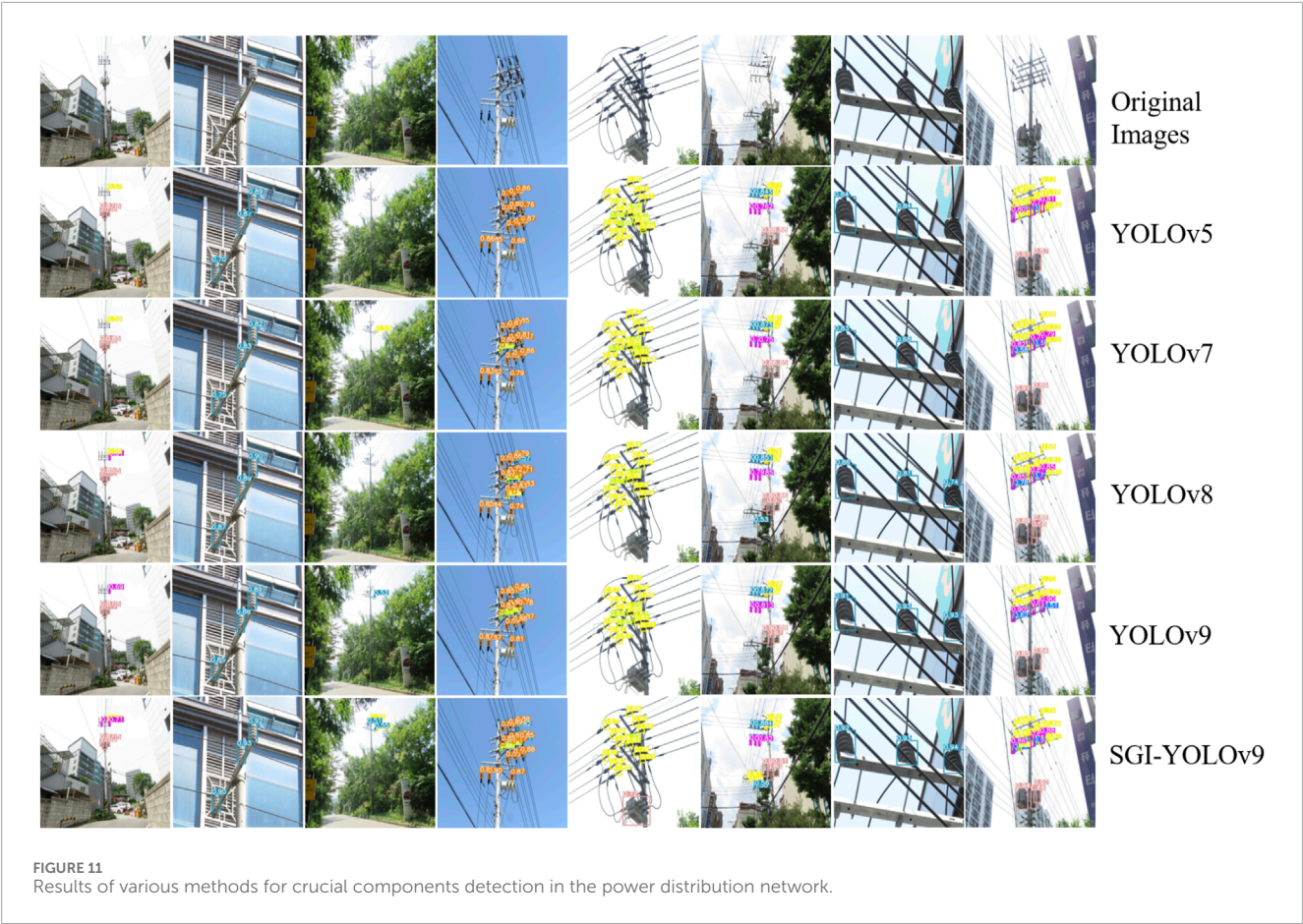
To ensure a fair comparison between our proposed SGI-YOLOv9 method and other mainstream object detection methods on the CPDN dataset, we train all models without loading any pre-trained models. Table 3 presents the comparative experimental results of SGI-YOLOv9 and other mainstream object detection methods on the test set. As shown in the table, SGI-YOLOv9 achieves the highest scores across all evaluation metrics on the CPDN test set, with a mAP@50 of 79.1%, which is a 3.9% improvement over the original YOLOv9. This demonstrates that our SGI-YOLOv9 method offers a significant advantage in detecting crucial components within the complex background of the power distribution network.

To further evaluate the effectiveness of the SGI-YOLOv9 algorithm in detecting different types of components in the power distribution network, we record the AP@50 for seven component types in the dataset, as shown in Figure 10. It is evident that the SGI-YOLOv9 model consistently outperforms the original YOLOv9 in terms of overall AP@50. Specifically, SGI-YOLOv9 demonstrates stable performance improvements when detecting larger components such as COS and Transformers, with increases of 1.7% and 1.3%, respectively. For smaller components, such as Arresters and Insulators, the improvements are even more significant. Notably, SGI-YOLOv9 achieves a 6.1% increase in AP@50 for Arresters, marking the most substantial gain. Additionally, the mAP@50 for the four types of Insulators increases



by 4.58%. These results confirm the significant improvement of SGI-YOLOv9 in detecting small objects, highlighting its enhanced ability to focus on and handle small objects in complex scenes.

To comprehensively validate the effectiveness of the proposed SGI-YOLOv9 method, we compare its visualization results for crucial components detection with those of other mainstream object detection algorithms, as shown in Figure 11. It is evident that YOLOv5, YOLOv7, YOLOv8, and YOLOv9 all exhibit varying degrees of omission and false detections. This is especially pronounced when the crucial components are small or occluded, where other mainstream models demonstrate low confidence in



their predicted bounding boxes, leading to numerous missed detections and false positives. In contrast, our SGI-YOLOv9 model shows higher detection accuracy when handling small and occluded crucial components. These findings demonstrate that SGI-YOLOv9 is highly effective for crucial component detection tasks in the complex environments of the power distribution network.

3.4 Compare with state-of-arts on the PASCAL VOC 2007 dataset

To further validate the effectiveness of the proposed SGI-YOLOv9 model in object detection tasks, we conducted training experiments on the PASCAL VOC 2007 dataset and systematically compared its performance on the test set with several mainstream object detection algorithms. Notably, all models utilized in the comparison were lightweight versions. As shown in Table 4, the SGI-YOLOv9 model achieved a significant performance improvement, attaining a mAP@50 value of 63.3%, which represents a 1.6% increase compared to the original YOLOv9. Additionally, the precision improved by 1.4%, and the recall increased by 1.3% over the original YOLOv9. These results demonstrate that SGI-YOLOv9 not only delivers superior accuracy in insulator defect detection tasks but also excels in general-purpose object detection tasks. This highlights the model's robustness, algorithmic superiority, and strong generalization capability across diverse application scenarios.

TABLE 4 Experimental Results of Different Models on the PASCAL VOC 2007 dataset.

	Precision (%)	Recall (%)	mAP@50 (%)
Faster-RCNN	34.1	54.7	57.5
Mask-RCNN	33.9	69.1	57.2
YOLOv5	69.4	52.9	60.3
YOLOv7	66.8	52.5	58.3
YOLOv8	68.8	53.0	56.5
YOLOv9	66.7	54.1	61.7
SGI-YOLOv9	68.1	55.4	63.3

4 Conclusion

This paper presents an improved method based on YOLOv9 to address the challenges of small objects detection and complex scenarios in the detection of crucial components in the power distribution network. By designing the SPDConv++ module, we reduce the loss of fine-grained feature information and improve the accuracy in detecting small objects. Simultaneously, the proposed

global context fusion module models global information from both spatial and channel dimensions, effectively handling complex backgrounds and occlusion issues. Additionally, we optimized the loss function of IoU in YOLOv9 by proposing the Inner-PIoU method, which combines the advantages of Powerful-IoU and Inner-IoU to enhance the regression performance of the bounding boxes, thereby improving the model's generalization ability and detection accuracy for crucial components in the power distribution network. Experimental results demonstrate the effectiveness of SGI-YOLOv9, achieving an mAP@50 of 79.1% on the CPDN dataset, an improvement of 3.9% over the original YOLOv9, and an mAP@50 of 63.3% on the PASCAL VOC 2007 dataset, surpassing YOLOv9 by 1.6%. The proposed method provides effective technical support for detecting crucial components in the power distribution network under complex scenarios, contributing to the safety and reliability of power grid. Future research may focus on further optimizing the model's computational efficiency and applying it to more power system scenarios to promote the development of smart grids.

Data availability statement

The original contributions presented in the study are included in the article/supplementary material, further inquiries can be directed to the corresponding author.

Author contributions

MY: Conceptualization, Methodology, Writing—original draft. BC: Methodology, Software, Writing—original draft. CL: Validation, Writing—review and editing. WY: Funding acquisition, Writing—review and editing. YL: Visualization, Writing—review and editing.

References

1. Yang L, Fan J, Liu Y, Li E, Peng J, Liang Z. A review on state-of-the-art power line inspection techniques. *IEEE Trans Instrumentation Meas* (2020) 69:9350–65. doi:10.1109/tim.2020.3031194
2. Siddiqui ZA, Park U. A drone based transmission line components inspection system with deep learning technique. *Energies* (2020) 13:3348. doi:10.3390/en13133348
3. Zhang K, Yang L. *Insulator segmentation algorithm based on k-means*. Chinese Automation Congress CAC (2019) p. 4747–51.
4. Yu Y, Cao H, Wang Z, Li Y, Li K, Xie S. Texture-and-shape based active contour model for insulator segmentation. *IEEE Access* (2019) 7:78706–14. doi:10.1109/access.2019.2922257
5. Zhao Z, Liu N, Wang L. Localization of multiple insulators by orientation angle detection and binary shape prior knowledge. *IEEE Trans Dielectrics Electr Insul* (2015) 22:3421–8. doi:10.1109/tdei.2015.004741
6. He X, Qi G, Zhu Z, Li Y, Cong B, Bai L. Medical image segmentation method based on multi-feature interaction and fusion over cloud computing. *Simulation Model Pract Theor* (2023) 126:102769. doi:10.1016/j.simpat.2023.102769
7. Zhu Z, Wang S, Gu S, Li Y, Li J, Shuai L, et al. Driver distraction detection based on lightweight networks and tiny object detection. *Math biosciences Eng* (2023) 20:18248–66. doi:10.3934/mbe.2023811
8. Huang X, Wang S, Qi G, Zhu Z, Li Y, Shuai L, et al. Driver distraction detection based on cloud computing architecture and lightweight neural network. *Mathematics* (2023) 11:4862. doi:10.3390/math11234862
9. Girshick R, Donahue J, Darrell T, Malik J. Rich feature hierarchies for accurate object detection and semantic segmentation. In: *Proceedings*

Funding

The author(s) declare that financial support was received for the research, authorship, and/or publication of this article. The author(s) declare financial support was received for the research, authorship, and/or publication of this article. This work was supported by The State Grid Corporation Headquarters Science and Technology Project: Research on key Technologies of Aerial Vehicle Dock Replenishment for Transmission Line(5500-202321166A-1-1-ZN). The funder was not involved in the study design, collection, analysis, interpretation of data, the writing of this article, or the decision to submit it for publication.

Conflict of interest

Authors MY, BC, CL, WY and YL declare that they were employed by State Grid Corporation.

Generative AI statement

The author(s) declare that no Generative AI was used in the creation of this manuscript.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

- of the IEEE conference on computer vision and pattern recognition (2014) p. 580–7.
10. Ren S, He K, Girshick R, Sun J. Faster r-cnn: towards real-time object detection with region proposal networks. *IEEE Trans pattern Anal machine intelligence* (2016) 39:1137–49. doi:10.1109/tpami.2016.2577031
11. He K, Gkioxari G, Dollár P, Girshick R. Mask r-cnn. In: *Proceedings of the IEEE international conference on computer vision* (2017) p. 2961–9.
12. Zhao Z, Zhen Z, Zhang L, Qi Y, Kong Y, Zhang K. Insulator detection method in inspection image based on improved faster r-cnn. *Energies* (2019) 12:1204. doi:10.3390/en12071204
13. Odo A, McKenna S, Flynn D, Vorstius JB. Aerial image analysis using deep learning for electrical overhead line network asset management. *IEEE Access* (2021) 9:146281–95. doi:10.1109/access.2021.3123158
14. Dong C, Zhang K, Xie Z, Shi C. An improved cascade rcnn detection method for key components and defects of transmission lines. *IET Generation, Transm & Distribution* (2023) 17:4277–92. doi:10.1049/gtd2.12948
15. Redmon J. You only look once: unified, real-time object detection. In: *Proceedings of the IEEE conference on computer vision and pattern recognition* (2016).
16. Redmon J. Yolov3: an incremental improvement. *arXiv preprint arXiv:1804.02767* (2018).
17. Bochkovskiy A, Wang CY, Liao HYM. Yolov4: optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934* (2020).

18. Wu W, Liu H, Li L, Long Y, Wang X, Wang Z, et al. Application of local fully convolutional neural network combined with yolo v5 algorithm in small target detection of remote sensing image. *PloS one* (2021) 16:e0259283. doi:10.1371/journal.pone.0259283
19. Wang CY, Bochkovskiy A, Liao HYM. YOLOv7: trainable bag-of-freebies sets new state-of-the-art for real-time object detectors. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2023) p. 7464–75.
20. Wang L, Zhang G, Wang W, Chen J, Jiang X, Yuan H A defect detection method for industrial aluminum sheet surface based on improved yolov8 algorithm. *Front Phys* (2024) 12:1419998. doi:10.3389/fphy.2024.1419998
21. Wang CY, Yeh IH, Liao HYM. YOLOv9: learning what you want to learn using programmable gradient information. In: *arXiv preprint arXiv:2402* (2024) p. 13616.
22. Liu W, Anguelov D, Erhan D, Szegedy C, Reed S, Fu CY, et al. Ssd: single shot multibox detector. In: *Computer vision—ECCV 2016: 14th European conference, Amsterdam, The Netherlands, October 11–14, 2016, proceedings, Part I* 14. Springer (2016) p. 21–37.
23. Lv H, Du Y, Ma Y, Yuan Y. Object detection and monocular stable distance estimation for road environments: a fusion architecture using yolov8 and abnormal jumping change filter. *Electronics* (2024) 13:3058. doi:10.3390/electronics13153058
24. Qi C, Chen Z, Chen X, Bao Y, He T, Hu S, et al. Efficient real-time detection of electrical equipment images using a lightweight detector model. *Front Energy Res* (2023) 11:1291382. doi:10.3389/fenrg.2023.1291382
25. Siddiqui ZA, Park U, Lee SW, Jung NJ, Choi M, Lim C, et al. Robust powerline equipment inspection system based on a convolutional neural network. *Sensors* (2018) 18:3837. doi:10.3390/s18113837
26. Liu Z, Wu G, He W, Fan F, Ye X. Key target and defect detection of high-voltage power transmission lines with deep learning. *Int J Electr Power & Energy Syst* (2022) 142:108277. doi:10.1016/j.ijepes.2022.108277
27. Qiu Z, Zhu X, Liao C, Shi D, Qu W. Detection of transmission line insulator defects based on an improved lightweight yolov4 model. *Appl Sci* (2022) 12:1207. doi:10.3390/app12031207
28. Liu M, Li Z, Li Y, Liu Y. A fast and accurate method of power line intelligent inspection based on edge computing. *IEEE Trans Instrumentation Meas* (2022) 71:1–12. doi:10.1109/tim.2022.3152855
29. Liu C, Ma L, Sui X, Guo N, Yang F, Yang X, et al. Yolo-csm-based component defect and foreign object detection in overhead transmission lines. *Electronics* (2023) 13:123. doi:10.3390/electronics13010123
30. Chen B, Zhang W, Wu W, Li Y, Chen Z, Li C. Id-yolov7: an efficient method for insulator defect detection in power distribution network. *Front Neurorobotics* (2024) 17:1331427. doi:10.3389/fnbot.2023.1331427
31. He M, Qin L, Deng X, Liu K. Mfi-yolo: multi-fault insulator detection based on an improved yolov8. *IEEE Trans Power Deliv* (2023) 39:168–79. doi:10.1109/tpwrd.2023.3328178
32. Zhu Z, Zheng R, Qi G, Li S, Li Y, Gao X. Small object detection method based on global multi-level perception and dynamic region aggregation. *IEEE Trans Circuits Syst Video Technology* (2024) 34:10011–22. doi:10.1109/tcsvt.2024.3402097
33. Qi G, Zhang Y, Wang K, Mazur N, Liu Y, Malaviya D. Small object detection method based on adaptive spatial parallel convolution and fast multi-scale fusion. *Remote Sensing* (2022) 14:420. doi:10.3390/rs14020420
34. Li Y, Zhou Z, Qi G, Hu G, Zhu Z, Huang X. Remote sensing micro-object detection under global and local attention mechanism. *Remote Sensing* (2024) 16:644. doi:10.3390/rs16040644
35. Zhang T, Zhang Y, Xin M, Liao J, Xie Q. A light-weight network for small insulator and defect detection using uav imaging based on improved yolov5. *Sensors* (2023) 23:5249. doi:10.3390/s23115249
36. University H. All image dataset (2023). Available from: <https://universe.roboflow.com/hanshin-university/allimage> (Accessed October 17, 2024).
37. Sunkara R, Luo T. No more strided convolutions or pooling: a new cnn building block for low-resolution images and small objects. In: *Joint European conference on machine learning and knowledge discovery in databases*. Springer (2022). p. 443–59.
38. Vaswani A. Attention is all you need. *Adv Neural Inf Process Syst* (2017).
39. Dosovitskiy A. An image is worth 16x16 words: transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).
40. Zheng Z, Wang P, Liu W, Li J, Ye R, Ren D. Distance-iou loss: faster and better learning for bounding box regression. *Proc AAAI Conf Artif intelligence* (2020) 34:12993–3000. doi:10.1609/aaai.v34i07.6999
41. Zhang YF, Ren W, Zhang Z, Jia Z, Wang L, Tan T. Focal and efficient iou loss for accurate bounding box regression. *Neurocomputing* (2022) 506:146–57. doi:10.1016/j.neucom.2022.07.042
42. Liu C, Wang K, Li Q, Zhao F, Zhao K, Ma H. Powerful-iou: more straightforward and faster bounding box regression loss with a nonmonotonic focusing mechanism. *Neural Networks* (2024) 170:276–84. doi:10.1016/j.neunet.2023.11.041
43. Zhang H, Xu C, Zhang S. Inner-iou: more effective intersection over union loss with auxiliary bounding box. *arXiv preprint arXiv:2311.02877* (2023).
44. Fu R, Hu Q, Dong X, Guo Y, Gao Y, Li B. Axiom-based grad-cam: towards accurate visualization and explanation of cnns. *arXiv preprint arXiv:2008.02312* (2020).



OPEN ACCESS

EDITED BY

Yu Liu,
Hefei University of Technology, China

REVIEWED BY

Marian Gaiceanu,
Dunarea de Jos University, Romania
Yi He,
Chengdu University of Information
Technology, China

*CORRESPONDENCE

Tongxin Yang,
✉ yangtx0704@163.com

RECEIVED 26 November 2024

ACCEPTED 03 March 2025

PUBLISHED 24 March 2025

CITATION

Ren S, Yang T, Luo J, Wu G, Mao K and Liu B
(2025) Performance evaluation of
photovoltaic scenario generation.
Front. Phys. 13:1534629.
doi: 10.3389/fphy.2025.1534629

COPYRIGHT

© 2025 Ren, Yang, Luo, Wu, Mao and Liu. This
is an open-access article distributed under
the terms of the [Creative Commons
Attribution License \(CC BY\)](https://creativecommons.org/licenses/by/4.0/). The use,
distribution or reproduction in other forums is
permitted, provided the original author(s) and
the copyright owner(s) are credited and that
the original publication in this journal is cited,
in accordance with accepted academic
practice. No use, distribution or reproduction
is permitted which does not comply with
these terms.

Performance evaluation of photovoltaic scenario generation

Siyu Ren¹, Tongxin Yang^{2*}, Jun Luo², Gang Wu^{3,4}, Kai Mao⁵ and Bowen Liu²

¹The Open University of Chengdu, Chengdu, Sichuan, China, ²School of Computer Science and Engineering, Chongqing University of Science and Technology, Chongqing, China, ³Chongqing Carbon Energy Technology Co., Ltd., Chongqing, China, ⁴Sichuan Aizhong Comprehensive Energy Technology Service Co., Ltd., Guangan, Sichuan, China, ⁵School of Artificial Intelligence, Chongqing Technology and Business University, Chongqing, China

Photovoltaic scenario generation plays a critical role in power systems characterized by high diversity and fluctuation. Despite recent theoretical advancements, effectively evaluating the performance of photovoltaic scenario generation remains a significant challenge. Existing studies predominantly rely on metrics such as mean, variance, and probability density functions for assessment. However, these approaches struggle to disentangle the underlying mechanisms of morphological features and environmental stochastic factors (e.g., cloud cover, seasonal variations) from individual or batch-generated samples. To address these limitations, this paper proposes an evaluation framework based on the wide-sense stationary process. By analyzing historical photovoltaic scenario data, a solar irradiance distribution model is first constructed to characterize its dynamic behavior. Subsequently, an autoregressive model is employed to quantify the influence of environmental randomness on photovoltaic scenarios. The proposed evaluation model not only comprehensively validates the reliability of various photovoltaic scenario generation techniques but also identifies the corresponding month or season of generated samples through scenario feature analysis. Experimental results demonstrate that, compared to conventional probability-based metrics, the proposed model more effectively reveals the performance characteristics of photovoltaic scenario generation technologies. This advancement provides a novel technical foundation for optimizing photovoltaic scenario generation in practical power systems.

KEYWORDS

photovoltaic scenario generation, wide-sense stationary process, autoregressive model, environmental randomness analysis, performance benchmarking for PV systems

1 Introduction

Electricity produced from solar photovoltaic (PV) panels is a vital source of clean energy, where much research has been done in recent years owing to its low pollution. As integration of PV powers into traditional power grid increases, a challenge surfaces due to regulation requirements of balancing existing supply-demand in energy markets [1, 2]. One solution to this problem is PV energy prediction in solar reception process using PV scenario generation to simulate real PV energy [3–5]. Key to informative analysis of PV scenario is accurate representation that describes the variability and uncertainty of solar generation from both spatial and temporal aspects [6, 7]. Two attributes of PV systems make it difficult to generate reliable PV scenario. One is that solar generation is mostly

dominated by the accessibility of solar irradiation that changes across a day and during a year following the movement round the Sun and the Earth's rotation [8–10]. Compared with traditional energy generation techniques, PV systems involve strong uncertainty and environmental variables (e.g., cloud and season). The other is PV systems contain highly scalable ranging from a few kilowatts (kW) to hundreds of megawatts (MW) [11], which indicates the diversity of PV scenarios. Thereby, exploring the most accurate information on the PV generation characteristics determines the reliability of generated PV scenarios.

A lot of studies have devoted to improve the effectiveness and precision of PV scenario generation [12], categorized as model-based and model-free approaches. The former represents solar reception process with a specific model of presuppositions, and the accuracy of model determines the reliability of PV scenario generation. Most of these techniques utilize probability models to simulate solar reception process. A modular statistical modeling approach is presented to predict power generation of both PV and wind power systems [13]. A pseudo-random number generation technique is proposed to reduce prediction error of PV scenario generation with considering uncertainty and variability indices [14]. Aggregated power curves are also analyzed and contribute to PV generation [15]. Gaussian copulas are established to produce multivariate PV scenarios [16]. The advantages of these techniques are simple models and easy implementation, yet the disadvantage is its limited ability of representing the uncertainty. Overcoming the weakness of model-based methods, model-free techniques learn the inherent distribution of solar reception process through analyzing existing PV scenario data. Support vector regression is used to forecast regional PV power generation with past PV data [17]. Artificial neural network is built to predict the solar irradiance in PV systems [18]. A recurrent neural network (RNN) is utilized to produce PV scenarios with month and weather information, which requires no mathematical modelling [19]. Conditional generative adversarial networks (CGAN) and Wasserstein GAN (WGAN) are constructed to generate PV scenarios with sufficient diversity and good representation of environmental uncertainty [20]. These techniques are capable of capturing the details of PV scenarios caused by particular operation of solar receptions. The drawbacks of model-free techniques are severe computation complexity and unexplained characteristics in operation processes. While either approach claims its accuracy and effectiveness, there are no standards evaluating them.

Although existing methods for PV scenario generation have demonstrated certain advantages, they also face notable limitations. Model-based approaches, such as AR and autoregressive moving average (ARMA) models, rely on fixed assumptions about the underlying process, making it challenging to capture non-linear and abrupt changes in PV power caused by sudden weather shifts. Model-free approaches, such as generative adversarial networks (GANs) and other deep learning models, have greater flexibility but require large datasets and significant computational resources. Their generalization performance may also be affected when applied to unseen weather or seasonal conditions.

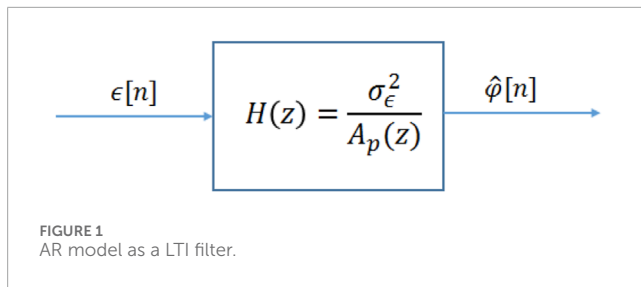
Hence, it is important to evaluate the effectiveness of PV scenario generation methods. Typically, probabilistic properties of produced PV scenarios (e.g., mean value, variance, and probability density function (pdf)) are common indicators of evaluating the

performance of PV scenario generation, because these indicators are easily calculated. In particular, hourly PV generation is used to estimate variable changes between times in the system and optimize future expansion plans. PV energy curves have high fluctuations due to the diversity and variety of solar reception process, i.e., the differences between any 2 sampling values in a PV scenario could be large. In this regard, mean value is weak to represent the uncertainty (e.g., cloud) in PV scenarios. Furthermore, variance is useful to evaluate the pattern of sunrise and sunset. Only averaging the variances of existing PV data is hard to determine which month or season that a generated sample belongs to. Additionally, pdf stands for purely solar energies from the Sun without environmental randomness, which is hard to assess the influence by cloud.

Therefore, unlike simple probability indicators such as mean and variance, we propose a novel evaluation model that explicitly assesses the reliability and effectiveness of photovoltaic (PV) scenario generation. While recent research, such as those relying on mean and variance, have been widely used, they are limited in their ability to capture the temporal dependencies and environmental randomness that significantly affect PV generation. These methods primarily focus on statistical summaries, which fail to account for abrupt changes caused by environmental factors like cloud movement, leading to less accurate predictions. In contrast, our autoregressive (AR)-based model provides a more comprehensive evaluation framework by explicitly modeling the temporal structure and stochastic fluctuations in PV scenarios. This AR model simulates the movement of clouds and its impact on solar reception, offering a more precise characterization of the underlying randomness in PV generation. Additionally, by incorporating month- and season-specific AR parameters, our approach is capable of categorizing PV scenarios into temporal categories, such as specific months and seasons, something that traditional mean and variance-based methods cannot achieve. This enhanced ability to classify and evaluate PV scenarios allows for more robust and context-aware PV energy predictions, making our approach not only more accurate but also more scalable and practical for real-world applications.

1. Proposed evaluation model has stronger ability of assessing the reliability of generated PV scenarios than simple probability indicators, which is conducive to improving the studies of PV scenario generation and promoting the application of PV systems.
2. Proposed evaluation model can estimate the corresponding month and season that a credible PV sample is geared to. To our knowledge, this is the first work evaluating the specific month and season of a generated PV scenario.
3. We discover the representative properties (e.g., the effect of cloud) of solar reception process and use the discovery to assess the inherent movement of cloud, which fills a gap of estimating the environmental randomness of PV scenario generation.

The rest of this paper is organized as follows. [Section 2](#) motivates the introduction of the details of AR model in both time and frequency domains. [Section 3](#) is devoted to the details of proposed approach. In [Section 4](#), experimental results are discussed and analyzed. Finally, [Section 5](#) concludes this paper.



2 Autoregressive model

In regard to time series events, correlations exist among the behaviors at certain intervals. Considering the correlations, an autoregressive (AR) model is able to predict future variables of interest according to their past values. AR model is basically a linear regression of current values against past values in the same time series [21]. An AR(p) model is mathematically defined as:

$$\varphi[n] = -c - a_1\varphi[n-1] - a_2\varphi[n-2] - \dots - a_p\varphi[n-p] + \epsilon[n] \quad (1)$$

where $\varphi[n]$ is n th value of variable observation; $\epsilon[n]$ is driving noise; p is the order of AR model; $\{a_1, a_2, \dots, a_p\}$ are AR parameters. c denotes the energy of the observation, which is defined as:

$$c = \left(1 - \sum_{t=1}^p a_t\right) \mu \quad (2)$$

where μ is the process mean.

AR parameters $\{a_t | t = 1, 2, \dots, p\}$ dominate the performance of AR model, which could be estimated by several techniques, e.g., Yule-Walker, burg method, Kalman filter, least-square, expectation-maximization, forward-backward. Changing the value of p and the parameter a_t leads to different time series patterns [22].

Taking Z-transforms of Equation 1, AR model as a LTI filter is depicted as:

$$H(z) = \frac{\sigma_\epsilon^2}{A_p(z)} = \frac{\sigma_\epsilon^2}{c + \sum_{t=1}^p a_t z^{-t}} \quad (3)$$

AR(p) model is used to model the observation $\varphi[n]$ as the response of a LTI filter with p order to an input $\epsilon[n]$ (illustrated as in Figure 1). The purpose is to discover the filter coefficients (AR parameters $\{a_t | t = 1, 2, \dots, p\}$) and the input $\epsilon[n]$ that make the estimated $\hat{\varphi}[n]$ as close to $\varphi[n]$ as possible. From Equation 3, AR model is an all-poles model. AR system could be unstable if the poles are outside the unit circle.

The AR parameters could be solved by Yule-Walker equations. Equation 1 can be reformed in a vector form.

$$(\varphi[n], \varphi[n-1], \dots, \varphi[n-p]) \begin{pmatrix} 1 \\ a_1 \\ \vdots \\ a_p \end{pmatrix} = \epsilon[n] - c \quad (4)$$

Multiple both sides by $\varphi[n]$ and take the expectation, we have:

$$E\varphi[n]x[n] = E \left\{ \varphi[n] (\varphi[n], \varphi[n-1], \dots, \varphi[n-p]) \begin{pmatrix} 1 \\ a_1 \\ \vdots \\ a_p \end{pmatrix} \right\} \\ = E\{\epsilon[n] - c\}$$

Then,

$$(r_0, r_1, \dots, r_p) \begin{pmatrix} 1 \\ a_1 \\ \vdots \\ a_p \end{pmatrix} = \sigma_\epsilon^2 \quad (5)$$

where $\{r_0, r_1, \dots, r_p\}$ are autocorrelation function. Then,

$$\begin{bmatrix} r_0 & r_1 & \dots & r_p \\ r_1 & r_0 & \dots & r_{p-1} \\ \vdots & \vdots & \ddots & \vdots \\ r_p & r_{p-1} & \dots & r_0 \end{bmatrix} \begin{bmatrix} 1 \\ a_1 \\ \vdots \\ a_p \end{bmatrix} = \begin{bmatrix} \sigma_\epsilon^2 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (6)$$

While deleting the first equation in Equation 6, the formula, termed Yule-Walker equations, is obtained as follows.

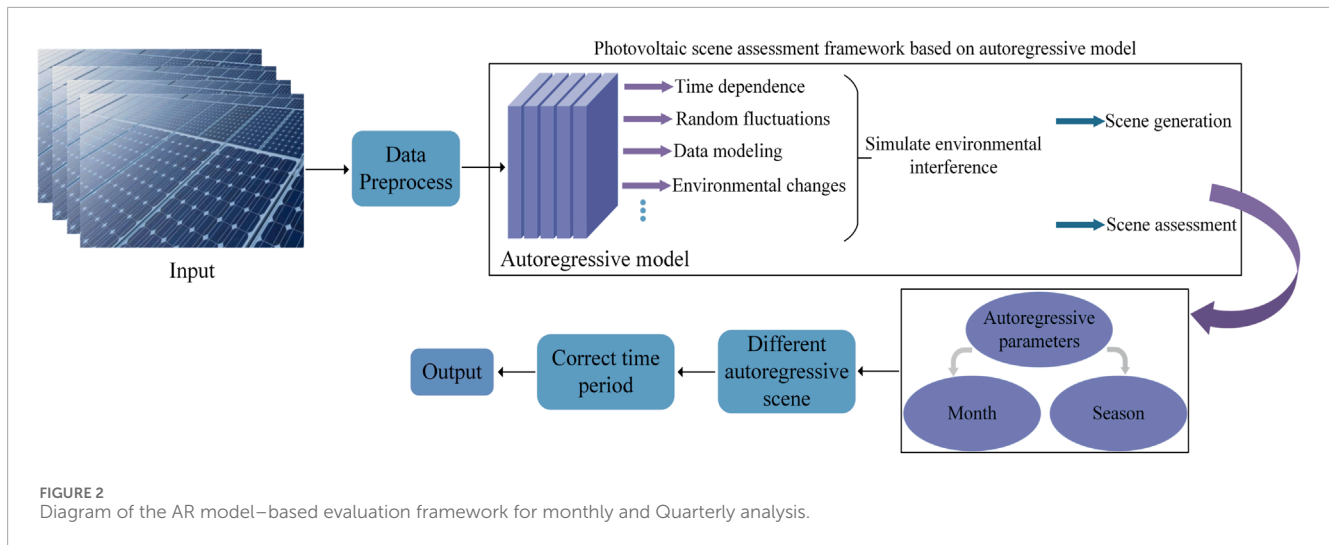
$$\begin{bmatrix} r_0 & r_1 & \dots & r_p \\ r_1 & r_0 & \dots & r_{p-1} \\ \vdots & \vdots & \ddots & \vdots \\ r_{p-1} & r_{p-2} & \dots & r_0 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_p \end{bmatrix} = - \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_p \end{bmatrix} \quad (7)$$

3 The proposed model

Figure 2 illustrates the overall methodology of our proposed autoregressive-based PV scenario assessment framework. By analyzing the time dependency and random fluctuations in PV data, the AR model can capture short-term variations caused by environmental factors, thereby providing a more robust evaluation of PV scenario reliability. This framework introduces an autoregressive process to model environmental disturbances in PV generation, taking into account the temporal dependence of PV output as well as abrupt fluctuations due to environmental changes. Meanwhile, the framework integrates AR model parameters with monthly and seasonal PV generation characteristics. By estimating AR parameters for different months and seasons, it assigns generated PV samples to the appropriate time period, enabling more accurate handling of both seasonal variations and longer-term fluctuations in the PV data.

3.1 PV dataset analysis

Take 4 data sets from existing PV data for observation, shown in Figure 3. To address the month-to-month and seasonal variations in solar radiation, we conducted an in-depth analysis of the



distribution characteristics of photovoltaic data for specific months and seasons. Solar radiation exhibits significant differences not only in its mean values but also in its variability and distribution patterns across months and seasons. For instance, summer months like July and August demonstrate higher average solar radiation and smoother temporal patterns due to stable weather conditions, while winter months like December and January are characterized by lower average radiation levels and more frequent fluctuations caused by cloud cover and shorter daylight hours. Seasonally, summer shows the highest consistency in solar radiation, while winter has the highest variability. In addition, the distribution characteristics of solar radiation, such as its skewness and kurtosis, were analyzed for each month and season to capture subtle temporal differences. Winter months generally exhibit a higher skewness due to irregular peaks in solar radiation, whereas summer months tend to have lower kurtosis, reflecting more consistent radiation patterns, as shown in Figures 3A, B.

While the clouds shade PV panels from the light, a downward peak occurs. The duration of these peaks depends on the moving speed of clouds. Additionally, as shown in Figures 3C, d, even though mean values of these 2 PV samples are the same, their representation are extremely different. Consequently, it is hard to utilize mean value and variance to evaluate the performance and diversity of PV scenario generation. The normalization is applied to the energy ratio $\lambda[n]$ and the normalized solar power $g[n]$. The energy ratio $\lambda[n]$ is normalized to follow a Gaussian distribution $G(\mu_\lambda, \sigma_\lambda^2)$, where μ_λ and σ_λ^2 are the mean and variance estimated from historical data. Similarly, $g[n]$ is transformed to follow $G(\mu_g, \sigma_g^2)$, with μ_g and σ_g^2 reflecting the peak solar power time and its variability within a day. These normalization steps standardize the feature distributions, ensuring consistent input for the evaluation models and improving the robustness of the analysis.

We first normalized the raw photovoltaic output power data in the data preprocessing phase. Specifically, each historical sample $x_o[n]$ is transformed into $x[n]$ using the following equation:

$$x[n] = \lambda[n] \cdot (g[n] + \varepsilon[n]) \quad (8)$$

Where $\lambda[n]$ represents the energy ratio, $g[n]$ is the normalized solar power, and $\varepsilon[n]$ accounts for the environmental interference, such as cloud cover and weather variations. All photovoltaic data are normalized to ensure consistency across the dataset following a Gaussian distribution. The energy ratio $\lambda[n]$ is normalized as follows:

$$\lambda[n] \sim G(\mu_\lambda, \sigma_\lambda^2) \quad (9)$$

Where μ_λ and σ_λ^2 are the mean and variance of the energy ratio, respectively. Similarly, the normalized solar power $g[n]$ follows a Gaussian distribution, given by:

$$g[n] \sim G(\mu_g, \sigma_g^2) \quad (10)$$

Where μ_g indicates the time of maximum solar power during the day, and σ_g^2 represents the time interval between sunrise and sunset. To model the environmental disturbances impacting the photovoltaic power, the interference term $\varepsilon[n]$ is modeled as an autoregressive process:

$$\varepsilon[n] = -\sum_{t=1}^p a_t \cdot \varepsilon[n-t] + u[n] \quad (11)$$

Where a_t are the AR model parameters, p is the order of the model, and $u[n]$ is the noise term, which follows a Gaussian distribution. These preprocessing steps ensure the data are appropriately normalized and standardized, providing consistent and reliable input for subsequent model training.

We establish the evaluation model according to month and season. In other words, we could assess a generated PV sample whether belongs to a specific month or season. For each x in historical data, a set of parameter values $\Theta_x: (\lambda^x, \mu_g^x, \sigma_g^2, a_1, a_2, \dots, a_p, \sigma_u^2)$ could be gained by Equation 1. Denote each generated sample as y , with a group of parameter values $\Theta_y: (\lambda^y, \mu_g^y, \sigma_g^2, a_1^y, a_2^y, \dots, a_p^y, \sigma_u^2)$ are also evaluated. While the proposed evaluation framework assumes WSS for month-specific and season-specific PV data, this assumption may not hold under significant non-stationary conditions, such as those caused by extreme weather events or rapid environmental changes. To

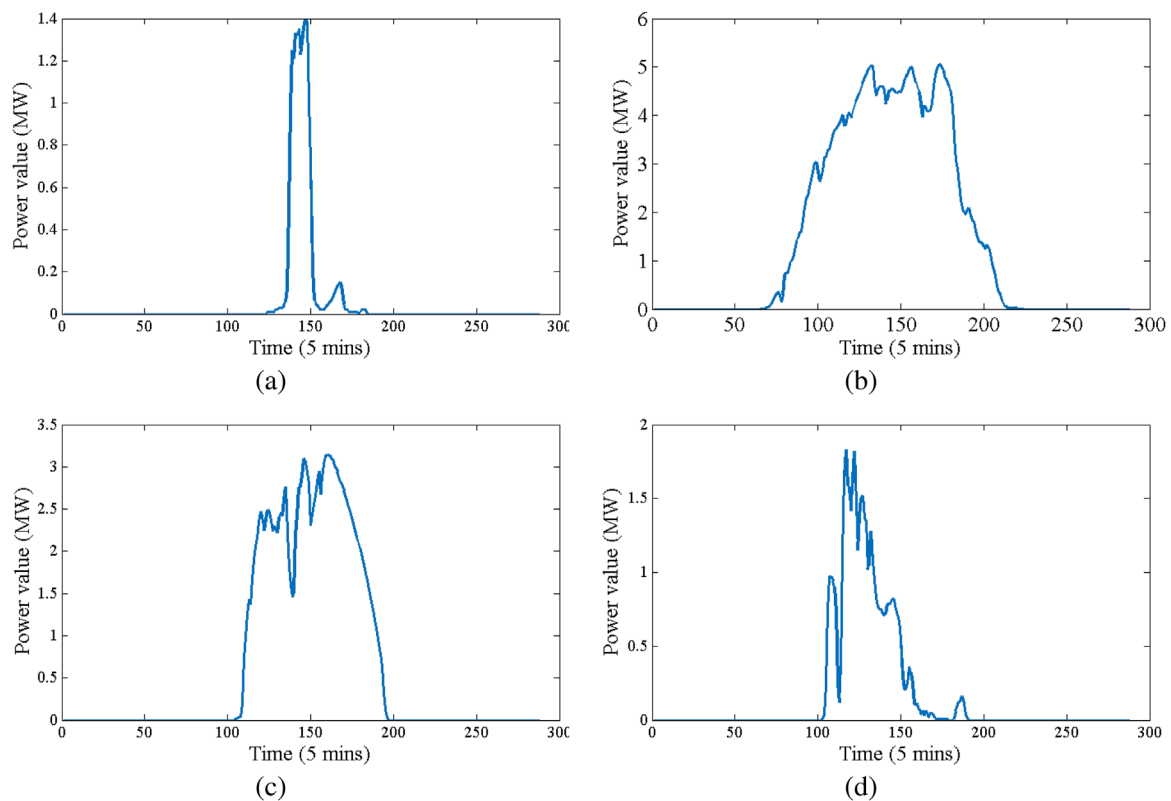


FIGURE 3

PV scenario samples, where data are sampled every 5 min (a) A daily PV output profile under relatively clear weather conditions. (b) A PV output profile exhibiting moderate fluctuations due to intermittent cloud cover. (c) A high-fluctuation scenario sharing the same average output as (d), but showing distinctly different variability patterns. (d) A contrasting scenario with the same mean power as (c), yet reflecting different environmental randomness in its generation curve.

mitigate this limitation, the evaluation process could incorporate adaptive mechanisms, such as recalibrating AR parameters for shorter time intervals to better capture transient dynamics. Furthermore, reflection coefficients and residual errors from the AR process could provide additional indicators of non-stationarity, allowing for more flexible evaluation criteria under such conditions.

The proposed model relies on several key assumptions for effective modeling and evaluation. The solar reception process is assumed to exhibit WSS within each month and season, allowing the AR model to capture temporal dependencies. Additionally, the energy ratio $\lambda[n]$ and normalized solar power $g[n]$ are assumed to follow Gaussian distributions, facilitating parameter estimation. Environmental fluctuations caused by weather changes, such as cloud movement, are modeled using a linear AR process, which enables the model to capture short-term dependencies. While these assumptions are reasonable for most PV scenarios, deviations, such as seasonal transitions or extreme weather conditions, may impact the model's generalization ability. The model is tested on PV data from multiple months and seasons to ensure robustness, reflecting a wide range of environmental conditions. This approach improves the model's generalization capacity and practical applicability.

3.2 Month-evaluation

For j th month in a year, evaluation model is determined by $\Theta_x^j: \{\tilde{\lambda}^j, \tilde{\mu}_g^j, \tilde{\sigma}_g^j, \tilde{d}_p^j, (\tilde{\sigma}_u^j)^2\}$, $j = \{1, 2, \dots, 12\}$. We arrange Θ_x^j into three parts in accordance with three indicators $\{\tilde{r}_\lambda^j[j], \tilde{r}_g^j[j], \tilde{r}_u^j[j]\}$ to evaluate a given generated PV sample y being “reliable”. If the sample is reliable, the indicators would obtain the month that this sample belongs to.

- 1) In the first part, $\tilde{r}_\lambda^j[j]$ manifests the probability that energy ratio of generated sample y follows the distribution of energy ratio of the j th month.

In term of the j th month, from Equation 2, the distribution of energy ratio is dominated by $\tilde{\mu}_\lambda^j$ and $(\tilde{\sigma}_\lambda^j)^2$ which are gained by the following formulas:

$$\tilde{\mu}_\lambda^j = \frac{1}{d} \sum_{i=1}^d \lambda^x[i] \quad (12)$$

$$(\tilde{\sigma}_\lambda^j)^2 = \frac{1}{d} \sum_{i=1}^d (\lambda^x[i] - \tilde{\mu}_\lambda^j)^2 \quad (13)$$

where d is the total days in the j th month. $\lambda^x[i]$ is the value of energy ratio in the i th day of j th month.

For each generated PV sample y , $\{r_\lambda^y[j] | j = 1, 2, \dots, 12\}$ are calculated separately as follows.

$$\hat{r}_\lambda^y[j] = \frac{1}{\sqrt{2\pi(\hat{\sigma}_\lambda^j)^2}} e^{-\frac{1}{2}\left(\frac{y - \hat{\mu}_\lambda^j}{\hat{\sigma}_\lambda^j}\right)^2} \quad (14)$$

2) In the second part, $\hat{r}_g^y[j]$ describes the difference of sunshine duration between generated samples and the j th month, which is evaluated by $\hat{\mu}_g^j$ and $(\hat{\sigma}_g^j)^2$.

$$\hat{r}_g^y[j] = \sqrt{(\mu_g - \hat{\mu}_g^j)^2 + (\sigma_g^2 - (\hat{\sigma}_g^j)^2)}, \quad (15)$$

$$\hat{\mu}_g^j = \frac{1}{d} \sum_{i=1}^d \mu_g^j[i] \quad (16)$$

$$(\hat{\sigma}_g^j)^2 = \frac{1}{d} \sum_{i=1}^d \sigma_g^2[i] \quad (17)$$

3) In the third part, $\hat{r}_u^y[j]$ represents the learning ability of generated samples for environmental randomness of j th month. An AR(p) model is designed to imitate the uncertainties. An importance of AR model is that all poles must be within unit circle to ensure system stability [23]. If $p > 2$, all $\{a_t | t = 1, 2, \dots, p\}$ are not necessary to be less than 1. Thereby, averaging a_t 's may result in poles falling outside the unit circle. Instead of averaging the values of a_t , we implement reflection coefficients $\{k_t | t = 1, 2, \dots, p\}$ to obtain the "averaged" AR parameters. k_t has a good property that it is bounded by 1. To find k_t by step-down (SD) procedure, we need to obtain $\{a_t[m] | m = 1, 2, \dots, t-1; t = p, p-1, \dots, 2\}$, where $a_t[m]$ is the m th AR parameter for model order $p = t$.

In SD procedure,

$$a_{t-1}[i] = \frac{a_t[i] - a_t[t] a_t^*[t-i]}{1 - |a_t[t]|^2} \quad (18)$$

where $*$ means transposition. Prediction error powers of each model order in the procedure is defined as:

$$v_{t-1} = \frac{v_t}{1 - |a_t[t]|^2} \quad (19)$$

SD procedure is completed while $v_t < (\hat{\sigma}_g^j)^2$. Furthermore, the procedure is initialized with $a_t[m] = a_t$ in Θ_x for $t = 1, 2, \dots, p$ and $v_p = \sigma_u^2$. After displaying SD procedure, reflection coefficients are obtained as $k_t = a_t[t]$, $t = 1, 2, \dots, p$. Averaged k_t is calculated as:

$$\hat{k}_t = \frac{1}{d} \sum_{i=1}^d k_t[i] \quad (20)$$

With \hat{k}_t , $\{\hat{a}_t | t = 1, 2, \dots, p\}$ is calculated using Levinson recursion.

$$\hat{a}_t[m] = \begin{cases} \hat{a}_{t-1}[m] + \hat{k}_t \hat{a}_{t-1}[t-m] & \hat{k} = 1, 2, \dots, t-1 \\ \hat{k}_t & m = t \end{cases} \quad (21)$$

$$\begin{aligned} (\hat{\sigma}_u^j)^2 &= \hat{r}[0] \prod_{t=1}^p (1 - \hat{k}_t^2) \\ \hat{r}[0] &= \frac{1}{n} \sum_{i=1}^n \varepsilon^2[n] \end{aligned} \quad (22)$$

$\hat{r}_u^y[j]$ is measured by Euclidean distance between AR parameters $\{a_t^y | t = 1, 2, \dots, p\}$ in a generated sample y and coefficients $\{\hat{a}_t^j | t = 1, 2, \dots, p\}$ in the j th month.

$$\hat{r}_u^y[j] = \sqrt{\sum_{t=1}^p (\hat{a}_t^j - a_t^y)^2 + ((\hat{\sigma}_u^j)^2 - (\sigma_u^y)^2)} \quad (23)$$

For each generated sample y , a set of $\{\hat{r}_\lambda^y[j], \hat{r}_g^y[j], \hat{r}_u^y[j]\}$ is obtained. Based on maximum likelihood theory, we search the largest values of these 3 coefficients with highest probability for finding the month that y belongs to.

$$\begin{cases} j_1 = \arg \max_j \hat{r}_\lambda^y[j], & j = 1, 2, \dots, 12 \\ j_2 = \arg \min_j \hat{r}_g^y[j], & j = 1, 2, \dots, 12 \\ j_3 = \arg \min_j \hat{r}_u^y[j], & j = 1, 2, \dots, 12 \end{cases} \quad (24)$$

To place the assessment, state null hypothesis and alternate hypothesis, respectively:

H_0 : y is a reliable sample for PV scenario.

H_1 : y is not a reliable sample for PV scenario.

Because two neighboring months have some similar solar receptions, we set a bias δ to adjust the evaluation. Considering the comparability between any two neighboring months, while $|j_1 - j_2| \leq \delta$, $|j_1 - j_3| \leq \delta$, and $|j_2 - j_3| \leq \delta$ are all satisfied, accept H_0 . Otherwise, accept H_1 . To further quantify the confidence in the evaluation results, we calculated the 95% confidence intervals (CIs) for the likelihood ratios $\hat{r}_\lambda^y[j]$, $\hat{r}_g^y[j]$, $\hat{r}_u^y[j]$ using the following formula:

$$CI = \hat{r} \pm z \cdot \frac{\sigma}{\sqrt{n}} \quad (25)$$

where \hat{r} is the maximum likelihood estimate, σ is the standard deviation of the likelihood ratios across the samples, n is the number of observations, and z is the critical value for a 95% confidence level. We use corresponding to a 95% confidence interval. This 5% significance level is a widely accepted convention in statistical inference, as it offers a practical balance between Type I and Type II errors in hypothesis testing [24–26]. Many studies in related fields (e.g., power systems, reliability analysis) similarly adopt the 95% CI when evaluating model performance or uncertainty quantification [27–29], ensuring that our approach remains consistent with standard practice. Additionally, statistical significance tests were conducted to validate the reliability of the classification results. We employed a one-sample t-test to assess whether the mean likelihood ratio for each month or season significantly differed from a pre-defined threshold μ_0 , representing unreliable PV scenarios. The test statistic is given by:

$$t = \frac{\bar{r} - \mu_0}{\frac{\sigma}{\sqrt{n}}} \quad (26)$$

where \bar{r} is the mean likelihood ratio, and σ is its standard deviation. The p-value corresponding to the test statistic determines whether to reject H_0 . A significance level of 0.05 was used as the cutoff for statistical significance, which is a widely adopted threshold in statistical hypothesis testing across various fields, including the social sciences, engineering, and environmental

studies. The choice of a 0.05 significance level, corresponding to a 95% confidence interval. This threshold means that there is a 5% chance of incorrectly rejecting the null hypothesis when it is actually true, which provides a practical balance between precision and practical decision-making. In our study, using this significance level allowed us to confidently assess whether the observed differences in likelihood ratios were statistically significant, thus providing reliable classification of PV scenarios. This choice of significance level aligns with commonly accepted practices in similar research studies [30–32]. By adopting this 0.05 significance level, we ensured that our results met established statistical standards, offering a reliable and robust evaluation of the PV scenario classification process.

Normally, $\delta = 1$. The parameter δ is used as a threshold to account for the natural similarity between neighboring months, ensuring that minor deviations do not lead to misclassification of PV scenarios. Based on empirical analysis, δ is set to 1, which balances classification accuracy and robustness. A larger δ increases the tolerance for classification, potentially leading to misclassifications, while a smaller δ may cause overly strict categorization, especially in transitional months with similar solar patterns (e.g., March and April or September and October). Algorithm 1 describes the details of Month-Evaluation model.

3.3 Season-Evaluation

This model is similar to Month-Evaluation model, in which s replaces d and represents the total days in each season. Season-Evaluation model is assessed by $\{\tilde{r}_\lambda^y[j], \tilde{r}_g^y[j], \tilde{r}_u^y[j]\}$. $h = \{1, 2, 3, 4\}$ represents spring, summer, autumn, and winter, respectively. Season-Evaluation model is used to identify whether a generated sample y as a reliable PV scenario belongs to a specific season, in which the three-step procedure is also designed with a season of $\Theta_x^h: \{\tilde{\lambda}^h, \tilde{\mu}_g^h, \tilde{\sigma}_g^h, \tilde{a}_t^h, (\tilde{\sigma}_u^h)^2\}$.

- 1) Firstly, $\tilde{r}_\lambda^y[j]$ is obtained through evaluation of distribution of energy ration according to Gaussian distribution which is determined by its mean value and variance.

$$\tilde{\mu}_\lambda^h = \frac{1}{s} \sum_{q=1}^s \lambda^x[q] \quad (27)$$

$$(\tilde{\sigma}_\lambda^h)^2 = \frac{1}{s} \sum_{q=1}^s (\lambda^x[q] - \tilde{\mu}_\lambda^h)^2 \quad (28)$$

where $\lambda^x[q]$ is q th day in h th season that contains s days.

Given a generated PV sample y , $\{\tilde{r}_\lambda^y[h] | h = 1, 2, 3, 4\}$ are evaluated separately using the following formula.

$$\tilde{r}_\lambda^y[h] = \frac{1}{\sqrt{2\pi(\tilde{\sigma}_\lambda^h)^2}} e^{-\frac{1}{2} \left(\frac{\lambda^y - \tilde{\mu}_\lambda^h}{\tilde{\sigma}_\lambda^h} \right)^2} \quad (29)$$

- 2) Secondly, owing to the particularity of sunshine duration in each season, the pattern of sunrise and sunset of seasons is depicted with $\tilde{\mu}_g^h$ and $(\tilde{\sigma}_g^h)^2$.

Input: Θ_x , Θ_y , generated PV sample y , d , δ and p

Output: H_0 and H_1 .

Initialization: $\tilde{r}_\lambda^y[j] = 0$, $\tilde{r}_g^y[j] = 0$, $\tilde{r}_u^y[j] = 0$.

for $j = 1$ to 12 **do**

for $i = 1$ to d **do**

1: Sum $\lambda^x[i]$ using all λ^x in a month;

end for

2: Calculate the mean value $\tilde{\mu}_\lambda^j$ and $(\tilde{\sigma}_\lambda^j)^2$ in

Equations 4, 5;

3: Compute $\tilde{r}_\lambda^y[j]$ in Equation 6.

end for

for $j = 1$ to 12 **do**

for $i = 1$ to d **do**

4: Sum μ_g^j using all μ_g in a month;

5: Sum $(\sigma_g^j)^2$ using all $(\sigma_g)^2$ in a month;

end for

6: Calculate the mean value $\tilde{\mu}_g^j$ and $(\tilde{\sigma}_g^j)^2$ in

Equations 7, 8;

7: Compute $\tilde{r}_g^y[j]$ in Equation 7.

end for

for $j = 1$ to 12 **do**

for $i = 1$ to d **do**

for $m = 1$ to p **do**

8: Evaluate $a_t[m]$ for each AR order;

end for

9: Obtain k_t using SD procedure with $a_t[m]$;

end for

10: Average all k_t in a month to get \hat{k}_t ;

11: Gain averaged \hat{a}_t using Levinson recursion

with \hat{k}_t ;

12: Calculate $(\tilde{\sigma}_u^j)^2$ with \hat{k}_t in Equation 14;

13: Compute $\tilde{r}_u^y[j]$ in Equation 15.

end for

14: Sort $\tilde{r}_\lambda^y[j]$, $\tilde{r}_g^y[j]$, and $\tilde{r}_u^y[j]$ to find j_1 , j_2 , and j_3 .

if $|j_1 - j_2| \leq \delta$ && $|j_1 - j_3| \leq \delta$ && $|j_2 - j_3| \leq \delta$. **then**

15: Accept H_0 .

else

16: Accept H_1 .

end if

Algorithm 1. Month-Evaluation Procedure.

$$\tilde{\mu}_g^h = \frac{1}{s} \sum_{q=1}^s \mu_g[q] \quad (30)$$

$$(\tilde{\sigma}_g^h)^2 = \frac{1}{s} \sum_{q=1}^s \sigma_g^2[q] \quad (31)$$

$\tilde{r}_g^y[h]$ is to evaluate the given sample y based on the discovery of the Gaussian distribution of solar reception process.

$$\tilde{r}_g^y[h] = \sqrt{(\mu_g - \tilde{\mu}_g^h)^2 + (\sigma_g^2 - (\tilde{\sigma}_g^h)^2)}, \quad (32)$$

- 3) Thirdly, reflection coefficients \tilde{k}_t are also used to gain averaged AR parameters \tilde{a}_t^h and $(\tilde{\sigma}_u^h)^2$ to guarantee system stability.

$$\tilde{k}_t = \frac{1}{s} \sum_{q=1}^s k_t[q] \quad (33)$$

$k_t[q]$ is obtained with a_t using SD procedure in Equation 10. Then, $\{\tilde{a}_t^h | t = 1, 2, \dots, p; h = 1, 2, 3, 4\}$ and $(\tilde{\sigma}_u^h)^2$ are evaluated with $\{\tilde{k}_t | t = 1, 2, \dots, p\}$ using Levinson recursion. $\tilde{r}_a^m[h]$ is designed as:

$$\tilde{r}_a^m[h] = \sqrt{\sum_{t=1}^p (\tilde{a}_t^h - a_t^y)^2 + ((\tilde{\sigma}_u^h)^2 - (\sigma_u^y)^2)} \quad (34)$$

With a series of $\{\tilde{r}_\lambda^y[h], \tilde{r}_g^y[h], \tilde{r}_u^y[h]\}$, the appropriate season that y belongs to is discovered with the highest probabilities using maximum likelihood theory.

$$\begin{cases} h_1 = \arg \max_h \tilde{r}_\lambda^y[h], & h = 1, 2, 3, 4 \\ h_2 = \arg \min_h \tilde{r}_g^y[h], & h = 1, 2, 3, 4 \\ h_3 = \arg \min_h \tilde{r}_u^y[h], & h = 1, 2, 3, 4 \end{cases} \quad (35)$$

To place the assessment, state null hypothesis and alternate hypothesis, respectively:

$$\begin{aligned} H_0: & y \text{ is a reliable PV sample in } h^* \text{ season} \\ H_1: & y \text{ is not a reliable sample for PV scenario.} \end{aligned}$$

Only while $h_1 = h_2 = h_3$, accept H_0 , where $h^* = h_1$. Otherwise, accept H_1 . Algorithm 2 describes the details of Season-Evaluation.

3.4 Correlation between month and seasonal evaluation

The monthly and seasonal evaluation methods differ in their temporal granularity and parameter estimation processes. The monthly evaluation classifies PV scenarios into 12 specific months, using month-specific parameters to capture fine-grained temporal differences in PV power generation. In contrast, the seasonal evaluation classifies scenarios into four broader seasonal categories, where seasonal parameters are aggregated from monthly data. This approach increases robustness to short-term fluctuations but reduces temporal resolution. Together, these methods provide a complementary framework, with the monthly evaluation offering higher precision and the seasonal evaluation providing more excellent stability in the presence of temporal variability.

4 Experimental results

4.1 Software settings

The experiments presented in this study were implemented using Matlab 2015, a powerful tool for numerical computations, statistical analysis, and time series modeling. Matlab's robust

environment allowed us to efficiently handle the large datasets required for evaluating photovoltaic (PV) scenarios and implementing the autoregressive (AR) model. For the AR and ARMA modeling, we utilized Matlab's built-in toolboxes, which provide comprehensive functions for time series analysis and statistical testing. Additionally, other techniques such as Generative Adversarial Networks (GAN) and Conditional GAN (CGAN) were implemented using TensorFlow to handle the deep learning-based generation of PV scenarios.

4.2 Experiment settings

With Solar Integration datasets [33], we choose solar data from both 32 solar power plants in the State of Washington to train all PV scenario generation techniques in the simulation. The Solar Integration dataset is selected for its wide acceptance and representativeness in PV scenario generation research. This dataset provides real-world operational data from 32 PV power plants, capturing the natural variability of solar power influenced by weather conditions, geographical differences, and seasonal changes. Unlike synthetic datasets, the Solar Integration dataset reflects the stochastic nature of PV power generation, allowing for a more comprehensive evaluation of the model's generalization ability. Its diverse feature distribution and realistic noise levels ensure that the model is tested under practical conditions, enhancing the credibility of the experimental results. Potential biases in the dataset are primarily introduced by weather-related randomness and seasonal shifts in solar irradiance. These factors may result in an imbalanced sample distribution, especially during months with more frequent weather disturbances. To mitigate this, the proposed model incorporates month-specific and season-specific parameter estimation, allowing it to account for these natural variations. The model aims to achieve more robust evaluation performance by explicitly modeling temporal and seasonal effects. Including these statistical characteristics ensures that the evaluation method's assumptions are transparent and justified.

All experiments of evaluation are operated in Matlab R2014b software with 8 GB memory. In order to generate PV scenario data with different techniques, we implement AR and autoregressive and moving average (ARMA) toolbox in Matlab and display Gaussian copula Matlab codes. Additionally, GAN, CGAN, and CGAN-filtering models are established in tensorflow with a single Nvidia TITAN Xp GPU to obtain new PV scenario samples. Furthermore, three popular data generation techniques are also used to produce PV scenarios: random oversample (ROS), synthetic minority over-sampling technique (SMOTE), and adaptive synthetic sampling (ADASYN). The parameter settings of these 9 techniques are as follows:

1. AR: $p = 6$, which is the order of AR model.
2. ARMA: $p = 6$, and $q = 5$. q is the order of MA model.
3. Gaussian copula: solar energy has a wide range between different seasons. Thus, we normalize the existing PV scenarios as pretreatment.
4. GAN: the generator is initialized as a 4-layer neural network in which the sizes of hidden layers are $\{128, 1024, 256, 288\}$. As the discriminator, 4 hidden layers whose sizes are $\{288, 256, 1024, 128\}$ and a softmax layer are involved.

Input: θ_x , θ_y , generated PV sample y , s , and p
Output: H_0 and H_1
Initialization: $\tilde{r}_\lambda^y[h] = 0$, $\tilde{r}_g^y[h] = 0$, $\tilde{r}_u^y[h] = 0$.
if y is not a reliable sample in Month-Evaluation model. **then**
 1: Accept H_1 .
else
 for $h = 1$ to 4 **do**
 for $q = 1$ to s **do**
 2: Sum $\lambda^x[q]$ using all λ^x in a month;
 end for
 3: Calculate the mean value $\bar{\mu}_\lambda^h$ and $(\bar{\sigma}_\lambda^h)^2$ in Equations 17, 18;
 4: Compute $\tilde{r}_\lambda^y[h]$ in Equation 19.
 end for
 for $h = 1$ to 4 **do**
 for $q = 1$ to s **do**
 5: Sum μ_g^h using all μ_g in a month;
 6: Sum $(\sigma_g^h)^2$ using all $(\sigma_g)^2$ in a month;
 end for
 7: Calculate the mean value $\bar{\mu}_g^h$ and $(\bar{\sigma}_g^h)^2$ in Equations 20 and 21;
 8: Compute $\tilde{r}_g^y[h]$ in Equation 22.
 end for
 for $j = 1$ to 4 **do**
 for $i = 1$ to s **do**
 for $m = 1$ to p **do**
 9: Evaluate $a_t[m]$ for each AR order;
 end for
 10: Obtain k_t using SD procedure with $a_t[m]$;
 end for
 11: Average all k_t in a month to get \bar{k}_t ;
 12: Gain averaged \bar{a}_t using Levinson recursion with \bar{k}_t ;
 13: Calculate $(\bar{\sigma}_u^h)^2$ with \bar{k}_t ;
 14: Compute $\tilde{r}_u^y[h]$ in Equation 24.
 end for
 15: Sort $\tilde{r}_\lambda^y[h]$, $\tilde{r}_g^y[h]$, and $\tilde{r}_u^y[h]$ to find h_1 , h_2 , and h_3 .
 if $h_1 = h_2 = h_3$. **then**
 16: $h^* = h_1$.
 17: Accept H_0 .
 else
 18: Accept H_1 .
 end if
end if

Algorithm 2. Season-Evaluation Procedure.

- CGAN: in addition to GAN parameters, mean values of powers are calculated for each sample, and the results are classified into 5 categories: $\mu(X_i) < 0.2$ (class 1), $0.2 < \mu(X_i) < 0.5$ (class 2), $0.5 < \mu(X_i) < 1$ (class 3), $1 < \mu(X_i) < 2$ (class 4), and $\mu(X_i) > 2$ (class 5).

- CGAN-filtering: in addition to CGAN settings, there are 2 parameters in the filtering. The orders of zeros and poles are initialized as 10 and 5, respectively.
- ROS: random state = 42, which represents the random number generator.
- SMOTE and ADASYN: $k = 3$, and $m = 5$. k is the number of nearest neighbours that construct synthetic samples, and m is the number of nearest neighbours that determine if a minority sample is in danger.

Normally, solar energy has a wide range between different seasons. Thus, we normalize the existing PV scenarios as pretreatment. 6 generated PV scenarios from these 8 generation techniques chosen by random are shown in Figure 4. Obviously, AR and ARMA methods generate PV scenarios with diversity, while Gaussian copula and GAN family specialize in learning the environment randomness (i.e., the shape of PV scenario), and SMOTE and ADASYN focus on the energy of PV scenarios.

4.3 Comparisons with evaluation metrics

4.3.1 Comparisons on generated PV samples

We compare the proposed model (termed as \mathbb{A}) with 2 popular evaluation metrics, i.e., mean value (denoted as \mathbb{M}) and variance (denoted as \mathbb{V}). 6 samples produced by 8 techniques (illustrated in Figure 3) are used to obtain evaluation results both for month and season estimation. We compare the proposed model (termed as \mathbb{A}) with 2 popular evaluation metrics, i.e., mean value (denoted as \mathbb{M}) and variance (denoted as \mathbb{V}). 6 samples produced by 8 techniques (illustrated in Figure 3) are used to obtain evaluation results both for month and season estimation. Moreover, the inability of mean value and variance-based metrics to capture these temporal and statistical characteristics often leads to higher evaluation errors, particularly in months or seasons with extreme variations in solar radiation. For instance, the AR-generated PV scenarios for January showed unrealistic consistency in solar energy levels during midday, which is against natural solar radiation patterns. While the mean value and variance metrics failed to identify such anomalies, the proposed AR-based evaluation model effectively detected these inconsistencies by analyzing the temporal correlations and stochastic fluctuations in the data.

In contrast, GAN-generated samples displayed distinct characteristics that influenced their evaluation. First, GANs excel in capturing the overall shape and variability of PV scenarios, as they learn the underlying data distribution from historical datasets. This allows GAN-generated samples to exhibit temporal dependencies that closely mimic real-world PV scenarios, particularly in months with stable solar radiation, such as July and August. However, the stochastic nature of GANs introduces noise into the generated samples, which can manifest as small, high-frequency fluctuations that deviate from natural solar radiation patterns. These noise-induced deviations are subtle and often escape detection by mean value and variance-based metrics but are effectively captured by the proposed AR-based evaluation model due to its sensitivity to temporal correlations. Additionally, the evaluation results revealed that the performance of GAN-generated samples varied with the

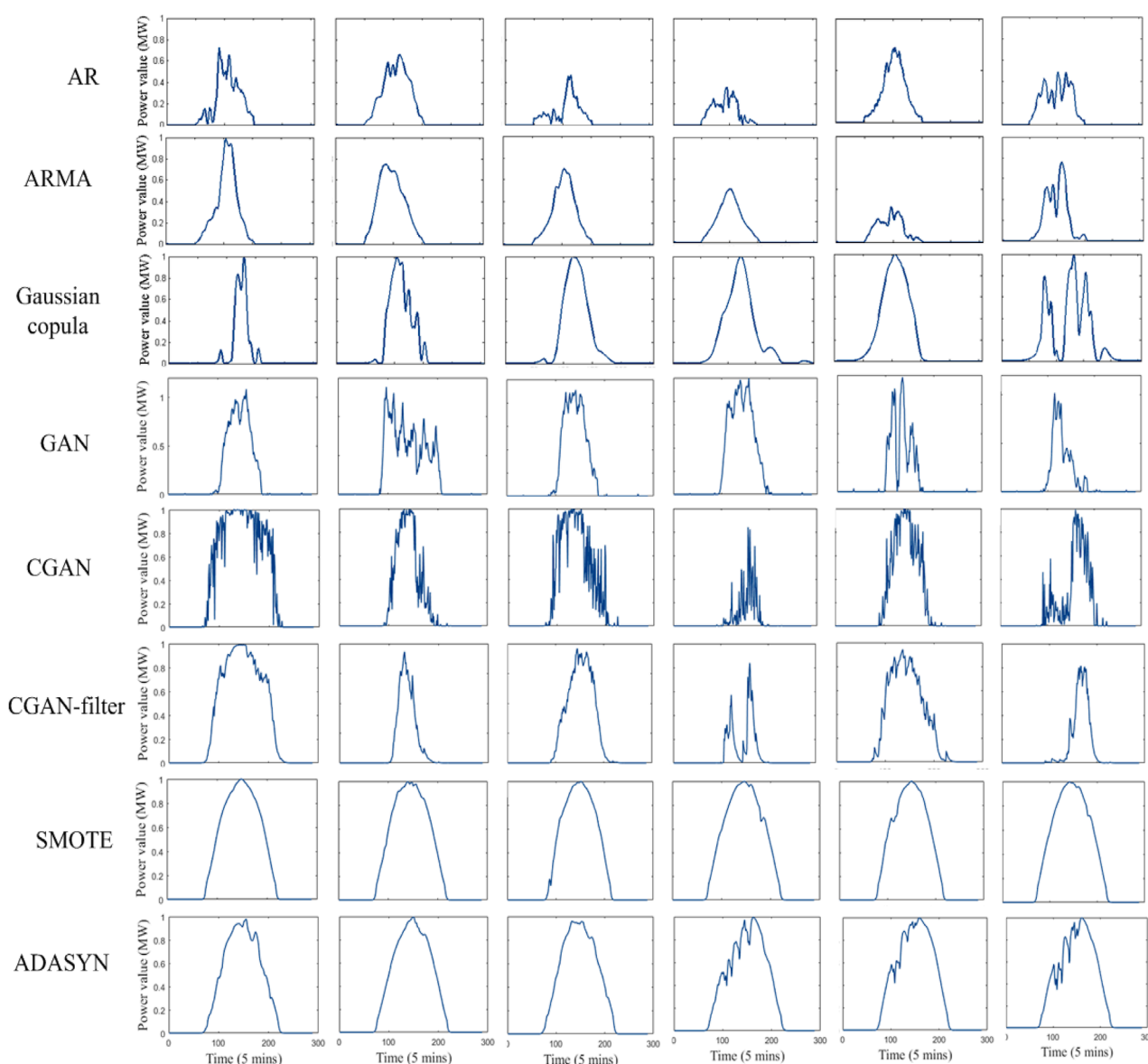


FIGURE 4
Generated PV scenarios by different approaches.

complexity of the temporal patterns in the target data. For example, in winter months like January, where solar radiation patterns are highly irregular due to frequent cloud cover, GAN-generated samples tended to exhibit over-smoothed temporal trends, failing to replicate the abrupt changes observed in real data. The proposed model successfully identified these limitations by analyzing the AR parameters and residuals, which highlighted the discrepancies in the stochastic dynamics between the generated samples and the actual data.

To observe the performance of evaluation methods, we specially choose some non-PV samples. Experimental results are shown in Tables 1, 2. ' i/j ' means that this PV scenario actually belongs to j th month or season which is evaluated as a sample in i th month or season. '×' indicates that this sample cannot be deemed as a reliable PV scenario in line with existing solar reception pattern. According

to the similarity between the generated samples and existing PV scenarios, we denote the month or season that generated samples belong to (i.e., the value of j).

In regard to month attribution of a PV sample set, the proposed model outperforms mean value and variance estimations. Table 1 shows mean value and variance have higher evaluation errors of the months that produced PV samples belongs to. In particular, while a generated PV sample does not follow solar reception principles, mean and variance have difficulties in observing it. Take the 5th sample from AR model for instance, it is impossible that solar energies remain the highest value between around 12 p.m. and 16 p.m., which is against nature law leading to an unreliable PV sample. This is because AR parameters in an AR model are sensitive and could result in unstable state. The proposed model is capable of estimating the unreliability

TABLE 1 Evaluation performance of generated PV scenario in specific months.

	AR	ARMA	Gaussian copula	GAN	CGAN	CGAN-filter	SMOTE	ADASYN
M	1/1	12/×	1/1	10/3	7/×	7/8	7/8	7/8
	12/×	2/×	2/×	7/8	2/×	1/1	7/8	7/8
	1/1	12/1	10/3	2/3	5/3	1/×	7/8	7/8
	1/1	1/1	3/3	3/3	1/×	2/3	7/8	7/8
	7/×	7/×	3/×	1/×	3/×	1/×	7/8	7/8
	1/1	1/1	3/4	11/×	12/4	5/4	7/8	7/8
V	1/1	12/×	1/1	6/3	8/×	8/8	8/8	8/8
	1/×	11/×	3/×	8/8	4/×	1/1	8/8	8/8
	1/1	1/1	4/3	3/3	8/×	1/×	8/8	8/8
	1/1	1/1	3/3	4/3	1/×	12/3	8/8	8/8
	8/×	8/×	5/×	1/×	9/×	1/×	8/8	8/8
	1/1	1/1	3/4	1/×	12/×	6/4	8/8	8/8
A	1/1	×/×	1/1	3/3	×/×	8/8	8/8	8/8
	×/×	×/×	×/×	8/8	×/×	1/1	8/8	8/8
	1/1	1/1	3/3	3/3	×/×	×/×	8/8	8/8
	1/1	1/1	3/3	3/3	×/×	3/3	8/8	8/8
	×/×	×/×	×/×	×/×	×/×	×/×	8/8	8/8
	1/1	1/1	4/4	×/×	×/×	4/4	8/8	8/8

of PV samples. On the opposite, mean and variance evaluation cannot identify the unreliability. Typically, PV scenarios among neighbouring months have analogous attributes, e.g., daylight hours that affects the width of PV scenarios. For example, the generated PV samples from SMOTE belong to August, yet they are deemed as scenarios in July by mean evaluation. Furthermore, the months in spring and autumn have similar solar conditions, e.g., solar reception amount in a day. Therefore, mean and variance are unable to tell the difference, e.g., October and March. At last, when learned representations of environmental changes are not precise, the fluctuations in PV samples are abnormal (e.g., PV samples generated by CGAN). In that case, these samples are unreliable, but mean value and variance cannot detect them.

For season attribution of a PV sample, the proposed model outperforms traditional mean value and variance estimations, as shown in Table 2). The improved performance for season classification is attributed to the autoregressive (AR) model used in the proposed approach, which captures temporal dependencies and seasonal variations more effectively than simple statistical summaries. The AR model was implemented using Matlab 2015 due to its robust capabilities for time-series modeling and statistical

analysis. In particular, we utilized Matlab's Econometrics Toolbox, which includes functions for autoregressive and moving average modeling, to analyze the PV data and assess the seasonal attributes of the generated scenarios. These basic statistical techniques were applied to evaluate the solar data and assess the seasonality of PV output based on average values and variances, which do not account for the underlying temporal correlations. In contrast to month evaluation, evaluation techniques display well on season attribute. This is because that the differences in solar movement and environment changes among seasons are more significant than months. To assess the seasonal performance, the confidence intervals for likelihood ratios were calculated for each season, using the one-sample t-test approach implemented in Matlab. The p-values corresponding to these tests were computed to determine the statistical significance of the seasonal differences in the PV samples.

4.3.2 Comparisons of evaluation performance

In order to verify the effectiveness of the proposed model, we implement mean value, variance, and the proposed model displaying on 100 samples generated by these 8 techniques, separately. Evaluation errors of these 3 measurements are shown

TABLE 2 Evaluation performance of generated PV scenario in specific seasons.

	AR	ARMA	Gaussian copula	GAN	CGAN	CGAN-filter	SMOTE	ADASYN
M	4/4	4/×	4/4	1/1	2/×	2/2	2/2	2/2
	4/×	1/×	1/×	2/2	1/×	4/4	2/2	2/2
	4/4	4/4	1/1	1/1	3/×	4/×	2/2	2/2
	4/4	4/4	1/1	1/1	4/×	1/1	2/2	2/2
	2/×	2/×	1/×	4/×	1/×	4/×	2/2	2/2
	4/4	4/4	1/1	4/×	4/×	2/1	2/2	2/2
V	4/4	4/×	4/4	3/1	2/×	2/2	2/2	2/2
	4/×	4/×	1/×	2/2	1/×	4/4	2/2	2/2
	4/4	4/4	1/1	1/1	2/×	4/×	2/2	2/2
	4/4	4/4	1/1	3/1	4/×	4/1	2/2	2/2
	2/×	2/×	3/×	4/×	2/×	4/×	2/2	2/2
	4/4	4/4	1/1	4/×	4/×	3/1	2/2	2/2
A	4/4	×/×	4/4	1/1	×/×	2/2	2/2	2/2
	×/×	×/×	×/×	2/2	×/×	4/4	2/2	2/2
	4/4	4/4	1/1	1/1	×/×	×/×	2/2	2/2
	4/4	4/4	1/1	1/1	×/×	1/1	2/2	2/2
	×/×	×/×	×/×	×/×	×/×	×/×	2/2	2/2
	4/4	4/4	1/1	×/×	×/×	1/1	2/2	2/2

in Figure 5, 6, in which the proposed model outperforms other 2 measurements. This validates that mean value and variance are unable to identify unreliable generated samples, leading to high evaluation errors with AR and ARMA approaches. By contrast, the proposed model exhibits strong ability of estimating the reliability of generated samples, resulting in low evaluation errors. Moreover, with increased diversity of generation by Gaussian copula, GAN, and CGAN-filter, mean value and variance measurements are easily trapped into misidentification between neighboring months or seasons. Furthermore, because the generated samples by SMOTE and ADASYN distributes in a narrow region, these 3 evaluation approaches also gain good performances. Compared with month evaluation, evaluation results for seasons achieve better performance.

4.3.3 Analysis of computational efficiency

To evaluate the computational efficiency of the proposed model, we analyze its time complexity, resource requirements, and runtime performance for datasets of varying sizes. The computational complexity of the key components is assessed to ensure the model's feasibility for large-scale and high-frequency PV data. Our experiments showed that the confidence intervals

for likelihood ratios across months and seasons were within a range of $\pm 2\%$ of the maximum likelihood estimates, indicating high reliability of the evaluation results. Furthermore, p-values from the t-tests confirmed the statistical significance of the classification results, with all p-values below the 0.05 threshold for reliable PV scenarios. These quantitative measures not only validate the evaluation results but also demonstrate the robustness of the proposed model under varying conditions. The AR model fitting process, with an order of p , has a computational complexity of $O(N \cdot p^2)$, where N denotes the number of samples. The month-based and season-based evaluation models, which involve parameter estimation and hypothesis testing, have a complexity of $O(N \cdot p)$. To provide practical insight into runtime performance, we conduct experiments on datasets with sizes ranging from 10,000 to 100,000 samples. The results show that the total runtime increases approximately linearly with the dataset size. This demonstrates that the proposed model maintains computational efficiency even when applied to large datasets, making it suitable for high-frequency PV data applications.

The superior performance of the proposed method, mainly when applied to GAN-generated PV scenarios, can be attributed to its ability to capture the temporal dependencies and stochastic

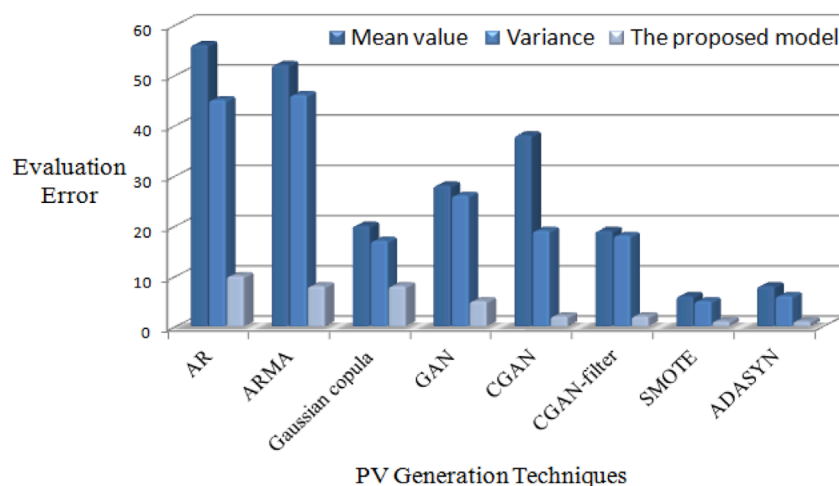


FIGURE 5
Comparisons of evaluation errors by 3 approaches for month.

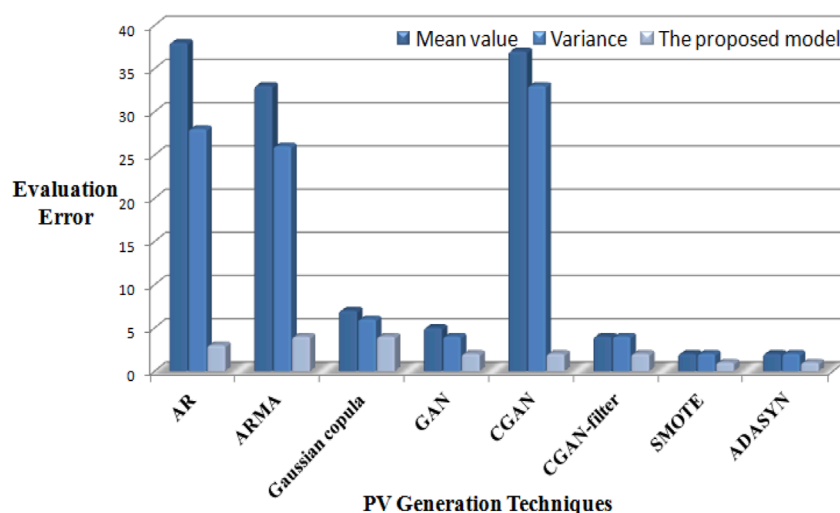


FIGURE 6
Comparisons of evaluation errors by 3 approaches for season.

nature of PV power fluctuations. GAN-generated samples often exhibit diversity in overall shape and randomness. However, they may need to accurately reflect the fine-grained temporal structure and dynamic changes caused by weather and cloud movements. Unlike traditional mean and variance-based evaluation methods, which only consider statistical summaries, the proposed AR-based model analyzes the temporal correlations within PV scenarios. By incorporating AR parameters and assessing the consistency of month- and season-specific temporal patterns, the proposed method can identify subtle deviations in the dynamic characteristics of GAN-generated scenarios. This capability allows for more precise detection of abnormal or unreliable samples that traditional metrics might overlook. The enhanced capacity to track and evaluate temporal patterns is a key reason for the superior performance of the proposed method when evaluating scenarios generated by GANs.

5 Conclusion

Solar photovoltaic had caught plenty of attentions due to its little pollution, and PV scenario generation was going to be an effective way to facilitate integrating solar energy into traditional energy systems. In order to effectively evaluate the performance of PV scenario generation, we proposed an evaluation model based on AR theory. After analyzing existing PV samples, we found out the shape of PV scenarios was an important representation of environmental randomness. In the simulation, we produced PV samples with 8 popular generation approaches. Compared with mean value and variance measurements, experiments showed the proposed model achieved better performance, especially in a unreliable PV scenario. Moreover, mean value and variance estimation confused with months that have similar solar movement

and environmental changes. With 100 generated PV scenarios, we simulated the evaluation among the proposed model and compared measurements. Simulations showed that the proposed model obtained better evaluation results than mean value and variance estimations.

In addition to its theoretical evaluation performance, the proposed model offered practical value for real-world applications, particularly power dispatching. Accurate classification of PV scenarios into month- and season-specific categories enabled system operators to predict solar power availability more effectively. By capturing the temporal patterns of solar power generation, the model supported power dispatching decisions, allowing grid operators to adjust dispatch schedules in response to seasonal and weather-induced fluctuations. The month-specific evaluation provided higher temporal resolution, enabling short-term dispatch adjustments, while the seasonal evaluation offered long-term insights for seasonal dispatch planning. This dual-level evaluation approach enhanced the robustness and flexibility of power-dispatching strategies.

Data availability statement

The raw data supporting the conclusions of this article will be made available by the authors, without undue reservation.

Author contributions

SR: Conceptualization, Writing—original draft. TY: Investigation, Software, Writing—review and editing. JL: Methodology, Supervision, Writing—review and editing. GW: Formal Analysis, Project administration, Validation, Writing—review and editing. KM: Resources, Visualization, Writing—review and editing. BL: Data curation, Funding acquisition, Writing—review and editing.

References

1. Esum T, Kimball JW, Krein PT, Chapman PL, Midya P. Dynamic maximum power point tracking of photovoltaic arrays using ripple correlation control. *IEEE Trans Power Electronics* (2008) 21:1282–91. doi:10.1109/tpe.2006.880242
2. Wai RJ, Wang WH, Lin CY. High-performance stand-alone photovoltaic generation system. *IEEE Trans Ind Electronics* (2008) 55:240–50. doi:10.1109/tie.2007.896049
3. Basore PA, Cole WJ. Comparing supply and demand models for future photovoltaic power generation in the usa. *Prog Photovoltaics Res Appl* (2018) 26:414–8. doi:10.1002/pp.2997
4. De Brito MAG, Galotto L, Sampaio LP, e Melo Gd. A, Canesin CA. Evaluation of the main mppt techniques for photovoltaic applications. *IEEE Trans Ind Electron* (2013) 60:1156–67. doi:10.1109/tie.2012.2198036
5. Renaudineau H, Donatantonio F, Fontchastagner J, Petrone G, Spagnuolo G, Martin J-P, et al. A pso-based global mppt technique for distributed pv power generation. *IEEE Trans Ind Electronics* (2015) 62:1047–58. doi:10.1109/tie.2014.2336600
6. Golestaneh F, Pinson P, Gooi HB. Very short-term nonparametric probabilistic forecasting of renewable energy generation— with application to solar energy. *IEEE Trans Power Syst* (2016) 31:3850–63. doi:10.1109/tpwrs.2015.2502423
7. De la Fuente DV, Rodríguez CLT, Garcera G, Figueres E, González RO. Photovoltaic power system with battery backup with grid-connection and islanded operation capabilities. *IEEE Trans Ind Electron* (2013) 60:1571–81. doi:10.1109/TIE.2012.2196011
8. Estébanez EJ, Moreno VM, Pigazo A, Liserre M, Dell'Aquila A. Performance evaluation of active islanding-detection algorithms in distributed-generation photovoltaic systems: two inverters case. *IEEE Trans Ind Electronics* (2011) 58:1185–93. doi:10.1109/TIE.2010.2044132
9. Yang T, Huang Q, Cai F, Li J, Jiang L, Xia Y. Vital characteristics cellular neural network (vcnn) for melanoma lesion segmentation: a biologically inspired deep learning approach. *J Imaging Inform Med* (2024) 1–18. doi:10.1007/s10278-024-01257-w
10. An Y, Zhang K, Chai Y, Zhu Z, Liu Q. Gaussian mixture variational-based transformer domain adaptation fault diagnosis method and its application in bearing fault diagnosis. *IEEE Trans Ind Inform* (2024) 20:615–25. doi:10.1109/TII.2023.3268750
11. Parsons H, Cochran S, Batra. Variability of power from large-scale solar photovoltaic scenarios in the state of Gujarat: preprint. In: *To be presented at the renewable energy world conference and expo-India, 5-7 may 2014*. New Delhi, India (2014).
12. Osório GJ, Lujano-Rojas JM, Matias JCO, Catalão JPS. A new scenario generation-based method to solve the unit commitment problem with high penetration of renewable energies. *Int J Electr Power Energy Syst* (2015) 64:1063–72. doi:10.1016/j.ijepes.2014.09.010
13. Ekstrom J, Koivisto M, Mellin I, Millar J, Lehtonen M. A statistical model for hourly large-scale wind and photovoltaic generation in new locations. *IEEE Trans Sustainable Energy* (2017) PP:1383–93. doi:10.1109/tste.2017.2682338
14. Luis LM (2018). Phdthesis: framework for scenario generation and reduction in photovoltaic-integrated generation commitment.

Funding

The author(s) declare that financial support was received for the research, authorship, and/or publication of this article. The funding for this research was provided by the Special Key Project for Technological Innovation and Application Development in Chongqing, under grant number NO.CSTB2024TIAD-KPX0093.

Conflict of interest

Author GW was employed by Chongqing Carbon Energy Technology Co., Ltd. and Sichuan Aizhong Comprehensive Energy Technology Service Co., Ltd.

The remaining authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Generative AI statement

The author(s) declare that no Generative AI was used in the creation of this manuscript.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

15. Nuno E, Cutululis N. Generation of large-scale pv scenarios using aggregated power curves. In: IEEE General Meeting Power & Energy Society, 2017 IEEE Power & Energy Society General Meeting (2017). p. 1–5.
16. Golestaneh F, Gooi HB. Multivariate prediction intervals for photovoltaic power generation (2018)
17. Junior JGDSF, Oozeki T, Ohtake H, Shimose KI, Takashima T, Ogimoto K. Forecasting regional photovoltaic power generation - a comparison of strategies to obtain one-day-ahead data. *Energ Proced* (2014) 57:1337–45. doi:10.1016/j.egypro.2014.10.124
18. Mellit A, Pavan AM. A 24-h forecast of solar irradiance using artificial neural network: application for performance prediction of a grid-connected pv plant at trieste, Italy. *Solar Energy* (2010) 84:807–21. doi:10.1016/j.solener.2010.02.006
19. Yona A, Senjyu T, Funabashi T. Application of recurrent neural network to short-term-ahead generating power forecasting for photovoltaic system. In: *Power engineering society general meeting* (2007). p. 1–6.
20. Chen Y, Wang Y, Kirschen DS, Zhang B. Model-free renewable scenario generation using generative adversarial networks. *IEEE Trans Power Syst* (2017) PP:1. doi:10.1109/TPWRS.2018.2794541
21. Kay SM *Fundamentals of statistical signal processing: practical algorithm development*, 3. Pearson Education IEEE Transactions on Signal Processing (2013).
22. Kirshner H, Unser M, Ward JP. On the unique identification of continuous-time autoregressive models from sampled data. *IEEE Trans Signal Process* (2014) 62:1361–76. doi:10.1109/tsp.2013.2296879
23. Kay S. Fundamentals of statistical signal processing: estimation theory. *Technometrics* (1993) 37:465–6. doi:10.2307/1269750
24. Lee S, Moon S, Kim K, Sung S, Hong Y, Lim W, et al. A comparison of green, delta, and Monte Carlo methods to select an optimal approach for calculating the 95 population-attributable fraction: guidance for epidemiological research. *J Prev Med Public Health = Yebang Uihakhoe chi* (2024) 45:78–89. doi:10.3961/jpmph.2012.45.2.78
25. Özkale MR, Hüsniye A. Bootstrap confidence interval of ridge regression in linear regression model: a comparative study via a simulation study. *Commun Stat - Theor Methods* (2023) 52:7405–41. doi:10.1080/03610926.2022.2045024
26. Yuichiro S, Takashi S, Hiroto H. Testing parallelism and confidence intervals of level difference in an intraclass correlation model with monotone missing data. *Commun Stat - Theor Methods* (2023) 52:6147–60. doi:10.1080/03610926.2022.2026961
27. Chitralok H, Mani K, Harsha B, Rashmi R. Application of isotonic regression in estimating ED_g and its 95% confidence interval by bootstrap method for a biased coin up-and-down sequential dose-finding design. *Indian J Anaesth* (2023) 67:828–31. doi:10.4103/ija.ija_431_23
28. Chittaranjan A. How to understand the 95 risk, odds ratio, and hazard ratio: as simple as it gets. *J Clin Psychiatry* (2023) 84. doi:10.4088/JCP.23f14933
29. Talsma PA. Estimation of median survival time and its 95 using sas proc lifetest. *J Biopharm Stat* (2023) 34:11–3.
30. Kolawole OJ, Oje MM, Betiku OA, Ijarotimi O, Adekanle O, Ndububa DA. Correlation of alanine aminotransferase levels and a histological diagnosis of steatohepatitis with ultrasound-diagnosed metabolic-associated fatty liver disease in patients from a centre in Nigeria. *BMC Gastroenterol* (2024) 24:147.
31. Rubanovich AV. Redefining the critical value of significance level (0.005 instead of 0.05): the bayes trace. *Biol Bull* (2019) 46:1449–57. doi:10.1134/s1062359019110086
32. Kenanidis P, Llompert M, Santos SF, Dabrowska E. Redundancy can hinder adult l2 grammar learning: evidence from case markers of varying salience levels. *Front Psychol* (2024) 15:1368080. doi:10.3389/fpsyg.2024.1368080
33. Draxl C, Clifton A, Hodge BM, Mccaa J. The wind integration national dataset (wind) toolkit. *Appl Energy* (2015) 151:355–66. doi:10.1016/j.apenergy.2015.03.121

Nomenclature

Indices

j	Index of month
t	Index of time (hours)
h	Index of season
P_{pv}	Photovoltaic power output (W)

Variables

E_{solar}	Solar energy received (J)
λ	Energy ratio
$g[n]$	Normalized solar power
$\epsilon[n]$	Environmental interference (e.g., cloud cover)

Models

AR	Autoregressive model
ARMA	Autoregressive moving average model

GAN	Generative adversarial network
CGAN	Conditional generative adversarial network
WGAN	Wasserstein generative adversarial network

Parameters

μ	Mean value of a variable
σ^2	Variance of a variable
a_t	Autoregressive model parameters
ρ	Friction index
k_t	Reflection coefficients in AR model
p	Order of AR model
σ_u^2	Noise variance in AR model
\mathcal{N}	Set of participants in scenario generation

Statistical distributions

Gaussian distribution	Normal distribution, often used to model randomness
------------------------------	---



OPEN ACCESS

EDITED BY

Zhiqin Zhu,
Chongqing University of Posts and
Telecommunications, China

REVIEWED BY

Xiaosong Li,
Foshan University, China
Zhihang Wu,
Zhejiang University, China

*CORRESPONDENCE

Tiande Ma,
✉ 20201303225@stu.xju.edu.cn

RECEIVED 12 February 2025

ACCEPTED 05 March 2025

PUBLISHED 02 April 2025

CITATION

Jia Y and Ma T (2025) Multi-focus image
fusion based on pulse coupled neural
network and WSEML in DTCWT domain.
Front. Phys. 13:1575606.
doi: 10.3389/fphy.2025.1575606

COPYRIGHT

© 2025 Jia and Ma. This is an open-access
article distributed under the terms of the
[Creative Commons Attribution License \(CC
BY\)](#). The use, distribution or reproduction in
other forums is permitted, provided the
original author(s) and the copyright owner(s)
are credited and that the original publication
in this journal is cited, in accordance with
accepted academic practice. No use,
distribution or reproduction is permitted
which does not comply with these terms.

Multi-focus image fusion based on pulse coupled neural network and WSEML in DTCWT domain

Yuan Jia¹ and Tiande Ma^{2*}

¹School of Statistics, Renmin University of China, Beijing, China, ²School of Computer Science and Technology, Xinjiang University, Urumqi, China

The goal of multi-focus image fusion is to merge near-focus and far-focus images of the same scene to obtain an all-focus image that accurately and comprehensively represents the focus information of the entire scene. The current multi-focus fusion algorithms lead to issues such as the loss of details and edges, as well as local blurring in the resulting images. To solve these problems, a novel multi-focus image fusion method based on pulse coupled neural network (PCNN) and weighted sum of eight-neighborhood-based modified Laplacian (WSEML) in dual-tree complex wavelet transform (DTCWT) domain is proposed in this paper. The source images are decomposed by DTCWT into low- and high-frequency components, respectively; then the average gradient (AG) motivate PCNN-based fusion rule is used to process the low-frequency components, and the WSEML-based fusion rule is used to process the high-frequency components; we conducted simulation experiments on the public Lytro dataset, demonstrating the superiority of the algorithm we proposed.

KEYWORDS

multi-focus image, image fusion, DTCWT, PCNN, WSEML

1 Introduction

Multi-focus image fusion is a technique in the field of image processing that combines multiple images, each focused on different objects or regions, into a single image that captures the sharp details from all focal points [1]. This approach is particularly useful in applications where the depth of field is limited, such as in macro photography, surveillance, medical imaging, and robotics [2, 3].

In typical photography, a single image can only present objects within a certain range of focus clearly, leaving objects closer or farther away blurry [4, 5]. However, by capturing several images with different focus points and then combining them through image fusion techniques, it is possible to create a final image that maintains sharpness across a wider range of depths [6–8].

The process of multi-focus image fusion generally involves several key steps: image alignment, where all the images are aligned spatially; focus measurement, where the sharpness of various regions in each image is assessed; and fusion, where the sharpest information from each image is retained [9–11]. Advanced fusion algorithms, including pixel-level, transform-domain, and machine learning-based methods, can be employed to optimize the fusion quality and preserve important features from all focused regions. This technology has a broad range of applications. In medical imaging, it helps to create clearer, more detailed visualizations of organs or tissues. In surveillance, it enhances the clarity of

objects at varying distances. In robotics, it contributes to improved perception by enabling robots to focus on multiple objects simultaneously [12, 13]. As computational power and algorithms continue to advance, multi-focus image fusion is expected to play an increasingly significant role in a variety of fields requiring high-quality visual information [14–17].

Currently, image fusion can be categorized into two types: traditional algorithms and deep learning algorithms [18–20]. Traditional algorithms typically rely on handcrafted features and conventional image processing techniques, such as Laplacian pyramid [21], wavelet transform [22], dual-tree complex wavelet transform (DTCWT) [23], contourlet [24–26], shearlet [27, 28] and gradient-based methods [29], to combine focused regions from multiple images. Mohan et al. [30] introduced the multi-focus image fusion method based on quarter shift dual-tree complex wavelet transform (qshiftN DTCWT) and modified principal component analysis (MPCA) in the Laplacian pyramid (LP) domain, and this method outperforms many state-of-the-art techniques in terms of visual and quantitative evaluations. Mohan et al. [31] introduced the image fusion method based on DTCWT combined with stationary wavelet transform (SWT). Lu et al. [32] introduced the multi-focus image fusion using residual removal and fractional order differentiation focus measure, and this algorithm simultaneously employs nonsubsampling shearlet transform and the sum of Gaussian-based fractional-order differentiation. These methods are generally effective in simpler scenarios, but they may struggle with more complex images, especially when dealing with varying levels of focus and noise. Pulse coupled neural network (PCNN) also has extensive applications in the field of image fusion, Xie et al. [33] proposed the multi-focus image fusion method based on sum-modified Laplacian and PCNN in nonsampled contourlet transform domain, and this method excellently improves the focus clarity.

On the other hand, deep learning has extensive applications in image fusion [34–37], image segmentation [38, 39], and video restoration [40–44], and image super-resolution [45, 46]. Deep learning algorithms leverage convolutional neural networks (CNNs), Transformer, Generative adversarial network (GAN), Mamba and other advanced models to automatically learn features and perform fusion in an end-to-end manner [47–49]. These methods can adapt to a wide range of image complexities, providing more accurate and visually appealing fused images, especially in challenging conditions like low light or high noise environments [50, 51]. Deep learning approaches have shown superior performance in recent years, particularly with the availability of large datasets and powerful computational resources [52, 53].

Inspired by the ideas from the algorithm in Reference [33], in this paper, a novel multi-focus image fusion method based on PCNN and weighted sum of eight-neighborhood-based modified Laplacian (WSEML) in DTCWT domain is proposed. The motivation behind this approach is to achieve a more robust and effective fusion method that can handle complex images with varying focus levels and noise, while also being computationally efficient. The source images are decomposed by DTCWT into low- and high-frequency components, respectively; then the average gradient (AG) motivate PCNN fusion rule is used to process the low-frequency components, and the WSEML-based fusion rule is used to process the high-frequency components. The algorithm's superiority is validated through comparative experiments on public Lytro dataset.

2 DTCWT

The dual-tree complex wavelet transform (DTCWT) is an advanced signal processing technique designed to overcome some of the limitations of the traditional discrete wavelet transform (DWT) [54]. It was introduced to provide better performance in tasks such as image denoising, compression, and feature extraction. The DTCWT is particularly useful for applications where directional sensitivity and shift invariance are important.

The DTCWT provides improved directional information compared to the traditional wavelet transforms. It uses two parallel trees of wavelet filters (hence “dual-tree”), one for the real part and one for the imaginary part. This structure allows for better representation of image features, especially edges and textures, in multiple orientations. Unlike the traditional DWT, which suffers from shift variance (i.e., small translations in the signal can cause large changes in the wavelet coefficients), the DTCWT provides a level of shift invariance [55, 56]. This makes it more robust to small shifts or distortions in the input signal, which is critical for many image and signal processing tasks. The transform uses complex-valued coefficients rather than real-valued coefficients. This allows for better capture of phase information in addition to amplitude, providing more detailed and richer representations of the signal or image. The DTCWT significantly reduces the aliasing effect, a common issue in wavelet transforms when high-frequency components mix with low-frequency ones. The dual-tree structure and the use of complex filters help mitigate this problem [57].

3 The proposed method

The multi-focus image fusion algorithm we proposed can be mainly divided into four steps: image decomposition, low-frequency fusion, high-frequency fusion, and image reconstruction. The structure of the proposed method is shown in Figure 1, and the specific process is as follows.

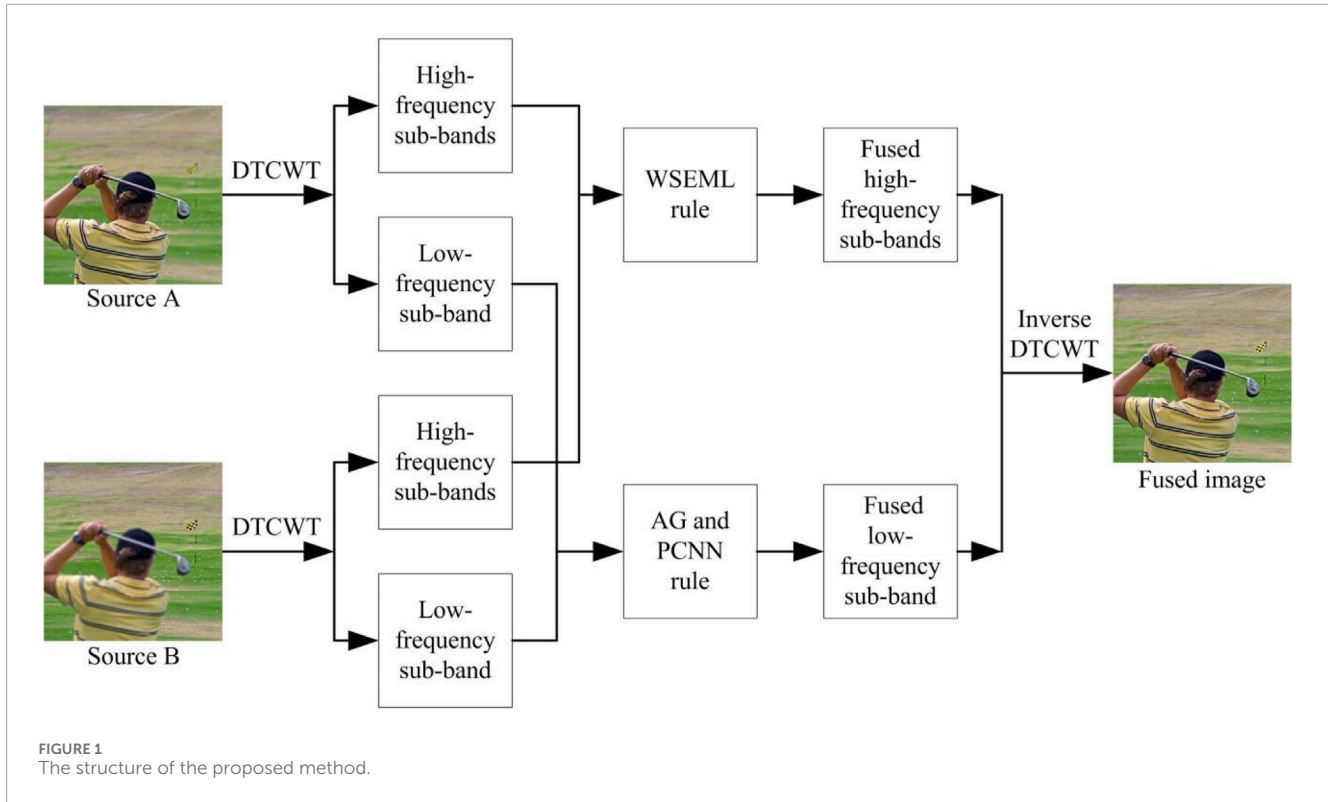
3.1 Image decomposition

The source images A and B are decomposed into low-frequency components $\{L^A, L^B\}$ and high-frequency components $\{H_{l,d}^A, H_{l,d}^B\}$ using DTCWT. The $L^X|X \in (A, B)$ shows the low-frequency, and $H_{l,d}^X|X \in (A, B)$ shows the high-frequency sub-bands l level in the d orientation.

3.2 Low-frequency fusion

The low-frequency component of the image contains the main background information of the image. The average gradient-based (AG) motivate PCNN fusion rule is used to process the low-frequency sub-bands, and the corresponding equations are defined as follows [58, 59]:

$$AG_{ij} = \frac{\sum_i \sum_j ((f(i,j) - f(i+1,j))^2 + (f(i,j) - f(i,j+1))^2)^{\frac{1}{2}}}{mn} \quad (1)$$



$$F_{ij}(n) = AG_{ij} \quad (2)$$

$$L_{ij}(n) = e^{-\alpha_L} L_{ij}(n-1) + V_L \sum_{pq} W_{ij,pq} Y_{ij,pq}(n-1) \quad (3)$$

$$U_{ij}(n) = F_{ij}(n) * (1 + \beta L_{ij}(n)) \quad (4)$$

$$\theta_{ij}(n) = e^{-\alpha_\theta} \theta_{ij}(n-1) + V_\theta Y_{ij}(n-1) \quad (5)$$

$$Y_{ij}(n) = \begin{cases} 1, & \text{if } U_{ij}(n) > \theta_{ij}(n) \\ 0 & \text{else} \end{cases} \quad (6)$$

$$T_{i,j} = T_{i,j}(n-1) + Y_{i,j}(n) \quad (7)$$

In Equation 1, the $f(i,j)$ is pixel intensity at (i,j) and $m \times n$ is the size of the image. In the mathematical model of PCNN in Equations 2–6, the feeding input F_{ij} is equal to the normalized AG_{ij} . The linking input L_{ij} is equal to the sum of neurons firing times in linking range. $W_{ij,pq}$ is the synaptic gain strength and subscripts p and q are the size of linking range in PCNN. α_L is the decay constants. V_L and V_θ are the amplitude gain. β is the linking strength. U_{ij} is total internal activity. θ_{ij} is the threshold. n denotes the iteration times. If U_{ij} is larger than θ_{ij} , then, the neuron will generate a pulse $Y_{ij} = 1$, also called one firing time. In fact, the sum of Y_{ij} in n iteration is often defined as Equation 7, called firing times, to represent image information. Rather than $Y_{ij}(n)$, one often analyzes $T_{ij}(n)$, because neighboring coefficients with similar features representing similar firing times in a given iteration times. AG is input to PCNN to motivate the neurons and generate pulse of

neurons with Equations 2–6. Then, firing times $T_{ij}(n)$ is calculates as Equation 7.

Get the decision map D_{ij} based on Equation 8 and select the coefficients with Equation 9, which means that coefficients with large firing times are selected as coefficients of the fused. The fusion rule is designed as follows:

$$D_{F,ij} = \begin{cases} 1 & \text{If } T_{A,ij}(n) \geq T_{B,ij}(n) \\ 0 & \text{else} \end{cases} \quad (8)$$

$$L^F(i,j) = \begin{cases} L_A(i,j) & \text{If } D_{ij}(n) = 1 \\ L_B(i,j) & \text{If } D_{ij}(n) = 0 \end{cases} \quad (9)$$

where L^F shows the fused low-frequency sub-band.

3.3 High-frequency fusion

The high-frequency component of the image contains the detailed information of the image. The weighted sum of eightneighborhood-based modified Laplacian (WSEML) is used to process the high-frequency sub-bands with Equations 10–12 [60]:

$$WSEML_X(i,j) = \sum_{m=-r}^r \sum_{n=-r}^r \Phi(m+r+1, n+r+1) \times EML_X(i+m, j+n) \quad (10)$$

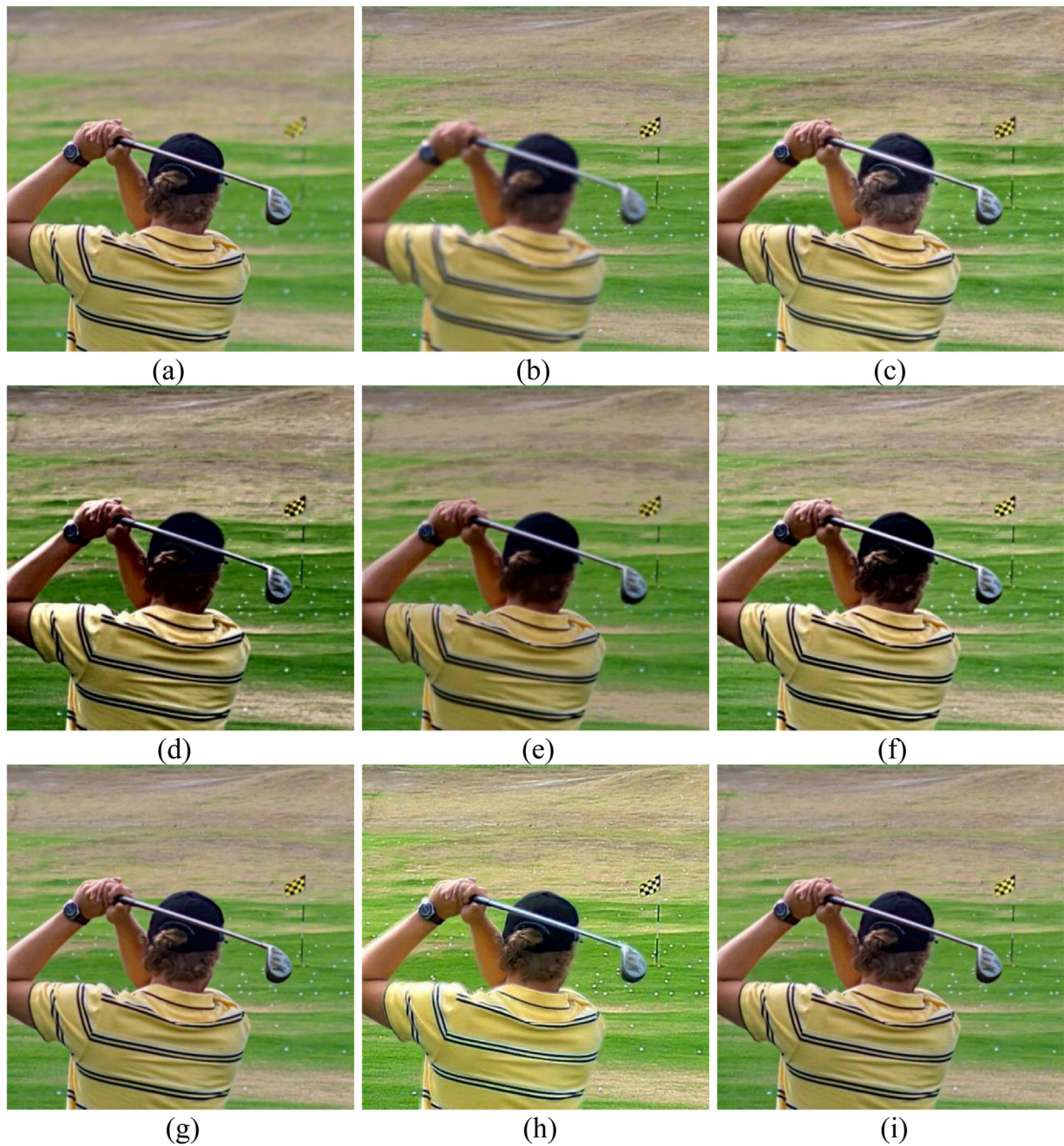


FIGURE 2 Fusion results on Lytro-01. (a) Source A; (b) Source B; (c) GD; (d) FusionDN; (e) PMGI; (f) U2Fusion; (g) ZMFF; (h) UUDFusion; (i) Proposed.

$$\begin{aligned}
 EML_X(i,j) = & |2X(i,j) - X(i-1,j) - X(i+1,j)| \\
 & + |2X(i,j) - X(i,j-1) - X(i,j+1)| \\
 & + \frac{1}{\sqrt{2}} |2X(i,j) - X(i-1,j-1) - X(i+1,j+1)| \\
 & + \frac{1}{\sqrt{2}} |2X(i,j) - X(i-1,j+1) - X(i+1,j-1)|
 \end{aligned} \quad (11)$$

where $X \in \{A, B\}$, and Φ is a $(2r+1) \times (2r+1)$ weighting matrix with radius r . For each element in Φ , its value is set to 2^{2r-d} , where d is its four-neighborhood distance to the center. As an example, the 3×3

normalized version of Φ is

$$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$

The fused high-frequency sub-bands are defined as follows:

$$H_{l,d}^F(i,j) = \begin{cases} H_{l,d}^A(i,j) & \text{if } WSEML_{H_{l,d}^A}(i,j) \geq WSEML_{H_{l,d}^B}(i,j) \\ H_{l,d}^B(i,j) & \text{else} \end{cases} \quad (12)$$

where $H_{l,d}^F(i,j)$ shows the fused high-frequency sub-bands.

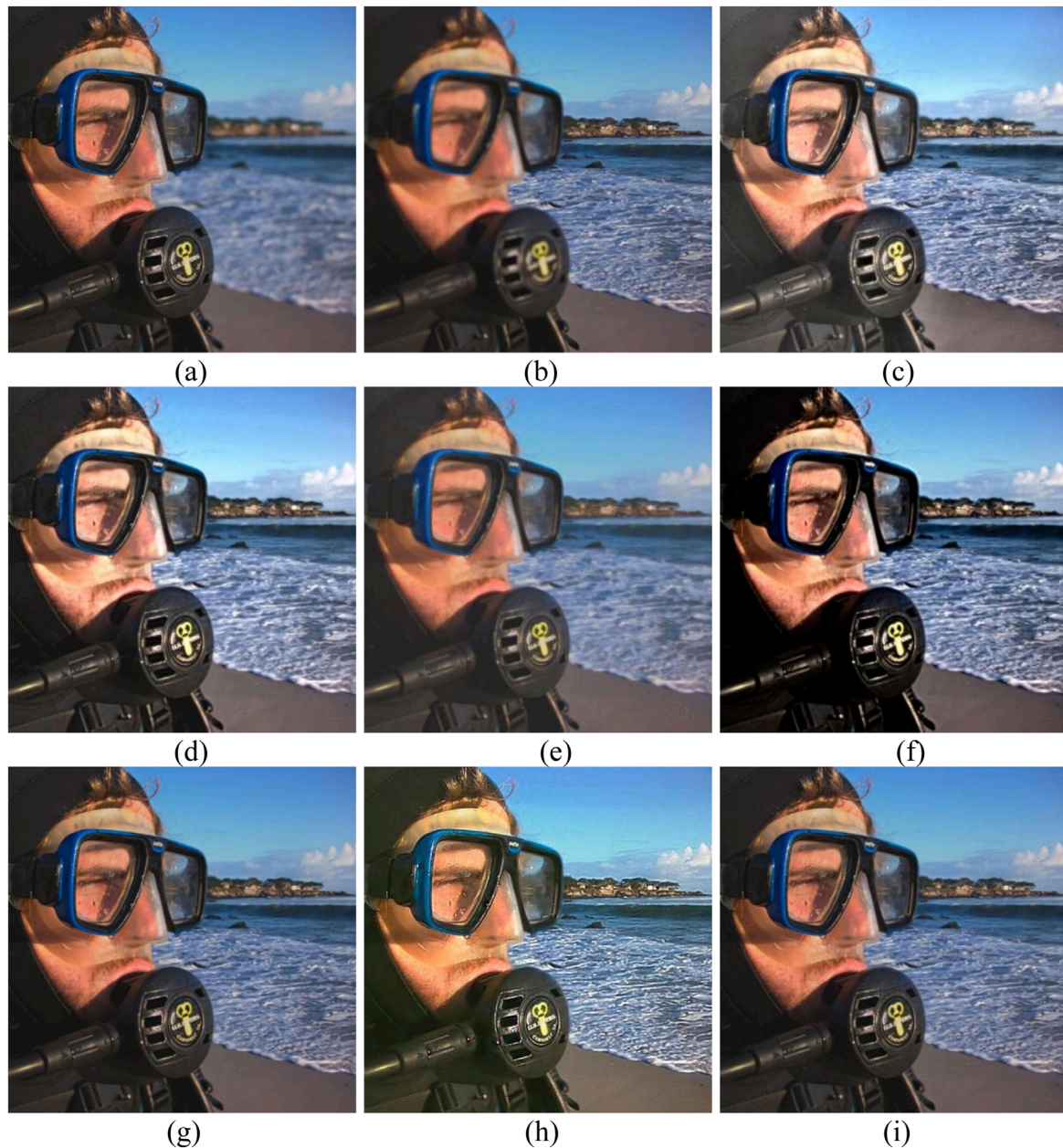


FIGURE 3
Fusion results on Lytro-02. (a) Source A; (b) Source B; (c) GD; (d) FusionDN; (e) PMGI; (f) U2Fusion; (g) ZMFF; (h) UUDFusion; (i) Proposed.

3.4 Image reconstruction

The fused image F is obtained by the inverse DTCWT on $L^F(i, j)$ and $H_{i,d}^F(i, j)$.

4 Experimental results and analysis

To demonstrate the effectiveness of our algorithm, we conducted simulation experiments on the commonly used public Lytro dataset [61] and compared it with six classic image fusion algorithms, namely, GD [29], FusionDN [62], PMGI [63], U2Fusion [64], ZMFF

[65], and UUDFusion [66]. Additionally, we employed six objective evaluation metrics to qualitatively assess the experimental results, namely, edge-based similarity measurement $Q_{AB/F}$ [59], mutual information metric Q_{MI} [59], nonlinear correlation information entropy Q_{NCIE} [67], Chen-Blum metric Q_{CB} [67], image fusion metric-based on phase congruency Q_p [67] and gradient-based fusion performance Q_G [67]. The higher these metric values, the better the fusion effect. We adopt a combined subjective and objective evaluation approach to measure the effectiveness of the algorithms. The parameters of the comparison algorithms were set according to the original papers, while in our algorithm, the decomposition level of DTCWT was set to 4 layers; parameters of

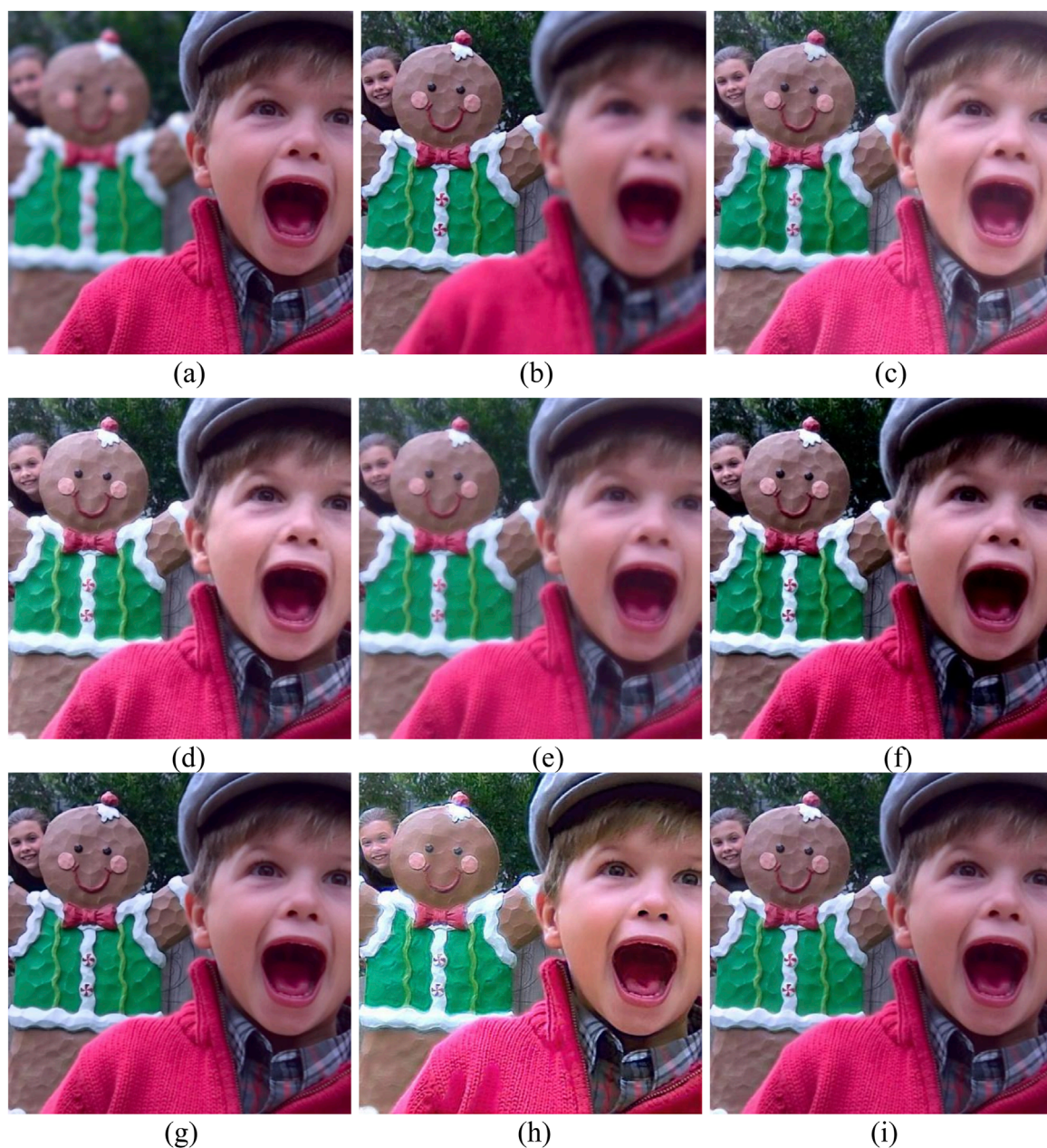


FIGURE 4
Fusion results on Lytro-03. (a) Source A; (b) Source B; (c) GD; (d) FusionDN; (e) PMGI; (f) U2Fusion; (g) ZMFF; (h) UUDFusion; (i) Proposed.

PCNN is set as $p \times q$, $\alpha_L = 0.06931$, $\alpha_\theta = 0.2$, $\beta = 0.2$, $V_L = 1.0$, $V_\theta = 20$, $\Phi = \begin{bmatrix} 0.707 & 1 & 0.707 \\ 1 & 0 & 1 \\ 0.707 & 1 & 0.707 \end{bmatrix}$, and the maximal iterative number is $n = 200$.

Figure 2 shows the fused results with different methods on Lytro-01. The GD method retains significant focus information from both the foreground and background. However, some blending artifacts are visible, and the focus transitions may not be smooth. The FusionDN algorithm preserves structural details well but exhibits some loss of sharpness in the golfer and background.

The fusion quality is moderate, with slight blurring at focus boundaries. The PMGI method achieves reasonable fusion but struggles with preserving contrast and sharpness, especially in the golfer's details. The background appears slightly oversmoothed. The ZMFF method performs well in maintaining the focus of both the foreground (golfer) and background. The details are well-preserved, but minor artifacts can be noticed in the focus transition areas. The UUDFusion method produces an average fusion result, with noticeable blurring in both the foreground and background. The image lacks the clarity and sharpness needed for an effective all-focus image. The proposed method delivers the best results. Both the golfer (foreground) and the background are sharply focused, with smooth

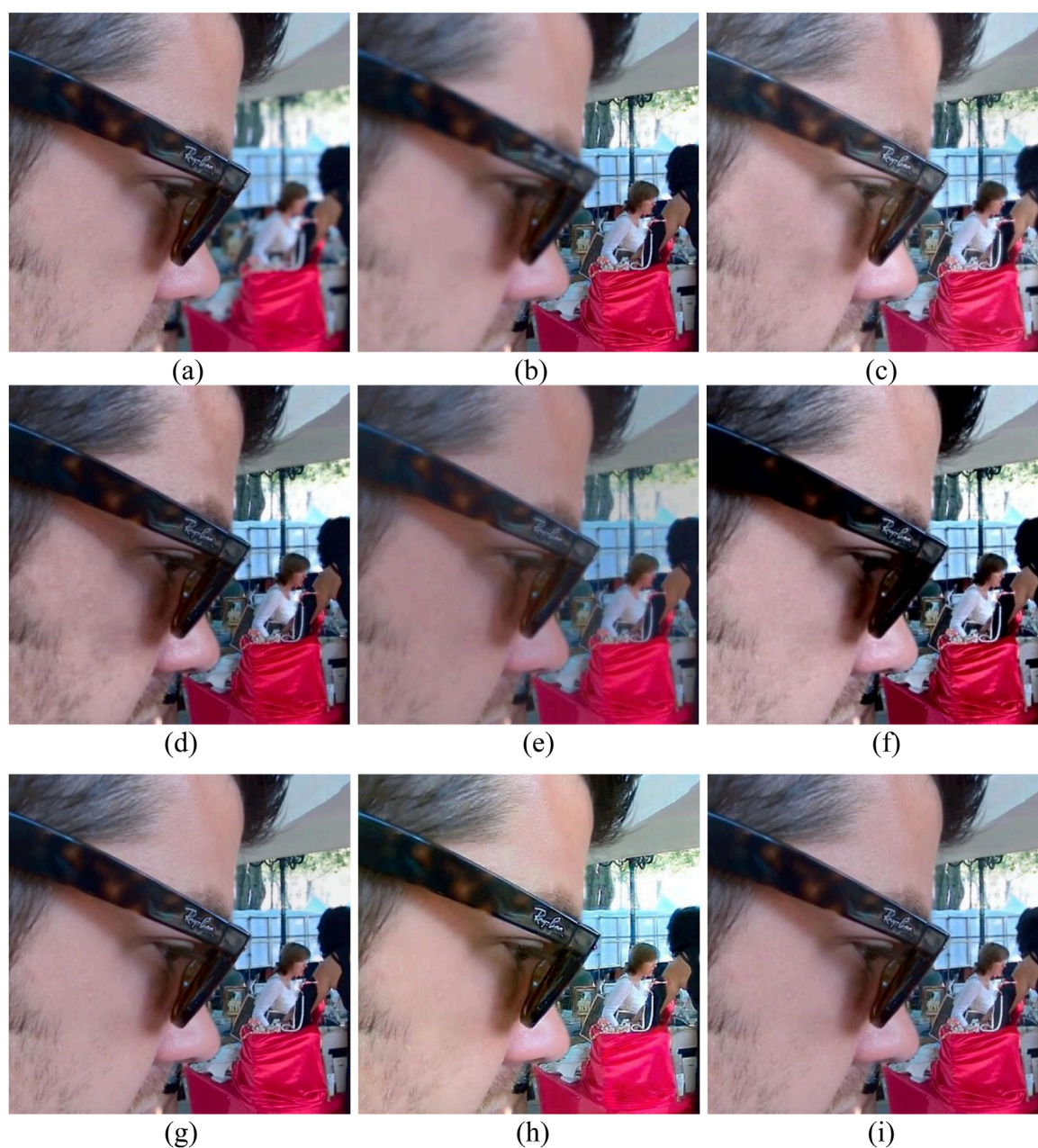


FIGURE 5

Fusion results on Lytro-04. (a) Source A; (b) Source B; (c) GD; (d) FusionDN; (e) PMGI; (f) U2Fusion; (g) ZMFF; (h) UUDFusion; (i) Proposed.

transitions between the focus regions. The image appears natural and well-balanced, with no noticeable artifacts.

Figure 3 presents fusion results for various algorithms applied to the Lytro-02 dataset, aiming to create an all-focus image by combining the near-focus (foreground) and far-focus (background) regions. The proposed method clearly outperforms all other methods, producing a sharp and balanced image where both the diver's face and the background are well-preserved. The transitions between focus regions are smooth and free of noticeable artifacts, resulting in a natural-looking image. ZMFF demonstrates competitive performance, preserving sharpness in both the diver's face and the background. However, slight artifacts and less refined transitions between focus

regions make it less effective than the proposed method. Similarly, FusionDN and U2Fusion provide moderate results, balancing focus between the foreground and background but lacking the sharpness and clarity of the best-performing algorithms. PMGI maintains good detail in the background but struggles with sharpness in the foreground, leading to an imbalanced fusion result. GD performs adequately, but the diver's face appears softened, and overall sharpness is inconsistent. Finally, UUDFusion produces the weakest fusion result, with significant blurring in both focus areas, making it unsuitable for generating high-quality all-focus images. In summary, the proposed method achieves the most visually appealing and technically superior fusion result, while ZMFF serves as a strong

TABLE 1 The average metric values of different methods on Lytro dataset.

	Year	$Q_{AB/F}$	Q_{MI}	Q_{NCIE}	Q_{CB}	Q_P	Q_G
GD	2016	0.7034	3.8521	0.8139	0.6115	0.7466	0.6987
FusionDN	2020	0.6018	5.7908	0.8221	0.6008	0.6221	0.5952
PMGI	2020	0.3901	5.8641	0.8225	0.5656	0.4620	0.3857
U2Fusion	2022	0.6143	5.7765	0.8221	0.5682	0.6657	0.6093
ZMFF	2023	0.7087	6.6271	0.8271	0.7412	0.7853	0.7030
UUDFusion	2024	0.5107	4.8412	0.8178	0.5989	0.5630	0.5055
Proposed		0.7409	7.1960	0.8313	0.7504	0.8137	0.7385

Notes: Bold font indicates the optimal values.

alternative with slight limitations. Other algorithms exhibit varying levels of performance but fall short of achieving the balance and detail provided by the proposed method.

Figure 4 compares the fusion results of multiple algorithms on the Lytro-03 dataset. Each algorithm demonstrates varying capabilities in handling multi-focus image fusion, balancing sharpness, color fidelity, and detail preservation. These are the two input images with distinct focal regions. Source A focuses on the foreground, while Source B highlights the background. The goal of fusion algorithms is to combine these focal regions into a single, sharp image. The GD method struggles with detail preservation and produces a fused image that appears slightly blurred, especially around the edges of the child's face. The colors also seem less vibrant, which detracts from the overall quality. As a deep learning-based approach, FusionDN performs well in preserving details and maintaining sharpness. The child's face and the Cartoon portrait are both clear, with vivid colors. However, minor edge artifacts are noticeable, which slightly impacts the naturalness of the result. The PMGI approach achieves a good balance between sharpness and detail integration. However, it slightly lacks precision in integrating the finest details. The U2Fusion provides decent sharpness and color fidelity but occasionally fails to balance focus across regions. For example, the child's face is slightly less sharp compared to the background, resulting in a less seamless fusion. Some areas also become very dark, resulting in severe information loss. This ZMFF method exhibits noticeable limitations. The fused image lacks sharpness, and the details in both the foreground and background are not well-preserved. The colors are also muted, leading to an overall decrease in visual quality. The image produced by UUDFusion exhibits severe distortion and artifacts, with significant color information loss and poor fusion performance. The proposed method outperforms all others in this comparison. It successfully combines the sharpness and details of both the child's face and the gingerbread figure. The colors are vibrant and natural, with no visible artifacts or blurriness. The transitions between the foreground and background are smooth, creating a visually seamless result.

Figure 5 compares the fusion results of various algorithms on the Lytro-04 dataset, focusing on how well the algorithms preserve details, manage focus regions, and maintain color fidelity. Figure 5a focuses on the foreground, specifically the man's face and sunglasses, while the background is blurred. Figure 5b focuses on the background (the person and chair) but blurs the foreground. Figures 5c-i represent

the fusion results of different algorithms. The GD exhibits moderate sharpness in both the foreground and background. However, some details in the man's sunglasses and the background elements appear slightly smoothed, reducing overall clarity. The color representation is acceptable but lacks vibrancy compared to other methods. As a deep learning-based method, FusionDN achieves good sharpness and color fidelity. The man's face and sunglasses are well-preserved, and the background details are clear. However, subtle edge artifacts are noticeable around the foreground and background transitions, slightly affecting the fusion quality. The PMGI fails to preserve sufficient details in both the foreground and background. The man's sunglasses appear blurred, and the background lacks clarity. The overall image looks less vibrant and exhibits significant information loss, making it one of the weaker methods in this comparison. The overall quality of the fused image is subpar. The U2Fusion method achieves decent fusion but struggles with focus balance. The foreground (sunglasses and face) is slightly less sharp, while the background elements are relatively clear. The ZMFF method produces relatively good fusion results, but the brightness and sharpness of the image still need improvement. The UUDFusion generates noticeable artifacts and distortions, particularly in the background. The details in the foreground (the man's face and sunglasses) are not clear, with significant color distortion, resulting in poor fusion performance. The proposed method demonstrates the best performance among the algorithms. Both the foreground (man's face and sunglasses) and the background (chair and person) are sharp, with vibrant and natural colors. The transitions between the focused regions are smooth, and there are no visible artifacts or distortions. It successfully preserves all critical details, making it the most effective fusion approach in this comparison.

Table 1 shows the average metric values of different algorithms in the simulation experiments on 20 data sets from the Lytro dataset. Table 1 compares the performance of various algorithms on the Lytro dataset across six evaluation metrics: $Q_{AB/F}$, Q_{MI} , Q_{NCIE} , Q_{CB} , Q_P and Q_G . Each metric highlights different aspects of image fusion quality. Among the listed methods, the proposed method demonstrates the best overall performance. It achieves the highest scores in all metrics, such as $Q_{AB/F} = 0.7409$, $Q_{MI} = 7.1960$, $Q_{NCIE} = 0.8313$, $Q_{CB} = 0.7504$, $Q_P = 0.8137$ and $Q_G = 0.7385$. These results suggest that the proposed method is highly robust and effective, delivering superior results across multiple dimensions of evaluation. ZMFF also shows competitive performance. The FusionDN and

U2Fusion maintain balanced performance but fail to excel in any particular metric. UUDFusion performs consistently lower across all metrics, indicating limited effectiveness compared to other algorithms. In summary, the proposed method clearly outperforms all other algorithms, providing the best fusion performance. The ZMFF and GD are strong competitors in specific metrics, but their inconsistencies in other areas limit their overall efficacy. This comparison highlights the superiority of the proposed method for image fusion tasks on the Lytro dataset. These results are consistent with the objective evaluation shown in Figures 2–5.

5 Conclusion

In this paper, a novel multi-focus image fusion method based on pulse coupled neural network and WSEML in DTCWT domain is proposed. The source images are decomposed by DTCWT into low- and high-frequency components, respectively; then the AG and pulse coupled neural network-based fusion rule is used to process the low-frequency components, and the WSEML-based fusion rule is used to process the high-frequency components. The experimental results show that our method achieves better performance in terms of both visual quality and objective evaluation metrics compared to several state-of-the-art image fusion algorithms. The proposed approach effectively preserves important details and edges while reducing artifacts and noise, leading to more accurate and reliable fused images. Future work will focus on further exploring its potential in other image processing tasks.

Data availability statement

The datasets presented in this study can be found in online repositories. The names of the repository/repositories and accession number(s) can be found in the article/supplementary material.

Ethics statement

Written informed consent was obtained from the individual(s) for the publication of any potentially identifiable images or data included in this article.

References

- Bai X, Zhang Y, Zhou F, Xue B. Quadtree-based multi-focus image fusion using a weighted focus-measure. *Inf Fusion* (2015) 22:105–18. doi:10.1016/j.inffus.2014.05.003
- Li H, Shen T, Zhang Z, Zhu X, Song X. EDMF: a new benchmark for multi-focus images with the challenge of exposure difference. *Sensors* (2024) 24:7287. doi:10.3390/s24227287
- Zhang Y, Bai X, Wang T. Boundary finding based multi-focus image fusion through multi-scale morphological focus-measure. *Inf Fusion* (2017) 35:81–101. doi:10.1016/j.inffus.2016.09.006
- Li X, Zhou F, Tan H, Chen Y, Zuo W. Multi-focus image fusion based on nonsubsampling contourlet transform and residual removal. *Signal Processing* (2021) 184:108062. doi:10.1016/j.sigpro.2021.108062
- Liu Y, Qi Z, Cheng J, Chen X. Rethinking the effectiveness of objective evaluation metrics in multi-focus image fusion: a statistic-based approach. *IEEE Trans Pattern Anal Machine Intelligence* (2024) 46:5806–19. doi:10.1109/tpami.2024.3367905
- Zheng K, Cheng J, Liu Y. Unfolding coupled convolutional sparse representation for multi-focus image fusion. *Inf Fusion* (2025) 118:102974. doi:10.1016/j.inffus.2025.102974
- Li X, Li X, Ye T, Cheng X (2024). Bridging the gap between multi-focus and multi-modal: a focused integration framework for multi-modal image fusion. In *Proceedings of the 2024 IEEE winter conference on applications of computer vision (WACV 2024)*, waikoloa, HI, United states, January 4–January 8, 2024, 4–8.
- Li X, Li X, Cheng X, Wang M, Tan H. MCDFFD: multifocus image fusion based on multiscale cross-difference and focus detection. *IEEE Sensors J* (2023) 23:30913–26. doi:10.1109/jsen.2023.3330871
- Zhou Z, Li S, Wang B. Multi-scale weighted gradient-based fusion for multi-focus images. *Inf Fusion* (2014) 20:60–72. doi:10.1016/j.inffus.2013.11.005
- Liu Y, Liu S, Wang Z. Multi-focus image fusion with dense SIFT. *Inf Fusion* (2015) 23:139–55. doi:10.1016/j.inffus.2014.05.004

Author contributions

YJ: Conceptualization, Data curation, Formal Analysis, Methodology, Software, Supervision, Writing–original draft, Writing–review and editing. TM: Data curation, Formal Analysis, Funding acquisition, Methodology, Software, Supervision, Writing–original draft, Writing–review and editing.

Funding

The author(s) declare that financial support was received for the research and/or publication of this article. This research was funded by the Tianshan Talent Training Project–Xinjiang Science and Technology Innovation Team Program (2023TSYCTD).

Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Generative AI statement

The author(s) declare that no Generative AI was used in the creation of this manuscript.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

11. Li S, Kang X, Hu J. Image fusion with guided filtering. *IEEE Trans Image Process* (2013) 22:2864–75. doi:10.1109/TIP.2013.2244222
12. Wang W, Deng L, Vivone G. A general image fusion framework using multi-task semi-supervised learning. *Inf Fusion* (2024) 108:102414. doi:10.1016/j.inffus.2024.102414
13. Wang W, Deng L, Ran R, Vivone G. A general paradigm with detail-preserving conditional invertible network for image fusion. *Int J Comput Vis* (2024) 132:1029–54. doi:10.1007/s11263-023-01924-5
14. Wu X, Cao Z, Huang T, Deng L, Chanussot J, Vivone G. Fully-connected transformer for multi-source image fusion. *IEEE Trans Pattern Anal Machine Intelligence* (2025) 47:2071–88. doi:10.1109/tpami.2024.3523364
15. Li J, Li X, Li X, Han D, Tan H, Hou Z, et al. Multi-focus image fusion based on multiscale fuzzy quality assessment. *Digital Signal Process.* (2024) 153:104592. doi:10.1016/j.dsp.2024.104592
16. Wan H, Tang X, Zhu Z, Xiao B, Li W. Multi-focus color image fusion based on quaternion multi-scale singular value decomposition. *Front Neurorobot* (2021) 15:695960. doi:10.3389/fnbot.2021.695960
17. Li X, Li X, Tan H, Li J (2024). SAMF: small-area-aware multi-focus image fusion for object detection. In *Proceedings of the IEEE international conference on acoustics, speech and signal processing (ICASSP)*, Seoul, Korea, pp. 3845–9.
18. Basu S, Singhal S, Singh D. Multi-focus image fusion: a systematic literature review. *SN COMPUT SCI* (2025) 6:150. doi:10.1007/s42979-025-03678-y
19. Li J, Chen L, An D, Feng D, Song Y. A novel method for CSAR multi-focus image fusion. *Remote Sens.* (2024) 16:2797. doi:10.3390/rs16152797
20. Zhang X. Deep learning-based multi-focus image fusion: a survey and a comparative study. *IEEE Trans Pattern Anal Mach Intell* (2022) 44:4819–38. doi:10.1109/tpami.2021.3078906
21. Liu Y, Liu S, Wang Z. A general framework for image fusion based on multi-scale transform and sparse representation. *Inf Fusion* (2015) 24:147–64. doi:10.1016/j.inffus.2014.09.004
22. Giri A, Sagan V, Alifu H, Maiwulanjiang A, Sarkar S, Roy B, et al. A wavelet decomposition method for estimating soybean seed composition with hyperspectral data. *Remote Sens* (2024) 16:4594. doi:10.3390/rs16234594
23. Wang F, Chen T. A dual-tree-complex wavelet transform-based infrared and visible image fusion technique and its application in tunnel crack detection. *Appl Sci* (2024) 14:114. doi:10.3390/app14010114
24. Vivone G, Deng L, Deng S, Hong D, Jiang M, Li C, et al. Deep learning in remote sensing image fusion methods, protocols, data, and future perspectives. *IEEE Geosci Remote Sensing Mag* (2024) 2:4–43. doi:10.1109/mgrs.2024.3495516
25. Wang G, Li J, Tan H, Li X. Fusion of full-field optical angiography images via gradient feature detection. *Front Phys* (2024) 12:1397732. doi:10.3389/fphy.2024.1397732
26. Zhu Z, Zheng M, Qi G, Wang D, Xiang Y. A phase congruency and local laplacian energy based multi-modality medical image fusion method in NSCT domain. *IEEE Access* (2019) 7:20811–24. doi:10.1109/access.2019.2898111
27. Chen H, Wu Z, Sun Z, Yang N, Menhas M, Ahmad B. CsdFusion: an infrared and visible image fusion method based on LatLRR-NSST and compensated saliency detection. *J Indian Soc Remote Sens* (2025) 53:117–34. doi:10.1007/s12524-024-01987-y
28. Ramakrishna Y, Agrawal R. Pan-sharpening through weighted total generalized variation driven spatial prior and shearlet transform regularization. *J Indian Soc Remote Sens* (2024) 53:681–91. doi:10.1007/s12524-024-02006-w
29. Paul S, Sevcenco I, Agathoklis P. Multi-exposure and multi-focus image fusion in gradient domain. *J Circuits Syst Comput* (2016) 25:1650123. doi:10.1142/s0218126616501231
30. Mohan CR, Chouhan K, Rout RK, Sahoo KS, Jhanjhi NZ, Ibrahim AO, et al. Improved procedure for multi-focus images using image fusion with qshiftN DTCWT and MPCA in Laplacian pyramid domain. *Appl Sci* (2022) 12:9495. doi:10.3390/app12199495
31. Mohan CR, Kiran S, Vasudeva. Improved procedure for multi-focus image quality enhancement using image fusion with rules of texture energy measures in the hybrid wavelet domain. *Appl Sci* (2023) 13:2138. doi:10.3390/app13042138
32. Lu J, Tan K, Li Z, Chen J, Ran Q, Wang H. Multi-focus image fusion using residual removal and fractional order differentiation focus measure. *Signal Image Video Process.* (2024) 18:3395–410. doi:10.1007/s11760-024-03002-w
33. Xie Q, Yi B. Multi-focus image fusion based on SML and PCNN in NSCT domain. *Computer Sci* (2017) 44:266–9.
34. Wu Z, Zhang K, Xuan H, Yuan X, Zhao C. Divide-and-conquer model based on wavelet domain for multi-focus image fusion. *Signal Processing: Image Commun* (2023) 116:116982. doi:10.1016/j.image.2023.116982
35. Liu Y, Shi Y, Mu F, Cheng J, Chen X. Glioma segmentation-oriented multi-modal MR image fusion with adversarial learning. *IEEE/CAA J Automatica Sinica* (2022) 9:1528–31. doi:10.1109/jas.2022.105770
36. Zhang K, Wu Z, Yuan X, Zhao C. CFNet: context fusion network for multi-focus images. *The Institution of Engineering and Technology*. 16 (2022) 499–508.
37. Shi Y, Liu Y, Cheng J, Wang Z, Chen X. VDMUFusion: a versatile diffusion model-based unsupervised framework for image fusion. *IEEE Trans Image Process* (2025) 34:441–54. doi:10.1109/tip.2024.3512365
38. Zhu Z, He X, Qi G, Li Y, Cong B, Liu Y. Brain tumor segmentation based on the fusion of deep semantics and edge information in multimodal MRI. *Inf Fusion* (2023) 91:376–87. doi:10.1016/j.inffus.2022.10.022
39. Zhu Z, Wang Z, Qi G, Mazur N, Yang P, Liu Y. Brain tumor segmentation in MRI with multi-modality spatial information enhancement and boundary shape correction. *Pattern Recognition* (2024) 153:110553. doi:10.1016/j.patcog.2024.110553
40. Wu Z, Sun C, Xuan H, Zhang K, Yan Y. Divide-and-conquer completion network for video inpainting. *IEEE Trans Circuits Syst Video Technology* (2023) 33:2753–66. doi:10.1109/tcsvt.2022.3225911
41. Wu Z, Sun C, Xuan H, Liu G, Yan Y. WaveFormer: wavelet transformer for noise-robust video inpainting. *AAAI Conf Artif Intelligence* (2024) 38:6180–8. doi:10.1609/aaai.v38i6.28435
42. Wu Z, Chen K, Li K, Fan H, Yang Y. BVINet: unlocking blind video inpainting with zero annotations. *arXiv* (2025). doi:10.48550/arXiv.2502.01181
43. Chen K, Wu Z, Hou W, Li K, Fan H, Yang Y. Prompt-aware controllable shadow removal. *arXiv* (2025). doi:10.48550/arXiv.2501.15043
44. Wang F, Guo D, Li K, Zhong Z, Wang M (2024). Frequency decoupling for motion magnification via multi-level isomorphic architecture. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)*, Seattle, WA. 18984–94.
45. Li J, Zheng K, Gao L, Han Z, Li Z, Chanussot J. Enhanced deep image prior for unsupervised hyperspectral image super-resolution. *IEEE Trans Geosci Remote Sensing* (2025) 63:1–18. doi:10.1109/tgrs.2025.3531646
46. Li J, Zheng K, Gao L, Ni L, Huang M, Chanussot J. Model-informed multistage unsupervised network for hyperspectral_enspnsuper-resolution. *IEEE Trans Geosci Remote Sensing* (2024) 62:5516117. doi:10.1109/TGRS.2024.3391014
47. Li S, Huang S. AFA-Mamba: adaptive feature alignment with global-local mamba for hyperspectral and LiDAR data classification. *Remote Sensing* (2024) 16:4050. doi:10.3390/rs16214050
48. Ouyang Y, Zhai H, Hu H, Li X, Zeng Z. FusionGCN: multi-focus image fusion using superpixel features generation GCN and pixel-level feature reconstruction CNN. *Expert Syst Appl* (2025) 262:125665. doi:10.1016/j.eswa.2024.125665
49. Liu Y, Chen X, Peng H, Wang Z. Multi-focus image fusion with a deep convolutional neural network. *Inf Fusion* (2017) 36:191–207. doi:10.1016/j.inffus.2016.12.001
50. Feng S, Wu C, Lin C, Huang M. RADFNet: an infrared and visible image fusion framework based on distributed network. *Front Plant Sci* (2023) 13:1056711. doi:10.3389/fpls.2022.1056711
51. Li H, Ma H, Cheng C, Shen Z, Song X, Wu X. Conti-Fuse: a novel continuous decomposition-based fusion framework for infrared and visible images. *Inf Fusion* (2025) 117:102839. doi:10.1016/j.inffus.2024.102839
52. Zhang Y, Liu Y, Sun P, Yan H, Zhao X, Zhang L. IFCNN: a general image fusion framework based on convolutional neural network. *Inf Fusion* (2020) 54:99–118. doi:10.1016/j.inffus.2019.07.011
53. Gao X, Liu S. BCMFIFuse: a bilateral cross-modal feature interaction-based network for infrared and visible image fusion. *Remote Sens* (2024) 16:3136. doi:10.3390/rs16173136
54. Selesnick I, Baraniuk R, Kingsbury N. The dual-tree complex wavelet transform. *IEEE Signal Process. Mag* (2005) 22:123–51. doi:10.1109/msp.2005.1550194
55. Jiang J, Zhai H, Yang Y, Xiao X, Wang X. Multi-focus image fusion method based on adaptive weighting and interactive information modulation. *Multimedia Syst* (2024) 30:290. doi:10.1007/s00530-024-01506-6
56. Vishwanatha JS, Srinivasa Pai P, D'Mello G, Sampath Kumar L, Bairy R, Nagaral M, et al. Image-processing-based model for surface roughness evaluation in titanium based alloys using dual tree complex wavelet transform and radial basis function neural networks. *Sci Rep* (2024) 14:28261. doi:10.1038/s41598-024-75194-7
57. Ghosh T, Jayanthi N. Multimodal fusion of different medical image modalities using optimised hybrid network. *Int J Ad Hoc Ubiquitous Comput* (2025) 48:19–33. doi:10.1504/ijahuc.2025.143546
58. Shreyamsha Kumar BK. Image fusion based on pixel significance using cross bilateral filter. *Signal Image Video Process.* (2015) 9:1193–204. doi:10.1007/s11760-013-0556-9
59. Qu X, Yan J, Xiao H, Zhu ZQ. Image fusion algorithm based on spatial frequency-motivated pulse coupled neural networks in nonsubsampling contourlet transform domain. *Acta Autom* (2008) 34:1508–14. doi:10.1016/s1874-1029(08)60174-3
60. Yin M, Liu X, Liu Y, Chen X. Medical image fusion with parameter-adaptive pulse coupled neural network in nonsubsampling shearlet transform domain. *IEEE Trans Instrumentation Meas* (2019) 68:49–64. doi:10.1109/tim.2018.2838778

61. Nejati M, Samavi S, Shirani S. Multi-focus image fusion using dictionary-based sparse representation. *Inf Fusion* (2015) 25:72–84. doi:10.1016/j.inffus.2014.10.004
62. Xu H, Ma J, Le Z, Jiang J, Guo X. (2020). FusionDN: a unified densely connected network for image fusion. *Proc Thirty-Fourth AAAI Conf Artif Intelligence (Aaai)*. 34. 12484–91. doi:10.1609/aaai.v34i07.6936
63. Zhang H, Xu H, Xiao Y, Guo X, Ma J. (2020). Rethinking the image fusion: a fast unified image fusion network based on proportional maintenance of gradient and intensity. *Proc AAAI Conf Artif Intelligence*. 34. 12797–804. doi:10.1609/aaai.v34i07.6975
64. Xu H, Ma J, Jiang J, Guo X, Ling H. U2Fusion: a unified unsupervised image fusion network. *IEEE Trans Pattern Anal Mach Intell* (2022) 44:502–18. doi:10.1109/tpami.2020.3012548
65. Hu X, Jiang J, Liu X, Ma J. ZMFF: zero-shot multi-focus image fusion. *Inf Fusion* (2023) 92:127–38. doi:10.1016/j.inffus.2022.11.014
66. Wang X, Fang L, Zhao J, Pan Z, Li H, Li Y. UUD-Fusion: an unsupervised universal image fusion approach via generative diffusion model. *Computer Vis Image Understanding* (2024) 249:104218. doi:10.1016/j.cviu.2024.104218
67. Liu Z, Blasch E, Xue Z, Zhao J, Laganier R, Wu W. Objective assessment of multiresolution image fusion algorithms for context enhancement in night vision: a comparative study. *IEEE Trans Pattern Anal Mach Intell* (2012) 34:94–109. doi:10.1109/tpami.2011.109



OPEN ACCESS

EDITED BY

Huafeng Li,
Kunming University of Science and
Technology, China

REVIEWED BY

Bo Jiang,
Northwest A&F University, China
Chaoxun Guo,
The Chinese University of Hong Kong,
Shenzhen, China
Tianhao Peng,
Moutai College, China

*CORRESPONDENCE

Yuansen Zhang,
✉ zhangyuansen06@163.com

RECEIVED 24 February 2025

ACCEPTED 14 March 2025

PUBLISHED 07 April 2025

CITATION

Zhang Y, Zhuang M, Chen W, Wu X and
Song Q (2025) GLI-Net: A global and local
interaction network for accurate classification
of gastrointestinal diseases in endoscopic
images.
Front. Phys. 13:1582245.
doi: 10.3389/fphy.2025.1582245

COPYRIGHT

© 2025 Zhang, Zhuang, Chen, Wu and Song.
This is an open-access article distributed
under the terms of the [Creative Commons
Attribution License \(CC BY\)](#). The use,
distribution or reproduction in other forums is
permitted, provided the original author(s) and
the copyright owner(s) are credited and that
the original publication in this journal is cited,
in accordance with accepted academic
practice. No use, distribution or reproduction
is permitted which does not comply with
these terms.

GLI-Net: A global and local interaction network for accurate classification of gastrointestinal diseases in endoscopic images

Yuansen Zhang*, Mengxiao Zhuang, Wenjun Chen, Xiaoqiu Wu and Qingqing Song

Department of gastroenterology, The Third Affiliated Hospital of Wenzhou Medical University, Wenzhou, Zhejiang, China

The accurate classification of gastrointestinal diseases from endoscopic images is essential for early detection and treatment. However, current methods face challenges in effectively integrating both global and local features, which limits their ability to capture both broad semantic information and subtle lesion details, ultimately affecting classification performance. To address this issue, this study introduces a novel deep learning framework, the Global and Local Interaction Network (GLI-Net). The GLI-Net consists of four main components: a Global Branch Module (GB) designed to extract global image features, a Local Branch Module (LB) focused on capturing detailed lesion features, an Information Exchange Module (LEM) that facilitates bidirectional information exchange and fusion between the global and local features, and an Adaptive Feature Fusion and Enhancement Module (AFE) aimed at optimizing the fused features. By integrating these modules, GLI-Net effectively captures and combines multi-level feature information, which improves both the accuracy and robustness of endoscopic image classification. Experiments conducted using the Kvasir and Hyper-Kvasir public datasets demonstrate that GLI-Net outperforms existing state-of-the-art models across several metrics, including accuracy, F1 score, precision, and recall. Additionally, ablation studies confirm the contribution of each module to the overall system performance. In summary, GLI-Net's advanced feature extraction and fusion techniques significantly enhance medical endoscopic image classification, highlighting its potential for use in complex medical image analysis tasks.

KEYWORDS

endoscopic image classification, deep learning, global and local feature fusion, global branch module, local branch module

1 Introduction

Gastrointestinal cancers are among the most common cancer types globally, affecting not only the United States but also many other countries. In 2023, it is estimated that there will be approximately 153,020 new cases of gastrointestinal cancer and 52,550 related deaths worldwide. Of these, colorectal cancer accounts for about 34.97% of gastrointestinal cancers. It is well-established that certain intestinal conditions, such as polyps and ulcers, play a significant role in the development of colorectal cancer. Early detection of cancer indicators is crucial for managing colorectal cancer, as it can notably improve patient outcomes and

survival rates. Therefore, early diagnosis is a critical component in the fight against this cancer, offering hope for better prognoses and higher survival chances.

Endoscopy remains a key method for the initial identification and evaluation of colorectal cancer, demonstrating its effectiveness in reducing mortality rates. This diagnostic tool captures numerous visual frames during gastrointestinal examinations, which are typically reviewed manually. This manual process is not only labor-intensive and repetitive but also subject to human error, as the accuracy of diagnosis depends on the endoscopist's expertise, experience, and mental acuity. Such variability can result in incorrect diagnoses or missed abnormalities. To address these challenges, there is an urgent need for a precise, advanced computer-assisted diagnostic system. This system would autonomously identify and flag suspicious images, reducing the significant manual workload for endoscopists and improving diagnostic accuracy. This technological innovation is poised to advance the early detection of colorectal cancer, potentially leading to better patient outcomes and increased survival rates.

For instance, Karargyris and Bourbakis [1] proposed a method using image processing techniques to detect polyps and ulcers in wireless capsule endoscopy videos, achieving improved detection rates. Mesejo et al. [2] developed a computer-aided system based on computer vision and machine learning for classifying gastrointestinal lesions in regular colonoscopy images, enhancing diagnostic accuracy. Charfi et al. [3] combined the local binary pattern variance and discrete wavelet transform to make texture extraction for wireless capsule endoscopy images. However, despite the fact that computer-aided diagnosis systems is beneficial for endoscopic image classification compared with human beings, it still encounters significant obstacles. Primarily, due to the high variability within the same class of samples, such as differences in size and shape of lesions, the extraction of consistent features from the same category is quite difficult. By contrast, the subtle differences between different classes also present a challenge in accurate classification, where the different samples from different classes may have the similar attributes. Furthermore, interference factors like bubbles, turbidity, and artifacts caused by the movement of the capsule camera during endoscopic procedures can also significantly reduce the detection rate of abnormal images. Obviously, these factors contribute to the overall difficulty in achieving high accuracy in endoscopic image classification, emphasizing the need for more advanced algorithms and techniques to address these challenges.

In recent years, deep learning, particularly convolutional neural networks (CNNs) [4–6], has made significant strides in the field of endoscopic image classification [7–9]. These technologies have automated medical image analysis, reducing the workload for physicians and enabling more efficient disease diagnosis through feature extraction and pattern recognition. Compared to traditional methods, deep learning models have demonstrated higher precision and recall. Deep learning's ability to learn from data has made it superior in tasks such as polyp detection, lesion classification, and region recognition, outperforming traditional algorithms in terms of speed and accuracy [10]. However, despite these advancements, deep learning models in endoscopic image classification have yet to reach a level suitable for widespread clinical application. There are still many challenges such as the requirement for large annotated datasets and the difficulty in achieving higher diagnostic precision

for rare or subtle pathologies. There is a need for more effective methods to enhance the classification accuracy of endoscopic images and address these limitations before deep learning can be fully integrated into clinical practice.

In this paper, we introduce a novel deep learning approach for classifying endoscopic images called GLI-Net (Global and Local Interaction Network). GLI-Net addresses the shortcomings of traditional methods in capturing both detailed features and global semantic information by effectively combining global and local features, leading to significant improvements in classification accuracy and robustness. The network is composed of four primary modules: the Global Branch Module (GB), which extracts global features and guides the Local Branch Module (LB); the Local Branch Module (LB), which focuses on extracting detailed features from lesion regions; the Information Interaction Module (LEM), which facilitates mutual information exchange and optimization between the global and local branches; and the Adaptive Feature Fusion and Enhancement Module (AFE), which adaptively fuses the global and local features, enhancing their representational power and boosting the model's discriminative performance. The synergistic interaction of these modules enables GLI-Net to achieve superior results in medical image classification.

2 Related work

In this section, we will briefly describe the related works about classification on the endoscopic images. Due to different approaches used in this field, we divide the related works into two branches: human-crafted feature based methods and deep learning based methods.

2.1 Human-crafted feature based methods

For the human-crafted feature based methods, many machine learning methods with different images features designed by human beings were studied. For instance, Charfi and El Ansari demonstrated that their computer-aided diagnosis system can effectively detect colon abnormalities in wireless capsule endoscopy images [3]. The system employed image preprocessing to enhance quality, extracted key features such as color and texture, and used a support vector machine (SVM) classifier for abnormality detection. Their results verified that integrating color and texture features with SVM significantly can improve detection accuracy compared to manual analysis. This approach highlights the potential of feature-based machine learning methods for automating gastrointestinal disorder diagnosis in clinical practice. Furthermore, Mesejo et al. [2] made a study on how to apply the computer technology to diagnose gastrointestinal lesions from regular colonoscopic videos. Specifically, it exploited both computer vision and machine learning methods, conducting a virtual biopsy to differentiate hyperplastic lesions, serrated adenomas, and adenomas. Karargyris and Bourbakis [1] conducted a study on the detection of small bowel polyps and ulcers using wireless capsule endoscopy videos. Specifically, they developed an algorithm that leveraged image processing techniques to identify and analyze these gastrointestinal

abnormalities, contributing to the advancement of non-invasive diagnostic methods.

Additionally, Li and Meng [11] developed an enhancement method based on adaptive contrast diffusion. This technique was designed to adjust the contrast in different regions of the image dynamically, which helped in highlighting the features of interest, particularly in the context of the gastrointestinal tract. By increasing the contrast, the method aimed to make it easier for medical professionals to identify and diagnose any abnormalities or pathologies within the small bowel. The enhancements are intended to facilitate a more accurate and reliable analysis of the endoscopic images, which is vital for effective clinical decision-making. The work by Souaidi and Ansari [12] delved into the detection of ulcer diseases from wireless capsule endoscopy images, employing a multi-scale analysis technique. Specifically, this approach involved examining images across various scales to identify ulcers of different sizes and shapes within the gastrointestinal tract, which enhanced the detection accuracy by capturing the nuances of ulcer appearances at multiple levels of detail.

2.2 Deep learning based methods

Different from human-crafted features based methods, deep learning based methods can automatically extract more semantic features for classification. For instance, Zhang et al. [13] focused on the automatic detection and classification of colorectal polyps by leveraging low-level CNN features from nonmedical domains. Specifically, the authors explored the transfer learning approach where pre-trained CNN models originally trained on nonmedical images were adapted for the task of polyp detection in endoscopic videos. The study aimed to demonstrate that features learned from large datasets in nonmedical domains could be effectively transferred to enhance the performance of medical image analysis tasks, particularly in the context of colorectal polyp identification. Shin and Balasingham [14] conducted a comparative study between a hand-crafted feature-based SVM and a CNN based deep learning framework for the automatic classification of polyps. They evaluated the performance of both methods in distinguishing polyps in endoscopic images, providing insights into the efficacy of deep learning versus traditional machine learning approaches for medical image classification. Zhao et al. [15] presented Adasan, an Adaptive Cosine Similarity Self-Attention Network for gastrointestinal endoscopy image classification, which integrated self-attention mechanisms with adaptive cosine similarity measures to enhance feature representation, improving classification accuracy of endoscopic images.

Furthermore, Zhu et al. [16] presented a method for lesion detection in endoscopy images leveraging features from CNNs. Also, a novel method for WCE video summarization was studied by using a Siamese neural network coupled with SVM, which condensed long WCE video sequences into shorter, representative summaries to facilitate faster and more efficient review by medical professionals. The Siamese network was employed to learn and compare image features, identifying similar frames within the video, while the SVM was utilized to classify these frames based on their medical relevance. Similarly [17], designed a network to identify and highlight potential lesions within the gastrointestinal

tract by analyzing WCE video frames. By extending the Siamese network, Guo et al. [18] introduced the Triple ANet, an Adaptive Abnormal-Aware Attention Network designed for the classification of WCE images. It included three main components: an abnormal region detection module, an attention mechanism to highlight these regions, and a classification module, in which the attention mechanisms was introduced to focus on abnormal regions within the gastrointestinal tract, being crucial for accurate diagnosis. And The paper probably detailed the architecture of the network, how it was trained on WCE images, and its effectiveness in classifying normal versus abnormal images. This approach aimed to improve the accuracy and efficiency of WCE image analysis, providing a valuable tool for medical professionals to detect gastrointestinal abnormalities. Similarly, an Effectively Fused Attention Guided Convolutional Neural Network was proposed to integrated attention mechanisms to enhance feature extraction from endoscopic images, focusing on discriminative regions indicative of gastrointestinal conditions [19,20].

In recent years, many deep learning-based approaches have been applied to classify colorectal cancer and WCE images, yielding promising outcomes. However, due to the inherent characteristics of these images, such as considerable intra-class variations and subtle inter-class differences, there is still a need for more robust models to improve the accuracy and reliability of these algorithms. To overcome these challenges, future research should focus on developing models that are better equipped to handle the complexity and variability of endoscopic images. This could involve exploring advanced network architectures, integrating multi-modal data, or utilizing sophisticated feature extraction methods to capture subtle pathological changes more effectively.

3 Methods

This section provides a detailed description of the overall architecture of GLI-Net (Global and Local Interaction Network). First, the main structure of the network and its global branch module (GB) and local branch module (LB) are introduced. Then, the structure and functionality of the Information Exchange Module (LEM) and the Adaptive Feature Fusion and Enhancement Module (AFE) are discussed in detail. The overall network architecture of GLI-Net is shown in Figure 1.

3.1 Overall network architecture

GLI-Net adopts a dual-branch global and local interaction network structure, as illustrated in Figure 1. The backbone of the network uses the Swin Transformer as a feature extractor, designed to extract both shallow and deep feature maps from the input endoscopic images and generate multi-scale feature representations. The feature sizes correspond to 1/4, 1/8, 1/16, and 1/32 of the input image size, as specified in Equation 1:

$$F_i = f_{Swin}(I), \quad i = 1, 2, 3, 4 \quad (1)$$

where f_{Swin} represents the Swin Transformer, $I \in \mathbb{R}^{H \times W \times 3}$ is the input image, and F_i represents the output multi-scale feature maps.

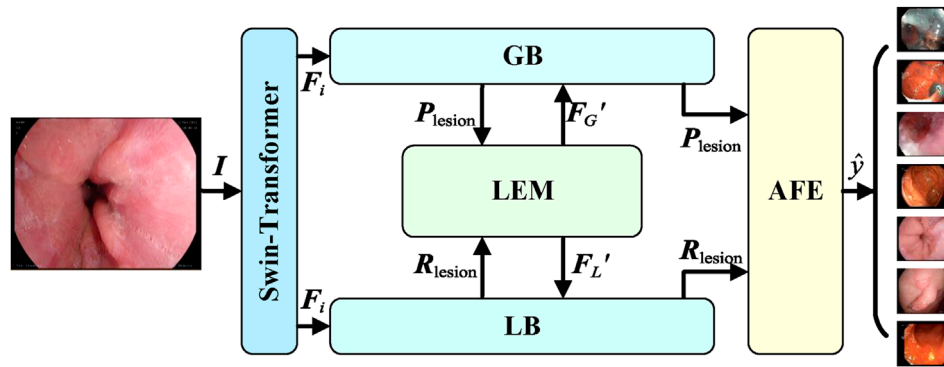


FIGURE 1
Overall architecture of GLI-Net.

These multi-scale features are then fed into the global and local branches, where global and local lesion features are extracted, as shown in Equation 2:

$$P_{\text{lesion}} = f_{GB}(F_i), R_{\text{lesion}} = f_{LB}(F_i) \quad (2)$$

where f_{GB} and f_{LB} denote the global and local branch modules, respectively, while P_{lesion} and R_{lesion} correspond to the outputs of the global and local branches. While the GB and LB modules extract the lesion features, the Information Exchange Module (LEM) facilitates the bidirectional information flow between the global and local features, ensuring their collaborative interaction. This enhances the comprehensiveness and accuracy of the features. The specific formulation is as follows:

$$(F_G', F_L') = f_{LEM}(P_{\text{lesion}}, R_{\text{lesion}}) \quad (3)$$

where F_G' and F_L' represent the enhanced global and local features, respectively. After obtaining the global feature P_{lesion} and the local feature R_{lesion} , the Adaptive Feature Fusion and Enhancement (AFE) module is responsible for fusing the enhanced global and local features, further enhancing their representational capability. Finally, the classifier outputs the corresponding class of the image. The specific formula is as follows:

$$\hat{y} = \text{Softmax}(f_{AFE}(P_{\text{lesion}}, R_{\text{lesion}})) \quad (4)$$

3.2 Global branch module (GB)

To effectively capture the overall lesion information in endoscopic images and guide the local branch module to focus on key regions, the Global Branch module (GB) is introduced. The goal of the GB module is to extract global lesion features from the deepest feature maps and generate lesion category prompts to guide the local branch, thereby enhancing the comprehensiveness of feature representations and improving classification accuracy. The GB module consists of convolutional layers, global adaptive pooling layers, and the Lesion Category Prompt Extractor (LCPE), with the specific structure shown in Figure 2.

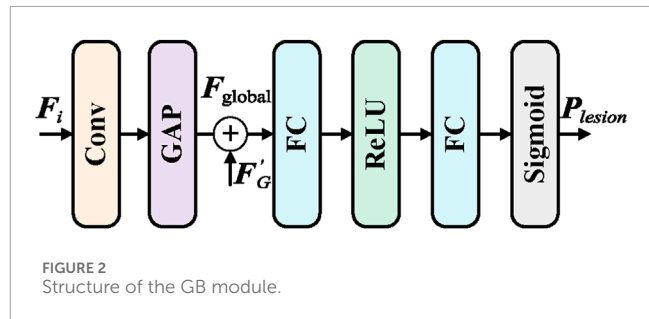


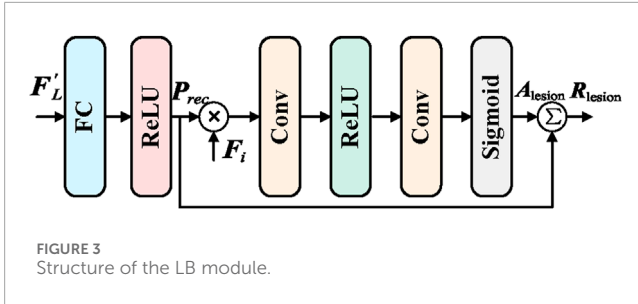
FIGURE 2
Structure of the GB module.

The input global feature map F_i is first processed through a series of convolutional layers to extract high-level semantic features. These convolutional layers effectively capture the global information in the image and enhance the expressive power of the features. Then, global adaptive pooling (GAP) is applied to aggregate the convolved feature map F_{conv} , generating a fixed-size global feature vector F_{global} . Global adaptive pooling automatically adjusts the pooling kernel size based on the input feature map's dimensions and shape, enabling more effective capture and aggregation of global feature information. This process is described by the following Equation 5:

$$F_{\text{global}} = f_{\text{GAP}}(f_{\text{conv}}(F_i)) + F_G' \quad (5)$$

The module f_{conv} contains multiple convolution operations, f_{GAP} represents the global adaptive pooling operation, and F_G' is the output of the LEM module. The GB module generates the lesion category prompt P_{lesion} from the global feature vector F_{global} using the LCPE module, which is used to guide the local branch to focus on the lesion regions. The LCPE module primarily consists of two fully connected layers and their corresponding activation functions. The global feature vector F_{global} is first mapped to the prompt space through the first fully connected layer f_{FC1} , and then the second fully connected layer f_{FC2} generates the final lesion category prompt P_{lesion} . The Sigmoid activation function is applied to ensure that the prompt values lie within the range of [0, 1]. The specific process is described by Equation 6:

$$P_{\text{lesion}} = \text{Sigmoid}(f_{FC2}(\text{ReLU}(f_{FC1}(F_{\text{global}})))) \quad (6)$$



3.3 Local branch module (LB)

In order to further capture detailed lesion information in endoscopic images, and integrate guidance from the global features, the Local Branch Module (LB) is proposed. The main objective of the LB module is to receive enhanced lesion category prompts from the Information Exchange Module (LEM) through the Lesion Category Prompt Receiver (LCPR). These prompts are then used by the Lesion Region Detector (LRD) to identify detailed lesion features. The output detailed features are fed back into the Information Exchange Module and the subsequent Adaptive Feature Fusion and Enhancement (AFE) module, enabling the collaborative enhancement of both global and local features. The structure of the LB module is shown in Figure 3.

The Lesion Category Prompt Receiver (LCPR) module is responsible for receiving the enhanced lesion category prompt F'_L from the Information Exchange Module (LEM) and applying it to the feature map of the local branch to guide the local branch in focusing on potential lesion regions. First, a fully connected layer along with an activation function modulates the prompt features, and these are element-wise multiplied with the initial feature map F_i of the local branch to generate the modulated local feature map P_{rec} . The specific calculation is as follows:

$$P_{rec} = F_i \otimes \text{ReLU}(f_{FC}(F'_L)) \quad (7)$$

where \otimes denotes the element-wise multiplication operation. After obtaining the modulated local feature map P_{rec} , the Lesion Region Detector (LRD) module is responsible for identifying and extracting the detailed lesion information. First, the modulated local feature map P_{rec} undergoes further convolution processing to extract higher-level detailed features. Then, through a series of convolutional layers and pooling layers, an attention map A_{lesion} for the lesion region is generated. Based on this attention map, the local feature map is weighted to extract the detailed feature vector R_{lesion} . The specific calculation is as follows:

$$\begin{cases} A_{lesion} = \text{Sigmoid}(f_{conv}(\text{ReLU}(f_{conv}(P_{rec}))) \\ R_{lesion} = \sum_{i=1}^H \sum_{j=1}^W A_{lesion}(i,j) \cdot P_{rec}(i,j) \end{cases} \quad (8)$$

3.4 Information exchange module (LEM)

In order to enable efficient collaboration between the global and local networks and enhance the overall feature representation

capability, an Information Exchange Module (LEM) has been proposed. The primary goal of the LEM module is to facilitate bidirectional information transfer and mutual supervision between the global branch (GB) and the local branch (LB), thereby improving the comprehensiveness of the features and the accuracy of classification. The detailed structure of the LEM module is shown in Figure 4.

The LEM module includes information transmission from global to local, feedback from local to global, and bidirectional information flow. The information transmission from global to local is responsible for passing the lesion category cue P_{lesion} generated by the GB module to the local branch module (LB) through the Information Exchange Module, guiding the local branch to focus on potential lesion areas. The feedback from local to global is responsible for sending the detailed lesion features R_{lesion} extracted by the local branch module (LB) back to the global branch module (GB), thereby enhancing the representation ability of the global features. The specific calculation details are provided in Equation 9.

$$\begin{cases} F'_G = P_{lesion} + \text{ReLU}(f_{FC}(P_{lesion}) + f_{FC}(R_{lesion})) \\ F'_L = R_{lesion} + \text{ReLU}(f_{conv}(R_{lesion}) + f_{conv}(F'_G)) \end{cases} \quad (9)$$

where F'_L refers to the transformed lesion category cue, and F'_G represents the enhanced global features.

3.5 Adaptive feature fusion and enhancement (AFE) module

To fully integrate global and local features and further enhance the feature representation capability, an Adaptive Feature Fusion and Enhancement (AFE) module has been proposed. The primary objective of the AFE module is to effectively fuse the enhanced global features F'_G with the local features F'_L , and to improve the expressiveness of the fused features through a feature enhancement mechanism, thereby achieving more accurate class predictions. The AFE module employs a learnable weighting mechanism, which dynamically adjusts the fusion ratio between the global and local features based on their relative importance. This mechanism ensures that features from both branches contribute appropriately to the final fused representation. Unlike traditional fusion methods that use fixed weights or simple averaging, this approach allows the model to prioritize more discriminative features from the global and local branches based on the task at hand, leading to enhanced feature representation and classification accuracy. The AFE module consists of feature fusion, feature enhancement, and the final classifier, with its detailed structure shown in Figure 5.

First, the feature fusion component is responsible for adaptively fusing the enhanced global features F'_G from the global branch module (GB) with the enhanced local features F'_L from the local branch module (LB). To achieve this, the AFE module employs a learnable weighting mechanism, as shown in Equation 10:

$$F_{fused} = \alpha \cdot F'_G + \beta \cdot F'_L \quad (10)$$

where α and β are learnable weight parameters obtained through the network, with the constraint $\alpha + \beta = 1$. This allows the model to dynamically adjust the fusion ratio based on the importance of different features, enabling effective integration of global and local

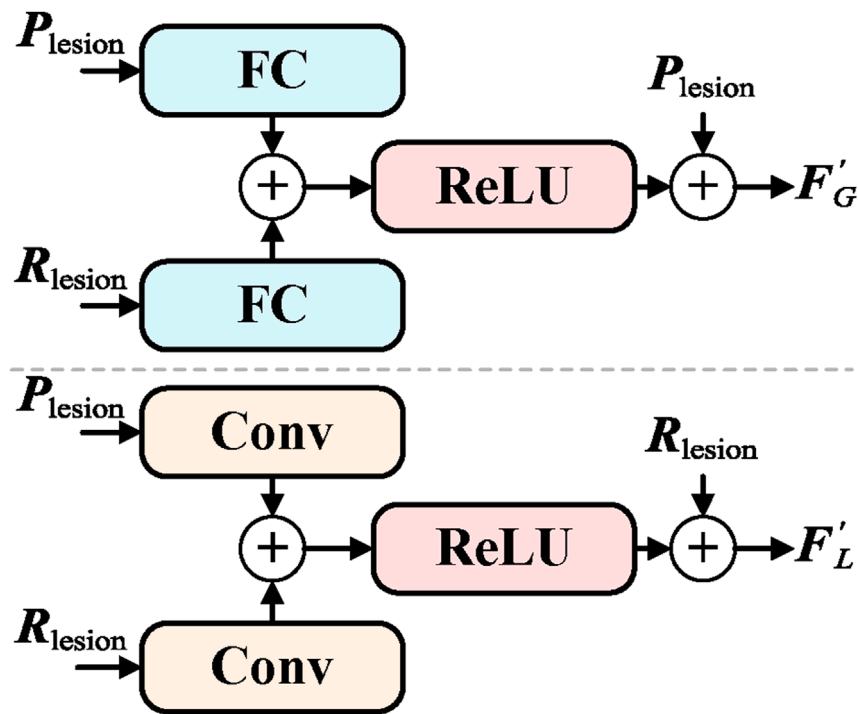


FIGURE 4
Structure of the LEM module.

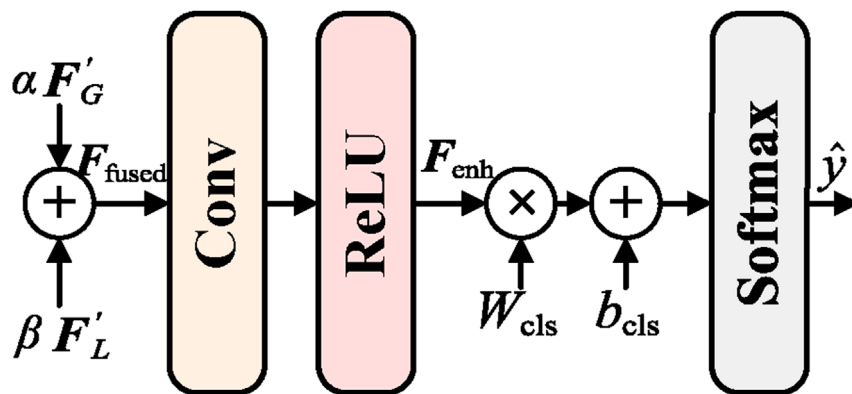


FIGURE 5
Structure of the AFE module.

features. After feature fusion, the AFE module further enhances the expressiveness of the fused features through a feature enhancement layer. The feature enhancement layer typically consists of a series of convolutional layers and activation functions to capture higher-level semantic information, producing the enhanced features F_{enh} . This is detailed in Equation 11:

$$F_{\text{enh}} = \text{ReLU}(f_{\text{conv}}(F_{\text{fused}})) \quad (11)$$

Finally, the enhanced features F_{enh} are input into the classifier component for the final class prediction, as shown in Equation 12:

$$\hat{y} = \text{Softmax}(W_{\text{cls}} \cdot F_{\text{enh}} + b_{\text{cls}}) \quad (12)$$

where W_{cls} and b_{cls} are the weight and bias parameters of the classifier, and \hat{y} represents the predicted class probability distribution.

3.6 Loss function

To effectively train GLI-Net, a comprehensive loss function has been designed, consisting of two main components: lesion region detection loss \mathcal{L}_{det} and classification loss \mathcal{L}_{cls} . The combination of these two loss functions is aimed at simultaneously optimizing the model's ability to identify lesion regions and its overall

classification performance, thereby improving the model's accuracy and robustness in endoscopic image classification tasks. The \mathcal{L}_{det} is designed to optimize the model's ability to detect lesion regions in the image. This loss function uses binary cross-entropy loss, and the main calculation formula is as follows:

$$\mathcal{L}_{\text{det}} = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] \quad (13)$$

where N is the total number of pixels or regions, y_i is the ground truth label for the i -th pixel or region (0 for non-lesion, one for lesion), and \hat{y}_i is the predicted lesion probability for the i -th pixel or region. The classification loss \mathcal{L}_{cls} is used to optimize the model's ability to predict the class of the entire image. This loss function employs categorical cross-entropy loss to measure the difference between the predicted class distribution and the true class labels. The specific calculation details are as follows:

$$\mathcal{L}_{\text{cls}} = -\frac{1}{M} \sum_{j=1}^M \sum_{k=1}^C y_{jk} \log(\hat{y}_{jk}) \quad (14)$$

where M is the number of samples, C is the total number of classes, y_{jk} is the ground truth label of the j -th sample for the k -th class, and \hat{y}_{jk} is the predicted probability of the j -th sample for the k -th class. The overall loss \mathcal{L} combines the lesion region detection loss and the classification loss to achieve simultaneous optimization of the model on both local and global features. The specific formula is as follows:

$$\mathcal{L} = \mathcal{L}_{\text{det}} + \lambda \mathcal{L}_{\text{cls}} \quad (15)$$

4 Experiments

4.1 Experimental details

4.1.1 Dataset

In our experiments, the Kvasir dataset and the Hyper-Kvasir dataset were used. The Kvasir dataset contains 4,000 endoscopic images of gastrointestinal diseases, covering eight categories, with 500 images per category. The dataset includes both anatomical landmarks (such as the Z-line, pylorus, cecum, etc.) and pathological findings (such as esophagitis, polyps, ulcerative colitis, etc.). The image resolutions range from 720×576 to 1920×1072 pixels. In the training and testing split of the dataset, considering the imbalance in the annotation of medical images, the labeled images are divided into a training set (70%), a validation set (15%), and a test set (15%). The Hyper-Kvasir dataset is a large multi-class public gastrointestinal dataset sourced from gastroscopy and colonoscopy exams conducted at the Baerum Hospital in Norway. All image annotations were provided by experienced radiologists. The dataset contains 110,079 images, covering both normal (healthy) and abnormal (unhealthy) patients, with 10,662 labeled images. Due to the scarcity of annotated samples and the large variation in the number of lesion samples across different categories, the dataset split follows the common strategy used in the medical field. Specifically, the 10,662 labeled images are divided into a training set (70%), a validation set (15%), and a test set (15%). These images cover a wide range of gastrointestinal abnormalities, including normal and abnormal conditions, with a particular focus on diseases such as polyps, ulcers, and colorectal cancer. The dataset is diverse,

featuring a variety of lesion shapes, sizes, and textures, which presents significant challenges for model training. The annotated images, provided by experienced radiologists, allow for a comprehensive evaluation of model performance across different disease categories and anatomical regions.

4.1.2 Evaluation metrics

We use accuracy (ACC), F1 score, precision, and recall as classification evaluation metrics. These metrics are all derived from the confusion matrix, where the symbols are defined as True Positive (TP), True Negative (TN), False Positive (FP), and False Negative (FN). The specific calculation formulas are as follows:

$$\text{Acc} = \frac{TP + TN}{TP + TN + FP + FN} \quad (16)$$

$$P = \frac{TP}{TP + FP} \quad (17)$$

$$R = \frac{TP}{TP + FN} \quad (18)$$

$$F1 = 2 \times \frac{P \times R}{P + R} = \frac{2TP}{2TP + FP + FN} \quad (19)$$

4.1.3 Implementation details

The experiments in this study were conducted on a computer equipped with an NVIDIA RTX 4090 GPU with 24 GB of memory. During training, the Adam optimizer was used, with specific parameters set as $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 10^{-6}$. The learning rate followed a cosine annealing strategy, with an initial value of 10^{-4} and a minimum value of 10^{-5} . The batch size was set to 32, and the maximum number of training epochs, T_{max} , was set to 100 to ensure training stability and eventual convergence.

4.2 Experimental comparison

4.2.1 Kvasir public dataset

To validate the outstanding performance of our proposed GLI-Net on the Kvasir public dataset, we compared it with current state-of-the-art models. Specifically, as shown in Table 1, compared to ConvNeXt-B, ViT-B/16, ViT-B/32, and Swin-B models, our method achieved improvements of 13.31%, 10.25%, 14.81%, and 9.03% in Acc, respectively; 13.55%, 10.82%, 15.36%, and 9.39% in F1 score; 13.76%, 10.94%, 15.43%, and 9.51% in P; and 14.31%, 11.57%, 16.07%, and 10.10% in R. Moreover, compared to the HiFuse model, GLI-Net improved accuracy, F1 score, precision, and recall by 3.29%, 3.54%, 3.70%, and 4.28%, respectively. These results demonstrate that GLI-Net is more effective in capturing and integrating both global and local features, significantly enhancing the accuracy and robustness of medical endoscopic image classification, and showcasing its superior performance in complex medical image analysis tasks.

To further demonstrate the superior performance of GLI-Net on the Kvasir dataset, we applied the Grad-CAM method to visualize the model's final layer, generating heatmaps that reflect the regions of the lesion the model focuses on. The specific details are shown in Figure 6. Compared to models such as ConvNeXt-B, ViT-B/16, ViT-B/32, Swin-B, and HiFuse, GLI-Net's heatmaps

TABLE 1 Comparative results of different methods on the Kvasir public dataset.

Method	Accuracy ↑	F1 score ↑	Precision ↑	Recall ↑
ConvNeXt-B	74.6	74.61	74.78	74.64
ViT-B/16	76.1	75.94	76.49	76.23
ViT-B/32	73.8	73.5	74.24	73.72
Swin-B	77.3	77.29	77.74	77.44
HiFuse	84.35	84.41	84.5	84.48
GLI-Net (Ours)	87.43	87.68	88.26	89.12

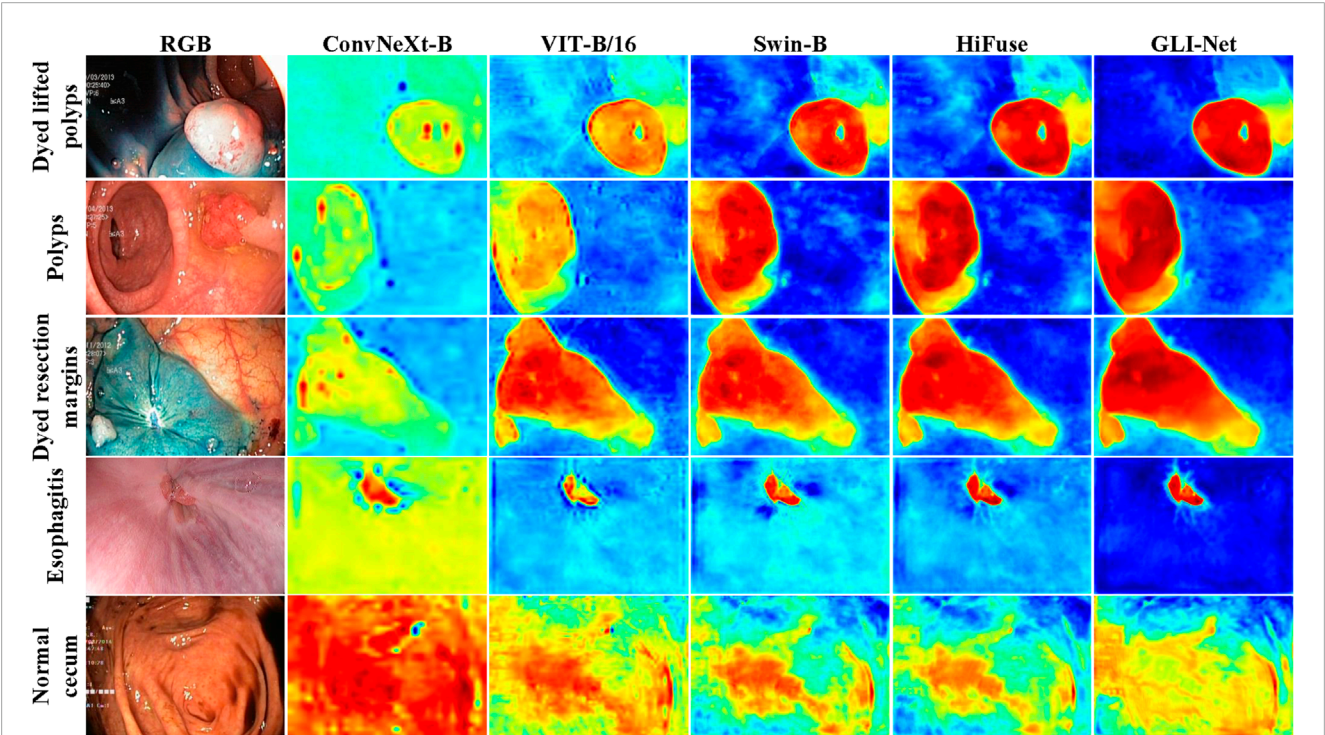


FIGURE 6 Heatmap visualization results on the Kvasir public dataset.

show higher focus and coverage of the lesion regions, allowing for more accurate localization of lesions in endoscopic images. While other models can recognize some lesion areas, they exhibit discrepancies in precise localization and coverage. For example, ConvNeXt-B and ViT-B/32 show relatively blurred recognition, while Swin-B and HiFuse incorrectly label many non-lesion areas. GLI-Net, by effectively covering lesion regions and minimizing background interference, demonstrates significant advantages in feature extraction and region localization. These visualization results prove GLI-Net's efficiency and reliability in medical image classification tasks.

4.2.2 Hyper-Kvasir public dataset

To validate the outstanding performance of our proposed GLI-Net on the Hyper-Kvasir public dataset, we compared it with current

state-of-the-art models. The specific results are shown in Table 2. Compared to ConvNeXt-B, ViT-B/16, ViT-B/32, Swin-B, and HiFuse models, GLI-Net achieved improvements of 13.31%, 10.25%, 14.81%, 9.03%, and 3.29% in Acc, respectively; 13.55%, 10.82%, 15.36%, 9.39%, and 3.54% in F1 score; 13.76%, 10.94%, 15.43%, 9.51%, and 3.70% in P; and 14.31%, 11.57%, 16.07%, 10.10%, and 4.28% in R. These significant performance improvements indicate that GLI-Net is more effective in capturing and integrating both global and local features, significantly enhancing the accuracy and robustness of medical endoscopic image classification, and showcasing its superior performance in complex medical image analysis tasks.

To demonstrate the superior performance of GLI-Net on the Hyper-Kvasir dataset, we used the Grad-CAM method to generate heatmaps that visualize the lesion regions the model focuses on.

TABLE 2 Comparative results of different methods on the Hyper-Kvasir public dataset.

Method	Accuracy ↑	F1 score ↑	Precision ↑	Recall ↑
ConvNeXt-B	72.53	72.61	72.8	72.71
VIT-B/16	75.59	75.34	75.62	75.45
VIT-B/32	71.03	70.8	71.13	70.95
Swin-B	76.81	76.77	77.05	76.92
HiFuse	82.55	82.62	82.86	82.74
GLI-Net (Ours)	85.84	86.16	86.56	87.02

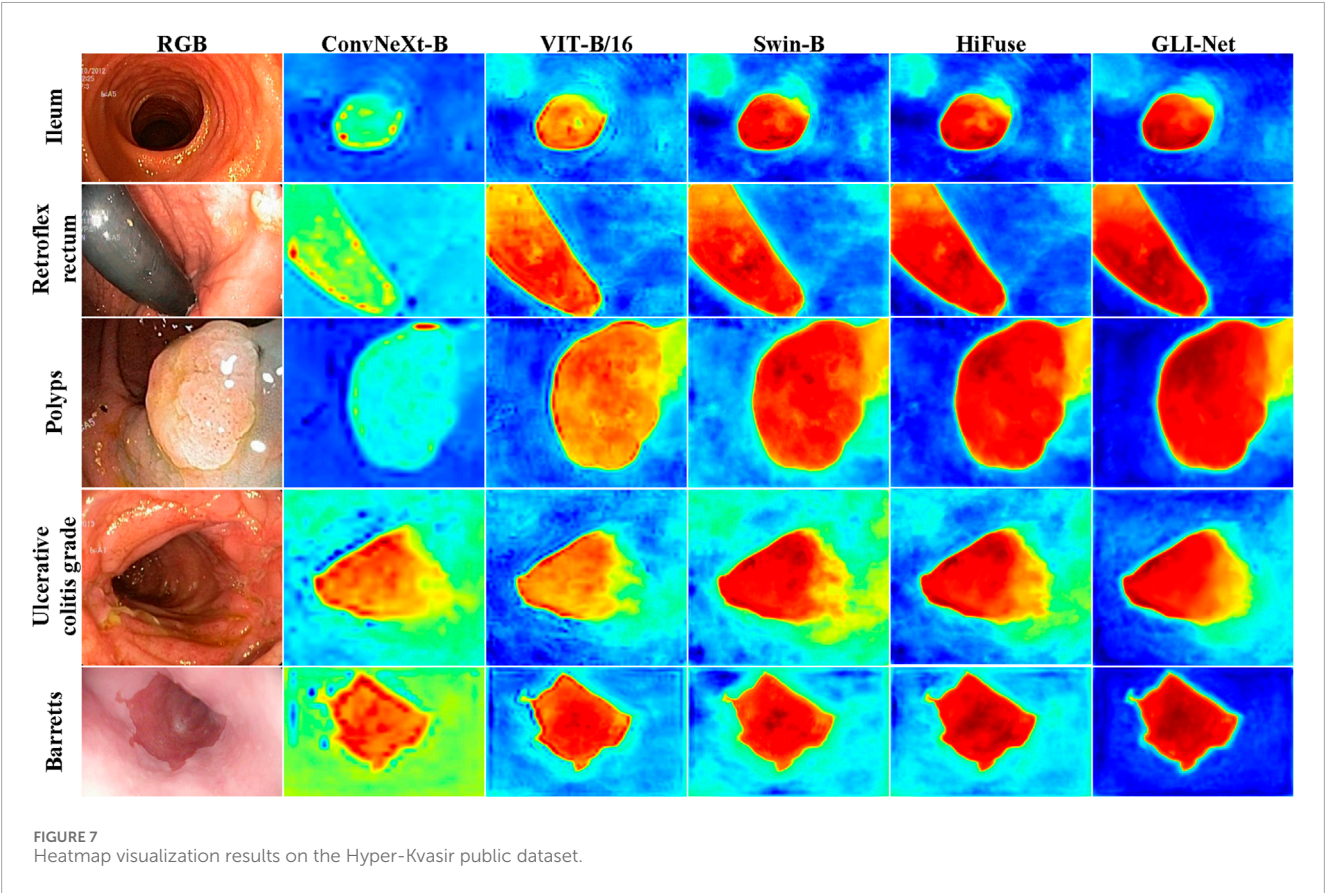


TABLE 3 Ablation study results of GLI-Net on the Kvasir public dataset.

Method	GB	LB	LEM	AFE	Acc \ %	F1 \ %	Prec \ %	Recall \ %
Case. S_1	X				81.12	81.23	81.52	81.04
Case. S_2		X			82.53	82.73	82.97	82.62
Case. S_3			X		83.86	83.9	84.16	83.83
Case. S_4				X	84.57	84.77	84.93	84.66
GLI-Net					87.43	87.68	88.26	89.12

The specific details are shown in Figure 7. Compared to models such as ConvNeXt-B, ViT-B/16, ViT-B/32, Swin-B, and HiFuse, GLI-Net’s heatmaps exhibit higher focus and coverage of the lesion areas, allowing for more accurate localization of lesions in endoscopic images. While other models can identify some lesion regions, they show discrepancies in localization and coverage. For example, ConvNeXt-B and ViT-B/32 exhibit relatively blurred recognition, while Swin-B and HiFuse incorrectly label non-lesion regions. GLI-Net, through more precise lesion region coverage and reduced background interference, demonstrates its advantages in feature extraction and region localization. These results show that GLI-Net can more effectively integrate global and local features, significantly improving the accuracy and robustness of medical endoscopic image classification, and proving its efficiency and reliability in real-world applications.

4.3 Ablation study

4.3.1 Ablation study of GLI-Net

To evaluate the performance of GLI-Net on the Kvasir dataset, we conducted ablation experiments by sequentially removing the GB, LB, LEM, and AFE modules from the model. The results

are shown in Table 3. The inclusion of each module significantly improved the model’s performance. The baseline model with the GB module removed achieved an accuracy of 81.12%. After removing the LB module, the accuracy increased to 82.53%, and further removal of the LEM module raised the accuracy to 83.86%. When the AFE module was removed, the accuracy reached 84.57%. Finally, the complete GLI-Net model achieved an accuracy of 87.43%, which is a 2.86% improvement over the model without the AFE module. In addition, GLI-Net also performed better in other evaluation metrics such as F1 score, precision, and recall, with improvements of 6.45%, 4.45%, 3.78%, 8.08%, 6.50%, and 5.29%, respectively. These experimental results demonstrate that the individual modules of GLI-Net play a critical role in enhancing feature extraction, feature fusion, and optimizing representation, significantly improving the accuracy and robustness of medical endoscopic image classification, and proving its superior performance in complex medical image analysis tasks.

To verify the role of each module in GLI-Net, we conducted ablation experiments on the Kvasir dataset by sequentially removing the modules and used the Grad-CAM method to visualize the lesion regions the model focuses on under different configurations. The specific details are shown in Figure 8. The experimental results

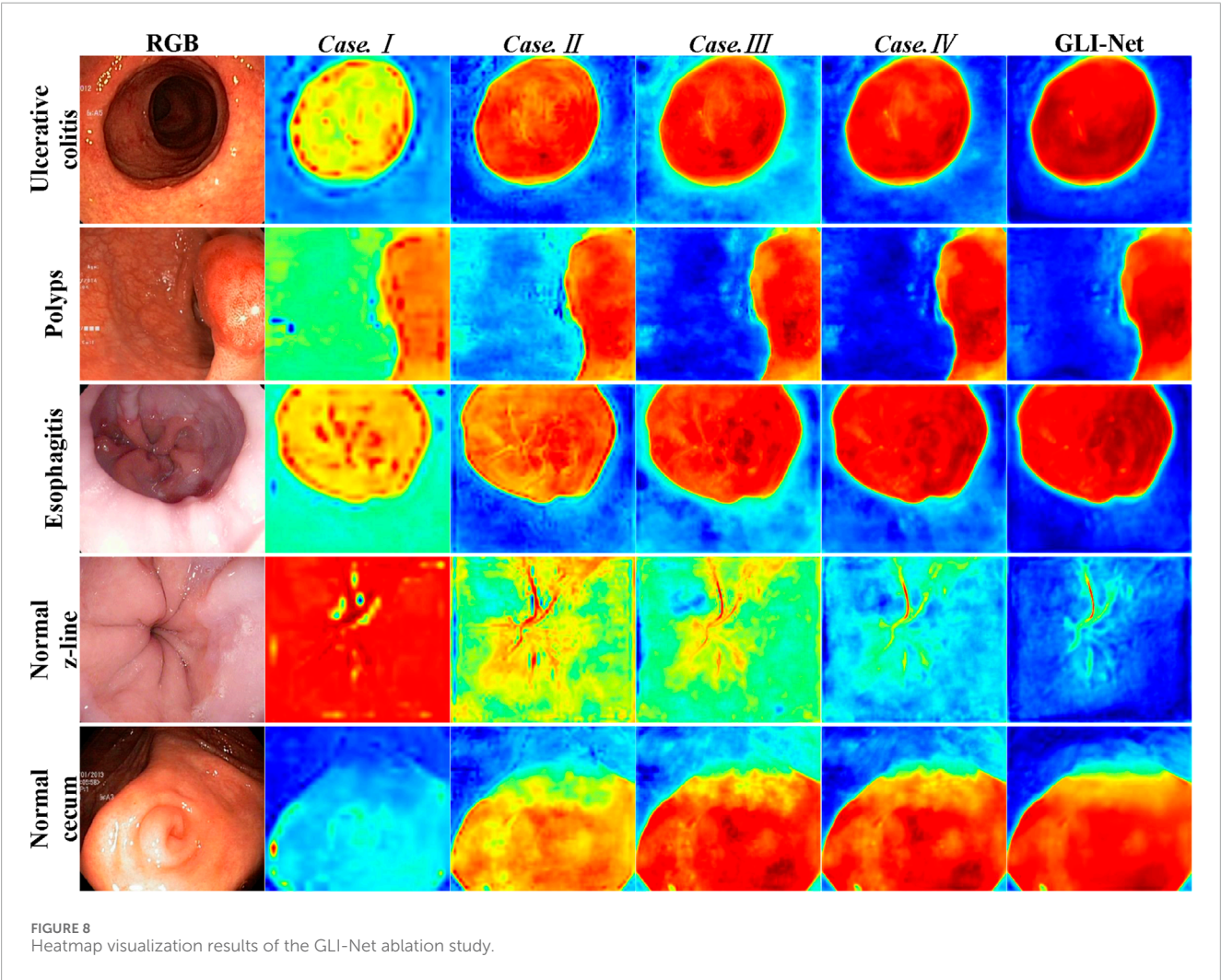


TABLE 4 Ablation study results of different losses on the Kvasir public dataset.

Method	\mathcal{L}_{det}	\mathcal{L}_{cls}	λ	Accuracy \uparrow	F1 score \uparrow	Precision \uparrow	Recall \uparrow
Case.S ₁	✓			72.32	71.51	72.35	72.13
Case.S ₂		✓		75.24	74.86	75.21	74.94
Case.S ₃	✓	✓		80.17	79.52	80.17	79.88
GLI-Net	✓	✓	✓	87.43	87.68	88.26	89.12

indicate that, after removing the Global Branch module (GB) (Case.S₁), the model's focus on lesion areas significantly decreased, revealing a deficiency in capturing global features. Removing the Local Branch module (LB) (Case.S₂) weakened the ability to extract detailed features, resulting in blurred lesion regions. After removing the Information Exchange Module (LEM) (Case.S₃), although the model could still detect lesion regions, the insufficient fusion of global and local features affected comprehensive coverage of the lesion areas. When the Adaptive Feature Fusion and Enhancement Module (AFE) (Case.S₄) was removed, although the focus on the lesion areas increased, feature expression and region optimization were insufficient, leading to residual background interference. In contrast, the complete GLI-Net model, through the synergistic action of all modules, accurately localized the lesion regions, significantly improving the model's accuracy and robustness in medical endoscopic image classification. The superior performance of GLI-Net can be attributed to the effective integration of global and local features, along with the bidirectional information exchange facilitated by the LEM module. The Global Branch (GB) extracts high-level semantic features that provide a broad context for the lesions, while the Local Branch (LB) captures fine-grained details of the lesions. The Information Exchange Module (LEM) allows for mutual enhancement of these features, ensuring that both global and local features are used in a complementary manner. This interaction mitigates the issues caused by intra-class variation and subtle inter-class differences, which are common in endoscopic images, and thus leads to more accurate and robust classification results. These results demonstrate the crucial roles of each module in feature extraction, region localization, and feature fusion.

4.3.2 Ablation study of loss function

To evaluate the contribution of each component of the loss function in GLI-Net, we conducted ablation experiments on the Kvasir dataset by sequentially removing the lesion region detection loss, classification loss, and weight coefficient. The results are shown in Table 4. When only the lesion region detection loss was used (Case.S₁), the accuracy was 72.32%. After adding the classification loss (Case.S₂), the accuracy increased to 75.24%. When both the lesion region detection loss and classification loss were used together (Case.S₃), the accuracy further improved to 80.17%. Finally, the complete GLI-Net model achieved an accuracy of 87.43%, 7.26% improvement over Case.S₃, highlighting the important role of the weight coefficient λ in balancing the loss function. Additionally, GLI-Net showed significant improvements in F1 score, precision, and recall, with increases of 6.78%, 4.81%, and 4.14%, respectively, compared to Case.S₃. These results indicate that the effective

combination of the lesion region detection loss and classification loss, along with the proper setting of the weight coefficient, significantly enhances the model's performance, confirming the key role of the loss function design in GLI-Net.

5 Conclusion

This paper presents GLI-Net, a novel network for medical endoscopic image classification, designed to enhance classification performance by effectively integrating both global and local features. GLI-Net utilizes a hierarchical multi-module architecture that includes a global branch module (GB), a local branch module (LB), an information exchange module (LEM), and an adaptive feature fusion and enhancement module (AFE) to facilitate comprehensive feature extraction and optimization. Evaluation on the Kvasir and Hyper-Kvasir public datasets showed that GLI-Net outperforms state-of-the-art models, including ConvNeXt-B, ViT-B/16, ViT-B/32, Swin-B, and HiFuse, across key metrics such as accuracy, F1 score, precision, and recall. Specifically, GLI-Net achieved accuracies of 87.43% and 85.84% on the Kvasir and Hyper-Kvasir datasets, respectively, surpassing the second-best models by 2.86% and 2.29%. Ablation studies confirmed the significant contribution of each module to the overall performance, as the removal of any module caused a notable performance decline, underscoring their synergistic interaction. Additionally, Grad-CAM visualization highlighted GLI-Net's improved ability to accurately localize lesion areas, with better focus and coverage compared to other models, effectively reducing interference from background and non-lesion regions. These results demonstrate GLI-Net's substantial advantages in feature extraction and region localization, leading to enhanced accuracy and robustness in medical endoscopic image classification.

Data availability statement

The original contributions presented in the study are included in the article/supplementary material, further inquiries can be directed to the corresponding author.

Author contributions

YZ: Writing – original draft, Writing–review and editing, Conceptualization, Data curation, Investigation, Methodology. MZ: Writing – original draft, Writing – review and editing,

Conceptualization, Investigation. WC: Writing – original draft, Writing – review and editing, Data curation, Methodology. XW: Writing – original draft, Writing – review and editing, Formal Analysis. QS: Writing – original draft, Writing – review and editing, Data curation.

Funding

The author(s) declare that no financial support was received for the research and/or publication of this article.

Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships

that could be construed as a potential conflict of interest.

Generative AI statement

The authors declare that no Generative AI was used in the creation of this manuscript.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

References

- Karargyris A, Bourbakis N. Detection of small bowel polyps and ulcers in wireless capsule endoscopy videos. *IEEE Trans Biomed Eng* (2011) 58:2777–86. doi:10.1109/tbme.2011.2155064
- Mesejo P, Pizarro D, Abergel A, Rouquette O, Beorchia S, Poincloux L, et al. Computer-aided classification of gastrointestinal lesions in regular colonoscopy. *IEEE Trans Med Imaging* (2016) 35:2051–63. doi:10.1109/tmi.2016.2547947
- Charfi S, Ansari ME. Computer-aided diagnosis system for colon abnormalities detection in wireless capsule endoscopy images. *Multimedia Tools Appl* (2018) 77:4047–64. doi:10.1007/s11042-017-4555-7
- LeCun Y, Bengio Y, Hinton G. Deep learning. *nature* (2015) 521:436–44. doi:10.1038/nature14539
- Simonyan K. *Very deep convolutional networks for large-scale image recognition* (2014). arXiv preprint arXiv:1409.1556.
- He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*; 27–30 June 2016; Las Vegas, NV, USA (2016). p. 770–8.
- Thambawita V, Strümke I, Hicks SA, Halvorsen P, Parasa S, Riegler MA. Impact of image resolution on deep learning performance in endoscopy image classification: an experimental study using a large dataset of endoscopic images. *Diagnostics* (2021) 11:2183. doi:10.3390/diagnostics11122183
- Mukhtov D, Rakhmonova M, Muksimova S, Cho Y-I. Endoscopic image classification based on explainable deep learning. *Sensors* (2023) 23:3176. doi:10.3390/s23063176
- Yue G, Wei P, Liu Y, Luo Y, Du J, Wang T. Automated endoscopic image classification via deep neural network with class imbalance loss. *IEEE Trans Instrumentation Meas* (2023) 72:1–11. doi:10.1109/tim.2023.3264047
- Bolhasani H, Jassbi SJ, Sharifi A. Dla-e: a deep learning accelerator for endoscopic images classification. *J Big Data* (2023) 10:76. doi:10.1186/s40537-023-00775-8
- Li B, Meng MQ-H. Wireless capsule endoscopy images enhancement via adaptive contrast diffusion. *J Vis Commun Image Representation* (2012) 23:222–8. doi:10.1016/j.jvcir.2011.10.002
- Souaidi M, Ansari ME. Multi-scale analysis of ulcer disease detection from wce images. *IET Image Process* (2019) 13:2233–44. doi:10.1049/iet-ipr.2019.0415
- Zhang R, Zheng Y, Mak TWC, Yu R, Wong SH, Lau JY, et al. Automatic detection and classification of colorectal polyps by transferring low-level cnn features from nonmedical domain. *IEEE J Biomed Health Inform* (2016) 21:41–7. doi:10.1109/jbhi.2016.2635662
- Shin Y, Balasingham I. Comparison of hand-craft feature based svm and cnn based deep learning framework for automatic polyp classification. In: *2017 39th annual international conference of the IEEE engineering in medicine and biology society (EMBC)*. IEEE (2017). p. 3277–80.
- Zhao Q, Yang W, Liao Q. Adasan: adaptive cosine similarity self-attention network for gastrointestinal endoscopy image classification. In: *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*; 13–16 April 2021; Nice, France. IEEE (2021). p. 1855–9.
- Zhu R, Zhang R, Xue D. Lesion detection of endoscopy images based on convolutional neural network features. In: *2015 8th International Congress on Image and Signal Processing (CISP)*; 14–16 October 2015; Shenyang, China. IEEE (2015). p. 372–6.
- Chen J, Zou Y, Wang Y. Wireless capsule endoscopy video summarization: a learning approach based on siamese neural network and support vector machine. In: *2016 23rd International Conference on Pattern Recognition (ICPR)*; 04–08 December 2016; Cancun, Mexico. IEEE (2016). p. 1303–8.
- Jeon Y, Cho E, Moon S, Chae S-H, Jo HY, Kim TO, et al. Deep convolutional neural network-based automated lesion detection in wireless capsule endoscopy. In: *International forum on medical imaging in asia 2019*, 11050. SPIE (2019). p. 64–296. doi:10.1117/12.2522159
- Guo X, Yuan Y. Triple anet: adaptive abnormal-aware attention network for wce image classification. In: *Medical image computing and computer assisted intervention—MICCAI 2019: 22nd international conference, shenzhen, China, october 13–17, 2019, proceedings, Part I* 22. Springer (2019). p. 293–301.
- Cao J, Yao J, Zhang Z, Cheng S, Li S, Zhu J, et al. Efag-cnn: effectively fused attention guided convolutional neural network for wce image classification. In: *2021 IEEE 10th Data Driven Control and Learning Systems Conference (DDCLS)*; 14–16 May 2021; Suzhou, China. IEEE (2021). p. 66–71.



OPEN ACCESS

EDITED BY

Yu Liu,
Hefei University of Technology, China

REVIEWED BY

Guohua Lv,
Qilu University of Technology, China
Min Li,
Xinjiang University, China

*CORRESPONDENCE

Yukui Che,
✉ 454983185@qq.com

RECEIVED 25 March 2025

ACCEPTED 06 May 2025

PUBLISHED 22 May 2025

CITATION

Wang K, Hu D, Cheng Y, Che Y, Li Y, Jiang Z,
Chen F and Li W (2025) Infrared and visible
image fusion driven by multimodal large
language models.
Front. Phys. 13:1599937.
doi: 10.3389/fphy.2025.1599937

COPYRIGHT

© 2025 Wang, Hu, Cheng, Che, Li, Jiang,
Chen and Li. This is an open-access article
distributed under the terms of the [Creative
Commons Attribution License \(CC BY\)](#). The
use, distribution or reproduction in other
forums is permitted, provided the original
author(s) and the copyright owner(s) are
credited and that the original publication in
this journal is cited, in accordance with
accepted academic practice. No use,
distribution or reproduction is permitted
which does not comply with these terms.

Infrared and visible image fusion driven by multimodal large language models

Ke Wang, Dengshu Hu, Yuan Cheng, Yukui Che*, Yuelin Li,
Zhiwei Jiang, Fengxian Chen and Wenjuan Li

Qujing Power Supply Bureau, Yunnan Power Grid Co., Ltd., Kunming, China

Introduction: Existing image fusion methods primarily focus on obtaining high-quality features from source images to enhance the quality of the fused image, often overlooking the impact of improved image quality on downstream task performance.

Methods: To address this issue, this paper proposes a novel infrared and visible image fusion approach driven by multimodal large language models, aiming to improve the performance of pedestrian detection tasks. The proposed method fully considers how enhancing image quality can benefit pedestrian detection. By leveraging a multimodal large language model, we analyze the fused images based on user-provided questions related to improving pedestrian detection performance and generate suggestions for enhancing image quality. To better incorporate these suggestions, we design a Text-Driven Feature Harmonization (Text-DFH) module. Text-DFH refines the features produced by the fusion network according to the recommendations from the multimodal large language model, enabling the fused image to better meet the needs of pedestrian detection tasks.

Results: Compared with existing methods, the key advantage of our approach lies in utilizing the strong semantic understanding and scene analysis capabilities of multimodal large language models to provide precise guidance for improving fused image quality. As a result, our method enhances image quality while maintaining strong performance in pedestrian detection. Extensive qualitative and quantitative experiments on multiple public datasets validate the effectiveness and superiority of the proposed method.

Discussion: In addition to its effectiveness in infrared and visible image fusion, the method also demonstrates promising application potential in the field of nuclear medical imaging.

KEYWORDS

infrared and visible image fusion, pedestrian detection, multimodal large language models, text-guided, model fine-tuning

1 Introduction

Multimodal sensor technology has facilitated the application of multimodal images across various fields. Among them, infrared and visible images have been widely used in diverse tasks due to the complementary nature of the information they contain. Specifically, infrared images provide thermal radiation information of objects and are

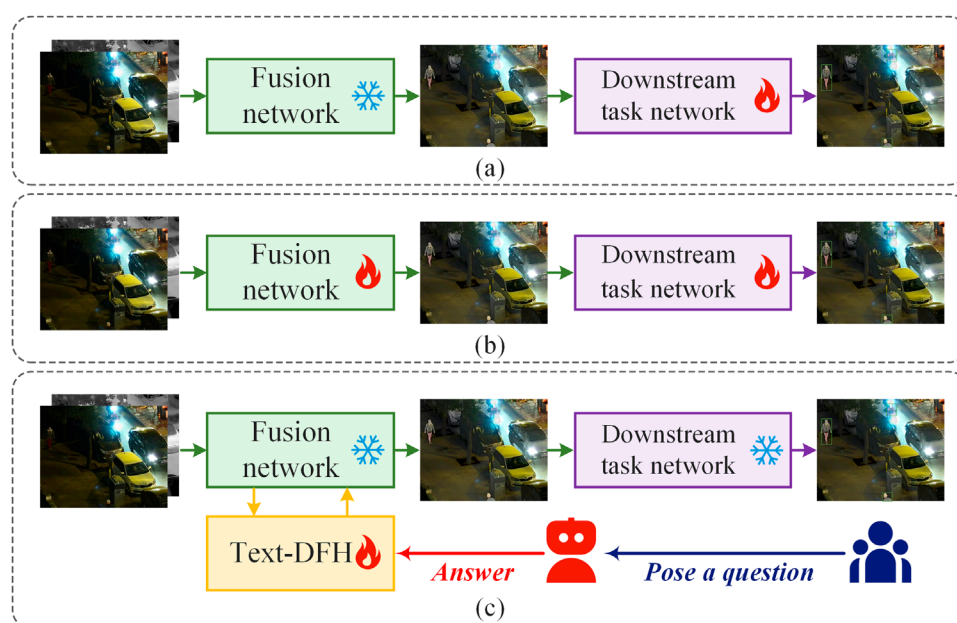


FIGURE 1
Comparison of different joint training strategies for image fusion and downstream tasks.

not affected by lighting conditions, but they lack detailed textures. In contrast, visible images capture rich texture details of the scene but are highly sensitive to lighting variations. Therefore, numerous methods [1–7] have focused on fusing infrared and visible images, aiming to integrate the complementary information from both modalities into a single, more informative fused image. This facilitates better decision-making and judgment in downstream tasks such as object detection [8–10] and semantic segmentation [11–14].

Current approaches that jointly train infrared-visible image fusion with downstream tasks can be broadly categorized into two types: independent optimization and joint optimization. Independent optimization methods first train a fusion network for infrared and visible images and then use the resulting fused images to train a downstream task network, as shown in Figure 1a. Consequently, most independent optimization methods focus on improving fusion quality, for example, by designing new network architectures [15–19] or introducing specific constraints [20–23]. However, such approaches neglect the potential guidance from downstream tasks and fail to establish a deep connection between fusion and task performance, often leading to suboptimal results. Simply chaining the fusion and downstream networks makes it difficult for the fused image to specifically cater to the downstream task's requirements. On the other hand, joint optimization methods use the downstream task network as a constraint to train the image fusion network, thereby forcing it to produce fused images that meet task-specific needs [24–28], as illustrated in Figure 1b. Nevertheless, the effectiveness of directly using high-level vision task supervision to guide fusion remains limited.

Recently, Multimodal Large Language Models (MLLMs) have gained popularity due to their strong capability in modeling data across different modalities, such as images and text. For instance,

Text-IF [29] and TeRF [30] leverage large models to encode user instructions and guide various types of fusion tasks. However, these methods do not consider the possibility of using large language models to feed back the specific needs of high-level vision tasks to the image fusion process, which could further improve the quality of fused images.

To address this challenge, we propose a novel infrared and visible image fusion method driven by a Multimodal Large Language Model, aiming to simultaneously enhance fusion quality and pedestrian detection accuracy, as shown in Figure 1c. By leveraging the deep semantic understanding and scene analysis capabilities of MLLMs, we provide precise guidance for improving fused image quality while ensuring better pedestrian detection performance. Specifically, our method analyzes the fused images based on user-provided questions related to pedestrian detection, then generates optimization suggestions using feedback from the language model. To fully utilize these suggestions, we design a Text-Driven Feature Harmonization (Text-DFH) module, which refines the fusion network's output features under the guidance of the MLLM, allowing the fused images to better meet the demands of pedestrian detection.

In summary, the main contributions of this paper are as follows:

- (1) We are the first to leverage Multimodal Large Language Models to provide feedback on the quality of fused images based on the specific requirements of downstream tasks, thus further improving infrared and visible image fusion.
- (2) We propose an effective Text-Driven Feature Harmonization (Text-DFH) module that enables text-based guidance to assist in enhancing image quality.
- (3) Our proposed method achieves excellent performance in infrared and visible image fusion, nuclear medical imaging, and pedestrian detection across multiple datasets.

The remainder of this paper is organized as follows. [Section 2](#) provides a brief overview of related work on multimodal large language models, infrared and visible image fusion, and pedestrian detection. [Section 3](#) presents our proposed method in detail. [Section 4](#) discusses the experimental results and analysis. [Section 5](#) concludes the paper.

2 Related work

In this section, we first briefly introduce multimodal large language models, and then review existing infrared and visible image fusion methods.

2.1 Multimodal large language models

With the advent of the multimodal data fusion era, the capability of unimodal systems is no longer sufficient to handle complex real-world tasks. As a result, multimodal large language models (MLLMs) have been proposed to integrate information from multiple data sources, enabling more comprehensive and accurate representations. These models have demonstrated significant practical value across various domains, including natural language processing, vision tasks, and audio tasks. In the visual domain, MLLMs enhance the performance of tasks such as image classification, object detection, and image captioning by combining textual descriptions with visual instructions. For example, GPT-4V [31] and Gemini [32] integrate image content with natural language descriptions to produce more vivid and accurate annotations. NExT-GPT [33] and Sora [34] are at the forefront of multimodal video generation, producing rich and realistic content by learning from multimodal data. Moreover, VideoChat [35] and Video-LLaVA [36] demonstrate excellent capabilities in analyzing and understanding video content in intelligent video understanding scenarios.

In the field of image fusion, Text-IF [29] and MGFusion [37] uses CLIP [38] to encode user requirement texts, guiding the model to fuse images. TeRF [30] utilizes LLaMA [39] to encode user instruction texts and generate prompts for guiding image fusion across different tasks. Although these methods employ MLLMs to tackle some challenges in image fusion, they do not consider the specific requirements of high-level downstream visual tasks for image fusion quality, which limits the application of infrared and visible image fusion in such tasks.

2.2 Infrared and visible image fusion

Conventional infrared and visible image fusion methods mainly focus on designing sophisticated feature extraction networks and fusion strategies to ensure the quality of the fused results. From the perspective of network design, these methods can be broadly categorized into CNN-based methods, CNN-Transformer hybrid methods, and GAN-based methods. CNN-based methods [40–45] typically apply convolution, activation, and pooling operations to extract features from the input images, then fuse and reconstruct the final result using the extracted features. However, since CNNs can only perceive local features within a limited receptive field,

they struggle to capture long-range contextual information, limiting their representational capacity. In contrast, Transformers [46] are better at modeling long-range dependencies and are more suited for capturing global features in images. ViT [47] was the first to introduce Transformer architectures into computer vision, achieving promising results. Subsequently, to combine the respective strengths of CNNs and Transformers, hybrid methods have gained increasing attention in the image fusion domain. For instance, CGTF [48], SwinFusion [16], YDTR [17], and DATFuse [49] insert Transformer layers after CNN layers to jointly leverage local and global feature extraction. CDDFuse [50] and EMMA [51] adopt dual-branch architectures combining CNNs and Transformers to simultaneously extract features from the input images and integrate them for fusion.

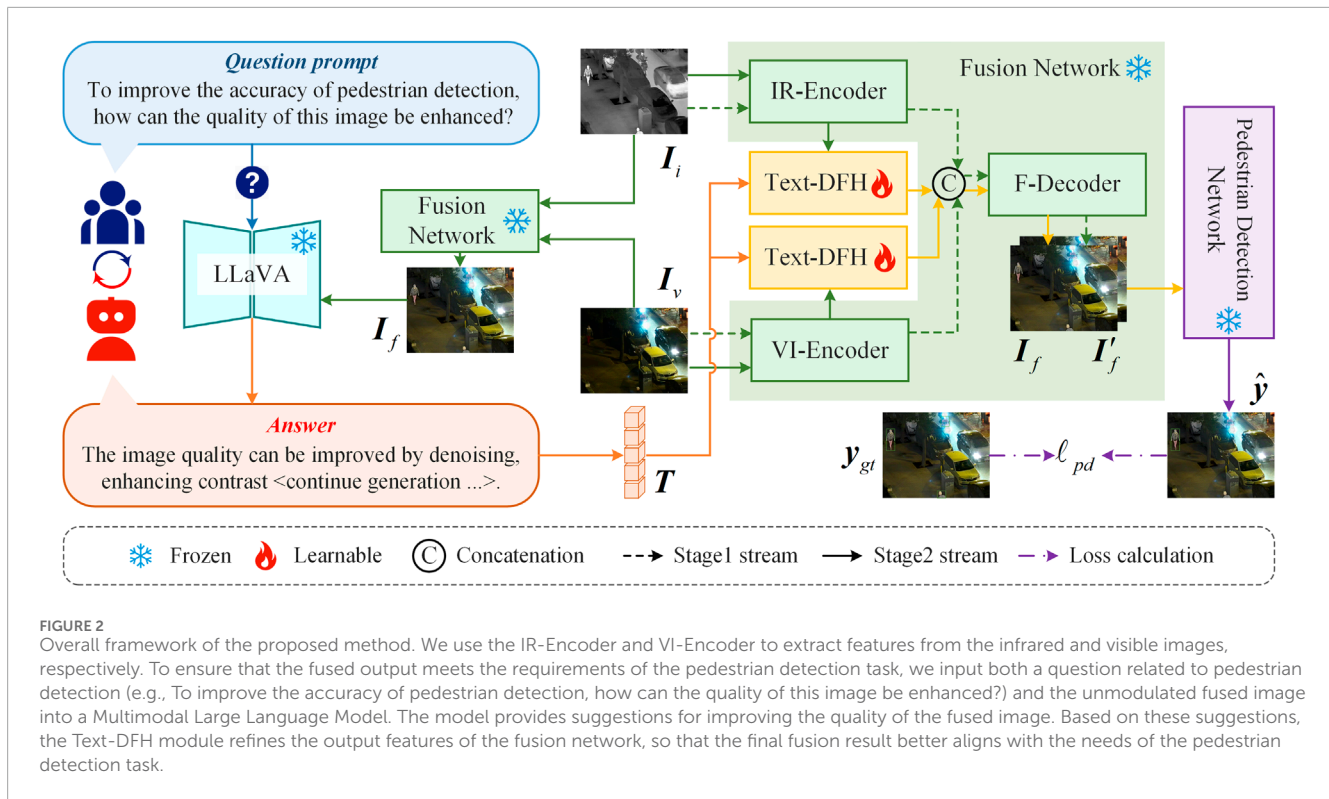
GAN-based methods enhance the model's feature extraction capabilities by introducing adversarial learning between generators and discriminators. Depending on the number of discriminators used, these methods can be classified into single-discriminator and dual-discriminator approaches. Single-discriminator methods [2, 52] tend to favor one modality over the other, potentially leading to information loss and reduced visual quality of the fusion results. To address this, dual-discriminator methods [53–56] are proposed to preserve important features from both source images simultaneously.

However, all of these methods primarily focus on designing effective feature extraction networks to produce high-quality fusion features and images. They overlook how fusion quality impacts downstream task performance, and fail to consider the potential feedback from downstream tasks that could help guide fusion more effectively.

2.3 Pedestrian detection

Pedestrian detection is a fundamental problem in computer vision with a wide range of applications. Cascade R-CNN [57] extends R-CNN [58] into a multi-stage framework, improving the ability to filter hard negative samples. Faster R-CNN [59] introduces a Region Proposal Network (RPN) that shares convolutional features with the detection network, making region proposals nearly cost-free. YOLO [60] reformulates object detection as a regression problem, allowing real-time inference directly on images through a convolutional neural network. SSD [61] uses multi-scale feature maps and predefined anchors for pedestrian detection, addressing YOLO's limitations in detecting small objects. DETR [62] adopts a Transformer-based encoder-decoder architecture for object detection. BAS Wu et al. [63] learns to represent the whole foreground region by leveraging foreground guidance and domain constraints. CREAM [64] proposes a clustering-based method to enhance activation within target regions. Group R-CNN [65] builds instance groups to perform pedestrian detection from point annotations.

However, most pedestrian detection methods are designed for unimodal images, which often leads to degraded detection performance due to incomplete scene information. In this work, we perform pedestrian detection on fused infrared and visible images, and incorporate task-specific prompts generated by large language



models. This not only improves the quality of the fused images but also enhances pedestrian detection performance.

3 Methods

3.1 Overview

As shown in Figure 2, the proposed method consists of two training stages. The first stage is dedicated to training the Fusion Network, enabling it to perform basic infrared and visible image fusion. In the second stage, the parameters of the pretrained fusion network are frozen, and a Text-Driven Feature Harmonization (Text-DFH) module is trained to refine the fusion results to better align with the requirements of pedestrian detection. The fusion network is composed of three main components: an Infrared Image Feature Encoder (IR-Encoder), a Visible Image Feature Encoder (VI-Encoder), and a Fusion Feature Decoder (F-Decoder). The IR/VI-Encoders are responsible for extracting features from the input infrared and visible images, respectively, while the F-Decoder reconstructs the fused image based on the combined features. The Text-DFH module adjusts the features extracted by the IR/VI-Encoders based on responses from a Multimodal Large Language Model (MLLM), ensuring that the resulting fused image better satisfies the needs of pedestrian detection. In this work, we adopt LLaVA [66] as the MLLM. LLaVA analyzes the unmodulated fused image and generates suggestions in response to user queries related to pedestrian detection tasks (e.g., To improve the accuracy of pedestrian detection, how can the quality of this image be enhanced?). More text examples of LLaVA answers are shown in Figure 3.

3.2 Feature extraction and fusion

In the first training stage, we train the fusion network to perform the basic task of infrared and visible image fusion. The fusion network primarily consists of three components: the IR-Encoder, VI-Encoder, and F-Decoder. Each of the IR-Encoder, VI-Encoder, and F-Decoder is composed of three feature extraction layers. Each layer is constructed by stacking a convolutional layer (kernel size = 3×3 , stride = 1), a Batch Normalization layer, and a LeakyReLU activation function. It is worth noting that the LeakyReLU activation function in the final feature extraction layer of the F-Decoder is replaced with a Tanh activation function to facilitate image reconstruction. We input the infrared image I_i and the visible image I_v into the IR-Encoder and VI-Encoder, respectively, to extract features F_i and F_v . To reconstruct the fused image, we concatenate F_i and F_v along the channel dimension and feed the result into the F-Decoder, which generates the final fused image I_f .

To encourage the fused image to retain as much scene information from the source images as possible, we introduce an intensity loss ℓ_{in} and an edge loss ℓ_{ed} , which together form the fusion loss ℓ_f :

$$\ell_f = \ell_{in} + \epsilon \ell_{ed}, \quad (1)$$

Here, ϵ denotes a hyperparameter used to balance the contribution of each sub-loss term. The intensity loss ℓ_{in} is defined as:

$$\ell_{in} = \frac{1}{HW} (\|I_f - I_i\|_1 + \|I_f - I_v\|_1), \quad (2)$$

The edge loss ℓ_{ed} is defined as:

$$\ell_{ed} = \frac{1}{HW} (\|\nabla I_f - \nabla I_i\|_1 + \|\nabla I_f - \nabla I_v\|_1), \quad (3)$$

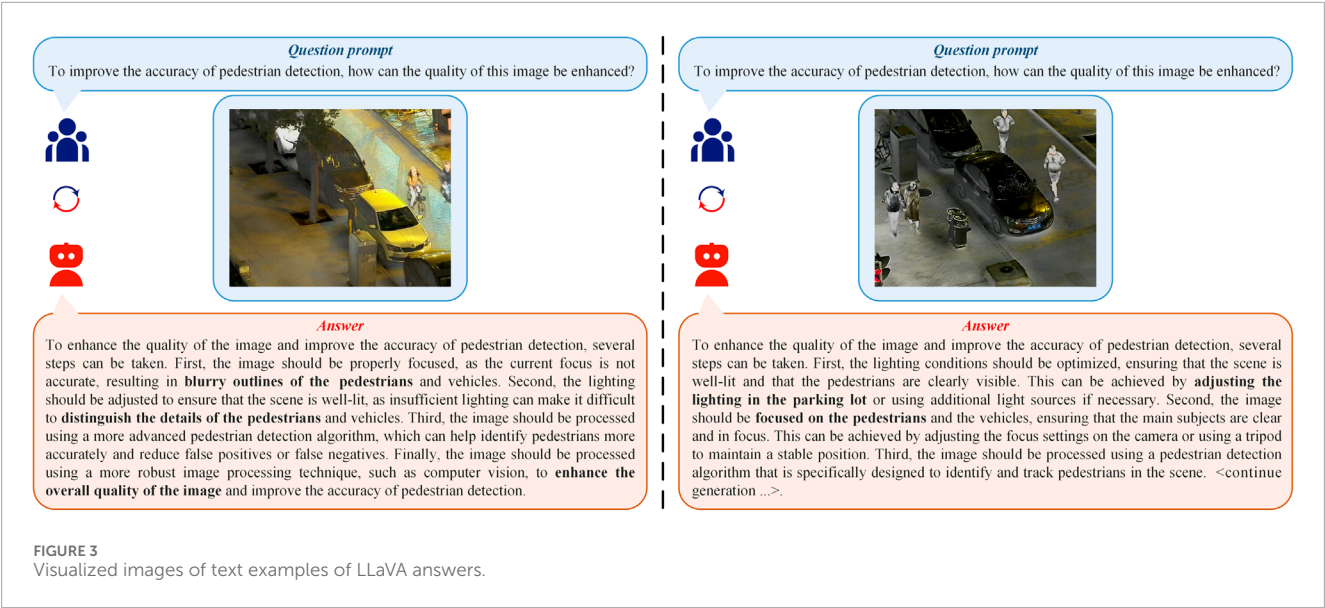


FIGURE 3
Visualized images of text examples of LLaVA answers.

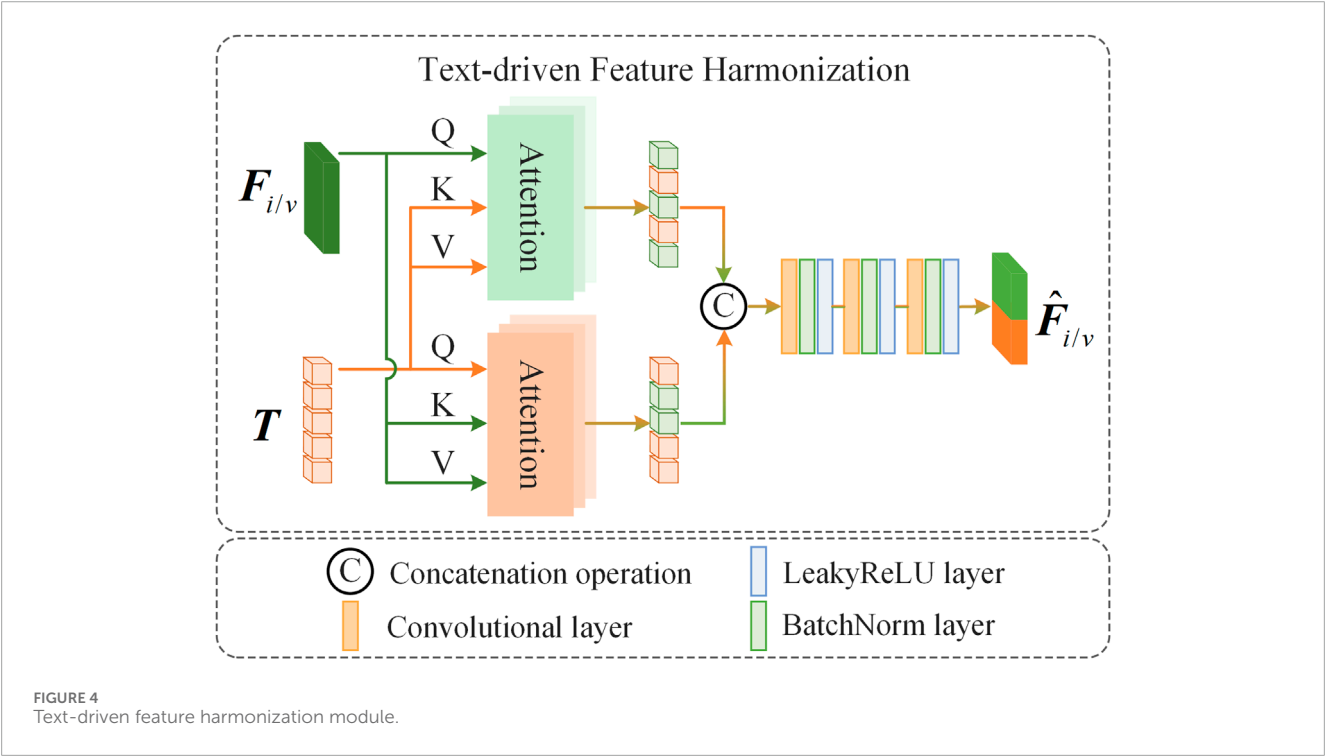


FIGURE 4
Text-driven feature harmonization module.

Here, H and W denote the height and width of the fused image, respectively; $\|\cdot\|_1$ represents the l1-norm, and ∇ denotes the Sobel edge extraction operator.

3.3 Text-driven feature harmonization

In the second training stage, we freeze the parameters of the pretrained fusion network and focus on training the Text-DFH module to ensure that the fusion results meet the requirements of the pedestrian detection task. Text-DFH refines the features

output by the IR/VI-Encoders in the fusion network based on the responses from the multimodal large language model, enabling the fused image to better align with the needs of pedestrian detection. As shown in Figure 4, Text-DFH mainly consists of a dual-branch Cross Attention (CA) module and three feature extraction layers. The dual-branch cross attention computes the cross-attention between the features extracted by the IR/VI-Encoders and the textual features, allowing the model to extract useful information from the text that can help improve pedestrian detection accuracy. Subsequently, the three feature extraction layers integrate this textual information with the image scene features to generate refined

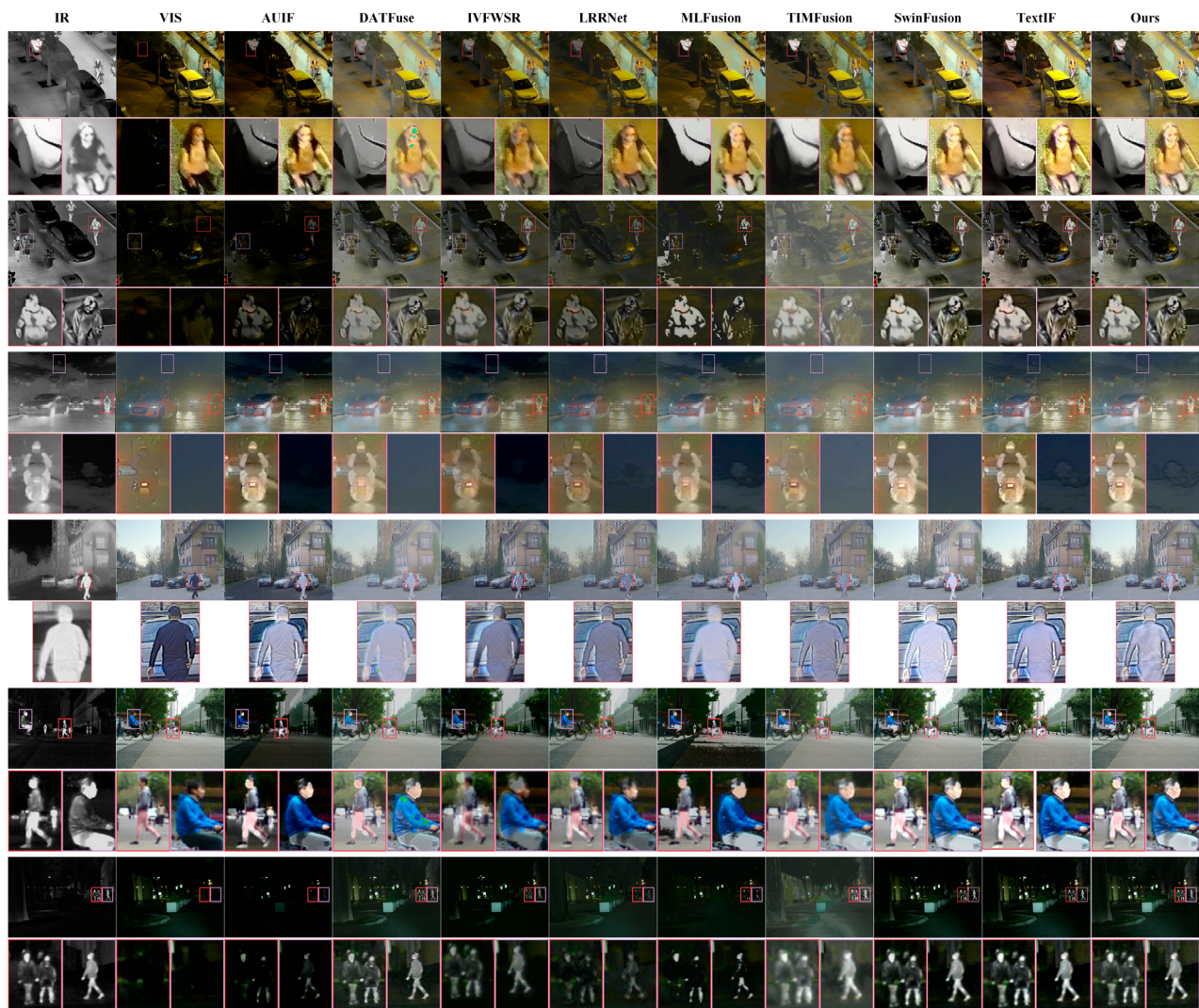


FIGURE 5

Visual comparison with SOTA methods. The top two rows, middle two rows, and bottom two rows of images are from the LLVIP, M³FD, and MSRS datasets, respectively. The first and second columns show the infrared and visible source images, while the third to ninth columns display the fusion results produced by the compared methods.

features. The structure of the CA module is similar to the Multi-Scale Attention (MSA) module used in DATFuse.

We input the infrared image I_i and visible image I_v into the pretrained fusion network with frozen parameters to obtain the fused image I_f . To obtain effective textual feedback that helps ensure the fused image meets the requirements of the pedestrian detection task, we input both I_f and the text prompt “To improve the accuracy of pedestrian detection, how can the quality of this image be enhanced?” into LLaVA, resulting in the textual feature T . We then input the outputs $F_{i/v}$ from the IR/VI-Encoders and the textual feature T into Text-DFH to harmonize the information in $F_{i/v}$. To comprehensively extract the task-relevant information from the textual features, we design a dual-branch processing strategy. In the first branch, we take $F_{i/v}$ as the Query (Q) and T as the Key (K) and Value (V) for cross-attention computation:

$$F_{i/v}^1 = \text{softmax} \left(\frac{Q_{i/v}^1 (K_{i/v}^1)^T}{\sqrt{d_1}} \right) V_{i/v}^1 \quad (4)$$

Here, $F_{i/v}^1$ represents the features injected with textual information in the first branch, d_1 denotes the dimensionality of $Q_{i/v}^1$, $Q_{i/v}^1 = W_{i/v}^{Q,1} F_{i/v}$, $K_{i/v}^1 = W_{i/v}^{K,1} T$, $V_{i/v}^1 = W_{i/v}^{V,1} T$. In the second branch, we use T as the Query (Q) and $F_{i/v}$ as the Key (K) and Value (V) for cross-attention computation:

$$F_{i/v}^2 = \text{softmax} \left(\frac{Q_{i/v}^2 (K_{i/v}^2)^T}{\sqrt{d_2}} \right) V_{i/v}^2 \quad (5)$$

Here, $F_{i/v}^2$ represents the features injected with textual information in the second branch, d_2 denotes the dimensionality of $Q_{i/v}^2$, and $Q_{i/v}^2 = W_{i/v}^{Q,2} T$, $K_{i/v}^2 = W_{i/v}^{K,2} F_{i/v}$, $V_{i/v}^2 = W_{i/v}^{V,2} F_{i/v}$. To comprehensively

TABLE 1 Quantitative results on the LLVIP dataset. The best and second-best values for each evaluation metric are highlighted in red and blue, respectively.

Methods	$Q_{AB/F}\uparrow$	$Q_{CV}\downarrow$	$Q_{SSIM}\uparrow$	$Q_{AG}\uparrow$	$Q_{SCD}\uparrow$
AUIF	0.3869	610.74	1.2016	3.5256	1.3413
DATFuse	0.4548	453.42	1.3130	3.1243	1.3351
IVFWSR	0.2925	512.77	1.2348	2.5252	1.1235
LRRNet	0.4426	534.89	1.3022	2.4625	0.9999
MLFusion	0.3239	523.41	1.2624	2.1613	0.9966
TIMFusion	0.2325	845.75	1.1742	2.1761	0.5368
SwinFusion	0.4266	598.53	1.2743	2.6346	1.3527
TextIF	0.5235	356.35	1.3056	3.4856	1.4527
Ours	0.5845	287.43	1.3441	3.9867	1.5462

TABLE 2 Quantitative results on the M³FD dataset. The best and second-best values for each evaluation metric are highlighted in red and blue, respectively.

Methods	$Q_{AB/F}\uparrow$	$Q_{CV}\downarrow$	$Q_{SSIM}\uparrow$	$Q_{AG}\uparrow$	$Q_{SCD}\uparrow$
AUIF	0.5425	852.56	1.3003	6.6735	1.5353
DATFuse	0.4854	563.57	1.3067	4.8326	1.3461
IVFWSR	0.4532	722.22	1.2735	3.5628	1.2452
LRRNet	0.5164	579.55	1.3735	4.5624	1.3461
MLFusion	0.4253	689.44	1.2835	4.4527	1.2687
TIMFusion	0.5352	616.16	1.2872	4.3336	1.2004
SwinFusion	0.5537	588.24	1.3086	6.0463	1.3456
TextIF	0.5423	534.21	1.2986	6.4026	1.5035
Ours	0.5856	454.45	1.3095	6.4561	1.6187

aggregate the textual information, we concatenate $F_{i/v}^1$ and $F_{i/v}^2$ along the channel dimension and feed the result into three feature extraction layers to obtain the harmonized features $\hat{F}_{i/v}$. We then concatenate \hat{F}_i and \hat{F}_v along the channel dimension and input the result into the F-Decoder to reconstruct the refined fused image I'_f . To ensure that the refined fused image I'_f meets the requirements of the pedestrian detection task, we introduce a pretrained pedestrian detection network with frozen parameters to supervise the fused image. We input I'_f into the detection network and obtain the pedestrian detection result \hat{y} . To make \hat{y} as close as possible to its ground truth y_{gt} , we constrain the Text-DFH module using the loss function ℓ_{pd} , which is the same as the one used during the training of YOLOv5.

TABLE 3 Quantitative results on the MSRS dataset. The best and second-best values for each evaluation metric are highlighted in red and blue, respectively.

Methods	$Q_{AB/F}\uparrow$	$Q_{CV}\downarrow$	$Q_{SSIM}\uparrow$	$Q_{AG}\uparrow$	$Q_{SCD}\uparrow$
AUIF	0.1736	799.97	0.9853	1.8844	1.1963
DATFuse	0.6326	416.67	1.2421	3.5481	1.5641
IVFWSR	0.3464	734.46	1.3462	2.1129	1.3581
LRRNet	0.4263	666.35	1.2952	2.5632	1.0854
MLFusion	0.2656	745.57	1.3457	2.6531	1.2053
TIMFusion	0.3346	1032.24	1.1003	2.6422	1.1783
SwinFusion	0.4527	439.46	1.3163	3.0042	1.4828
TextIF	0.6125	400.34	1.3357	3.6426	1.5457
Ours	0.6365	334.23	1.3537	3.5474	1.6854

4 Experiments

4.1 Datasets

The proposed method consists of two training stages. In both the first and second training stages, we train the fusion network and the text-driven feature harmonization module on the publicly available LLVIP dataset [67], respectively, in accordance with standard practices in the field [68–70]. Specifically, we randomly select 2,000 pairs of infrared and visible images from the LLVIP dataset as the training set. To enhance the diversity of training samples, we apply random flipping, random rotation, and random cropping as data augmentation techniques. For evaluation, we randomly select 200 pairs of infrared and visible images from each of the LLVIP, M³FD [71], and MSRS [3] datasets to form the test set, in order to assess both the fusion performance and pedestrian detection performance of the proposed method. Among them, LLVIP, M³FD, and MSRS are used to evaluate fusion performance, while LLVIP is specifically used to evaluate pedestrian detection performance.

4.2 Implementation details

The proposed method involves two training stages. In the first stage, the fusion network is trained. In the second stage, the parameters of the fusion network are frozen, and the text-driven feature harmonization module is trained. Both training stages use the Adam optimizer to update the network parameters, with a batch size of 16 and a learning rate of 1×10^{-3} . The total number of training epochs is set to 100 for the first stage and 200 for the second stage. In addition, the hyperparameter ε is set to 0.2. The proposed method is implemented based on the PyTorch framework and is trained on a single NVIDIA RTX A6000 GPU.

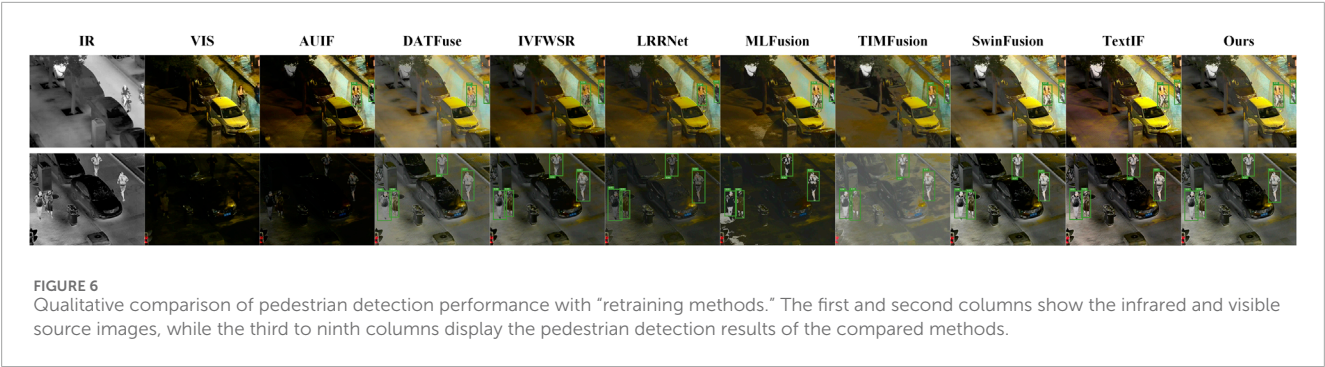


TABLE 4 Quantitative comparison of pedestrian detection performance with “retraining methods.” The best and second-best values for each evaluation metric are highlighted in red and blue, respectively.

Methods	mAP ₅₀ ↑	mAP ₇₅ ↑	mAP _{50→95} ↑
AUIF	98.2	91.8	74.4
DATFuse	99.0	91.5	74.3
IVFWSR	97.2	89.6	72.9
LRRNet	98.0	90.8	73.8
MLFusion	97.8	89.9	73.6
TIMFusion	97.9	88.4	74.0
SwinFusion	98.5	90.4	74.3
TextIF	98.9	91.7	74.6
Ours	99.1	92.8	75.0

4.3 Evaluation metrics

We adopt five commonly used objective evaluation metrics to quantitatively assess the fusion performance of the proposed method. These metrics include Edge Preservation Index ($Q_{AB/F}$) [72, 73], Chen-Varshney Index (Q_{CV}) [74], Structural Similarity Index (Q_{SSIM}) [75], Average Gradient (Q_{AG}) [76], and Sum of Correlations of Differences (Q_{SCD}) [77]. $Q_{AB/F}$ measures how well edge information from the source images is preserved in the fused image. $Q_{AB/F}$ higher value indicates less loss of texture details in the fused image. Q_{CV} evaluates fusion quality from the perspective of human visual perception; a lower value means the fused image aligns better with human visual preferences. Q_{SSIM} quantifies the similarity between the fused image and the source images in terms of luminance, contrast, and structure. A higher value indicates less information difference between the fused and source images. Q_{AG} measures the richness of gradient information in the fused image. A higher value means the fused image contains more detailed gradient content. Q_{SCD} assesses information loss during the fusion process by computing difference maps between the fused image and source images. A higher value indicates less distortion in the fused image. Among these, $Q_{AB/F}$, Q_{SSIM} , Q_{AG} and Q_{SCD} are

positive indicators, meaning a higher value indicates better fusion performance. Q_{CV} is a negative indicator, meaning a lower value represents better fusion performance. In addition, to objectively evaluate the effectiveness of the fused images in the pedestrian detection task, we adopt three widely used metrics in the pedestrian detection domain for quantitative analysis: Mean Average Precision (mAP) at IoU threshold of 0.5 (mAP_{50}), mAP at IoU threshold of 0.75 (mAP_{75}), and the averaged mAP at IoU threshold from 0.5 to 0.95 ($mAP_{50→95}$).

4.4 Comparison with state-of-the-art methods

In this study, we conduct a series of qualitative and quantitative comparisons between the proposed method and eight state-of-the-art (SOTA) methods to verify its superiority in both fusion performance and pedestrian detection performance. These methods include AUIF [78], DATFuse [49], IVFWSR [79], LRRNet [80], MLFusion [81], TIMFusion [82], SwinFusion [16], and TextIF [29]. The comparative experiments are divided into two distinct groups: In the first group, we compare the fusion performance of our method with that of the SOTA methods. In the second group, we freeze the fusion networks of the compared methods and retrain their pedestrian detection networks using the corresponding fused results. The retrained detection networks are then used to perform pedestrian detection on the fused images. This setup is designed to demonstrate that our proposed method can achieve strong pedestrian detection performance without requiring retraining of the detection network.

4.4.1 Fusion performance comparison

We conduct both quantitative and qualitative comparisons of the proposed method against AUIF, DATFuse, IVFWSR, LRRNet, MLFusion, TIMFusion, SwinFusion, and TextIF on the LLVIP, M³FD, and MSRS datasets to validate the superiority of our method in terms of fusion performance. As shown in the enlarged regions of Figure 5, our method effectively highlights the thermal radiation information from the infrared image while preserving fine texture details from the visible image. Compared to existing SOTA methods, the fused images produced by our method exhibit clearer local details as well as higher overall brightness and contrast at the global level. This not only improves visual quality but also facilitates better object recognition in downstream tasks. This

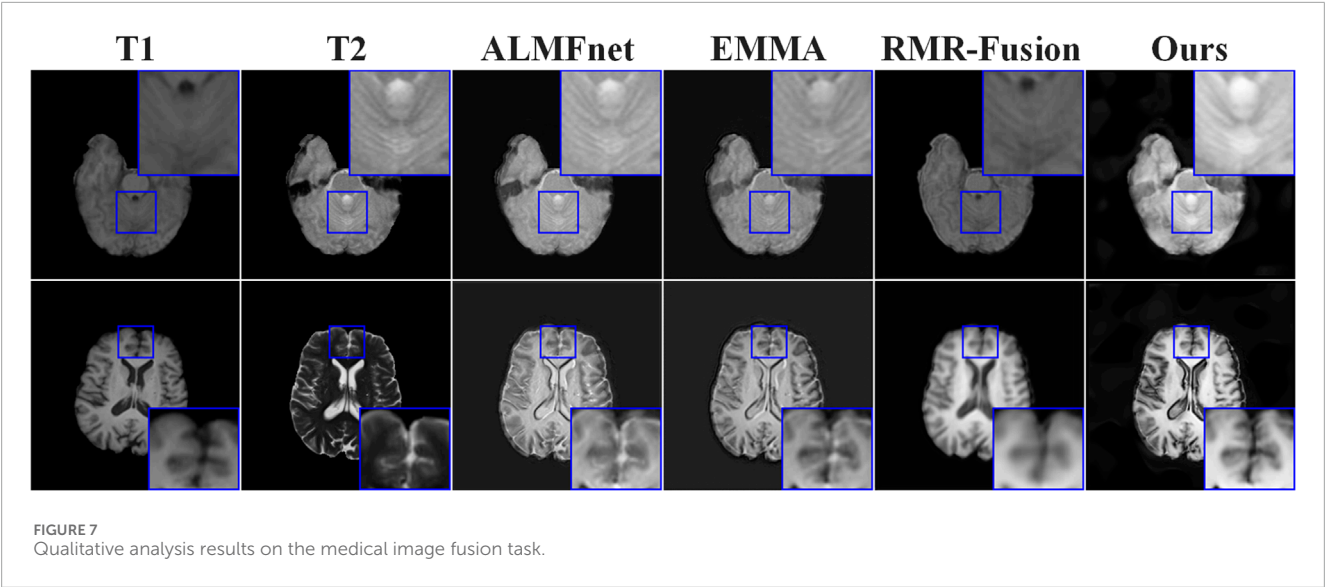


TABLE 5 Quantitative Analysis Results on the Medical Image Fusion Task. The best and second-best values for each evaluation metric are highlighted in red and blue, respectively.

Methods	$Q_{AB/F} \uparrow$	$Q_{CV} \downarrow$	$Q_{SSIM} \uparrow$	$Q_{AG} \uparrow$	$Q_{SCD} \uparrow$
ALMFnet	0.4700	1330.59	1.3432	3.5826	1.2991
EMMA	0.4682	1288.99	1.3232	3.1826	1.2999
RMR-Fusion	0.4419	1344.12	1.2967	3.2621	1.3781
Ours	0.4792	1203.12	1.3631	3.7521	1.3629

advantage is also reflected in the quantitative evaluation results, as shown in Tables 1–3. Specifically, our method achieves the lowest values in metric Q_{CV} , and ranks first in both metrics $Q_{AB/F}$ and Q_{AG} , indicating that the fused images contain richer edge information and are more consistent with human visual perception. In summary, both qualitative and quantitative results demonstrate that our proposed method offers significant improvements in fusion performance over the compared methods.

4.4.2 Pedestrian detection performance comparison

A common practice to improve the performance of fusion networks in downstream tasks is to freeze the parameters of the fusion network and retrain the downstream task network based on the generated fused results. Such approaches are referred to as “retraining methods.” To evaluate the effectiveness of our proposed method in pedestrian detection, we perform both quantitative and qualitative comparisons against these retraining methods. As shown in Figure 6, the pedestrian detection results of other methods often suffer from issues such as bounding boxes that fail to fully cover the pedestrians’ bodies, or boxes that include large amounts of irrelevant background, indicating insufficient detection accuracy. In contrast, the detection results produced by our method show significantly fewer irrelevant regions within the bounding boxes and more

accurate box placement. This advantage is also clearly reflected in the quantitative results, as shown in Table 4. Our method achieves the highest scores in metrics mAP_{50} , mAP_{75} , and $mAP_{50 \rightarrow 95}$, indicating superior performance in the pedestrian detection task compared to the other methods. In conclusion, our method demonstrates better performance than approaches that require retraining the pedestrian detection network, even without retraining. This highlights the effectiveness and advantage of our method in pedestrian detection tasks.

4.4.3 Analysis of application potential in medical image fusion

Furthermore, to validate the effectiveness and application potential of the proposed method in the field of nuclear medical imaging, we further deployed it in a medical image fusion task. Specifically, we conducted experiments on the BraTS2020 [83] dataset and performed both qualitative and quantitative analyses of the fusion results. As shown in Figure 7, compared with state-of-the-art methods such as ALMFnet [84, 85], and RMR-Fusion [86], the proposed method preserves more texture details and salient information in the fused medical images. As reported in Table 5, our method ranks first or second across most evaluation metrics. These results demonstrate the promising potential of the proposed method for applications in nuclear medical imaging.

4.5 Ablation study

The proposed method mainly consists of two core components: the Multimodal Large Language Model (MLLM) and the Text-Driven Feature Harmonization (Text-DFH) module. Within Text-DFH, both the text-guided cross-attention and the image-guided cross-attention play key roles. To validate the effectiveness of these components, we conduct a series of ablation experiments on the LLVIP dataset.

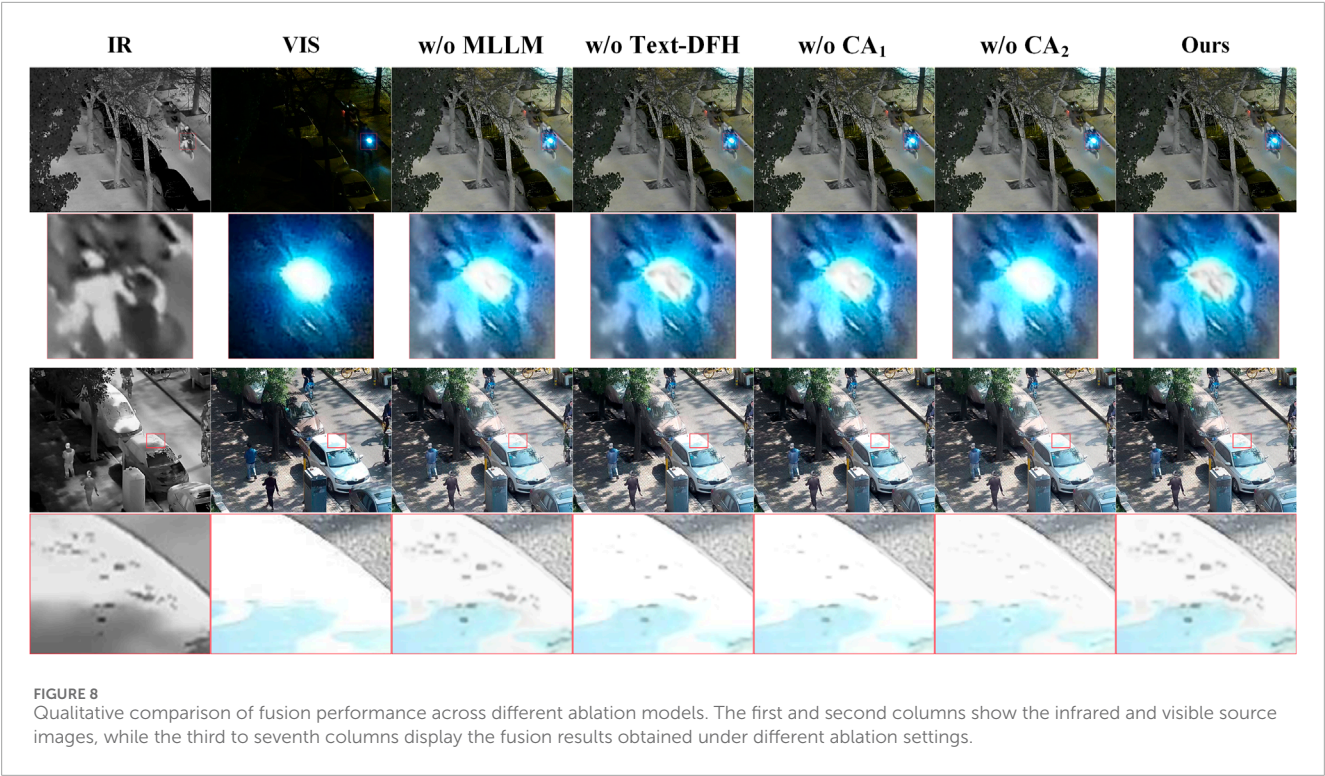


TABLE 6 Quantitative comparison of fusion performance across different ablation models. The best and second-best values for each evaluation metric are highlighted in red and blue, respectively.

Methods	$Q_{AB/F}\uparrow$	$Q_{CV}\downarrow$	$Q_{SSIM}\uparrow$	$Q_{AG}\uparrow$	$Q_{SCD}\uparrow$
w/o MLLM	0.5472	298.75	1.3244	3.6433	1.5367
w/o Text-DFH	0.5763	299.46	1.3321	3.4131	1.4992
w/o CA1	0.5834	305.92	1.3234	3.6362	1.5213
w/o CA2	0.5798	301.68	1.3401	3.6524	1.5123
Ours	0.5845	287.43	1.3441	3.9867	1.5462

TABLE 7 Quantitative comparison of pedestrian detection performance across different ablation models. The best and second-best values for each evaluation metric are highlighted in red and blue, respectively.

Methods	mAP ₅₀ \uparrow	mAP ₇₅ \uparrow	mAP ₅₀₋₉₅ \uparrow
w/o MLLM	98.5	91.6	73.9
w/o Text-DFH	98.8	92.1	74.0
w/o CA1	99.0	92.4	74.5
w/o CA2	98.9	91.8	74.4
Ours	99.1	92.8	75.0

4.5.1 Effectiveness of the multimodal large language model

We utilize the MLLM to analyze the fused images based on user-provided questions related to pedestrian detection performance and generate suggestions for improving image quality. To assess the contribution of the MLLM, we remove it and replace its feedback with a fixed text prompt: “Brighter brightness, higher contrast, and clearer texture details.” As shown in Figure 8, the fusion results from the ablation model without the MLLM are noticeably inferior in visual quality compared to the full model. To further validate this, we perform quantitative analysis as presented in Table 6. The results show that the full model outperforms the ablation model on all evaluation metrics. Additionally, we analyze the performance of pedestrian detection, as shown in Table 7 and Figure 9. Both the quantitative and qualitative results indicate that the fused images produced by the ablation model

without the MLLM lead to poorer detection performance. In contrast, the full model achieves better pedestrian detection results. In summary, both qualitative and quantitative analyses confirm the effectiveness of the Multimodal Large Language Model in our method.

4.5.2 Effectiveness of Text-DFH

Text-DFH refines the output features of the fusion network based on suggestions from the multimodal large language model, enabling the fused image to better meet the requirements of the pedestrian detection task. To verify the effectiveness of Text-DFH, we remove it from the architecture and instead concatenate the text features with the image features to be refined along the channel dimension. The combined features are then processed by CNNs to obtain the refined output. We conduct both quantitative and qualitative analyses of the fusion

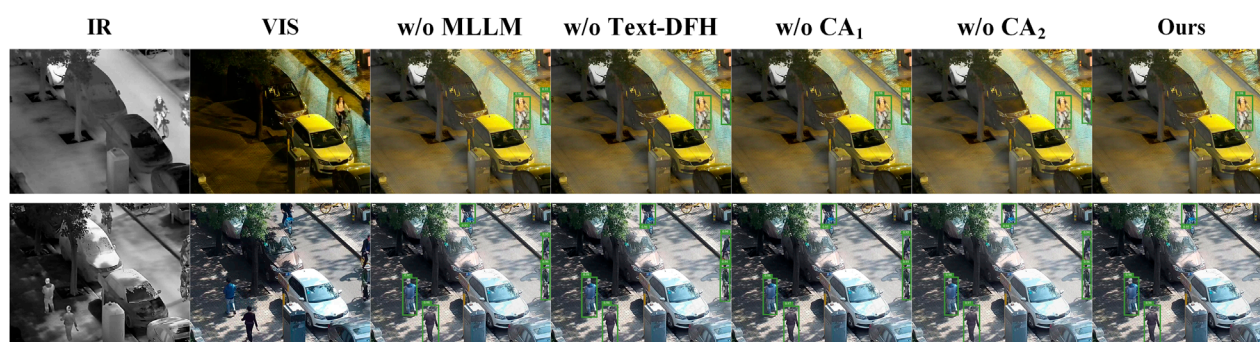


FIGURE 9

Qualitative comparison of pedestrian detection performance across different ablation models. The first and second columns show the infrared and visible source images, while the third to seventh columns display the pedestrian detection results under different ablation settings.

performance of the model without Text-DFH, as shown in Table 6 and Figure 8. As observed, the ablation model without Text-DFH performs worse than the full model across multiple evaluation metrics, and the visual quality of the fused images is also inferior. In addition, we evaluate pedestrian detection performance both quantitatively and qualitatively, as presented in Table 7 and Figure 9. The full model achieves higher scores compared to the ablation model without Text-DFH. In summary, a series of experiments clearly demonstrate the effectiveness of the Text-DFH module.

4.5.3 Effectiveness of dual-branch cross attention

In the Text-DFH module, we refine image features using text features through a dual-branch cross attention mechanism. To verify its effectiveness, we remove the cross attention from each branch individually, leaving only a single branch to refine the image features. These variants are referred to as CA1 and CA2, respectively. From the quantitative and qualitative results on fusion performance, it is evident that removing either branch of the cross attention leads to a significant drop in performance, as shown in Table 6 and Figure 8. Furthermore, to assess the impact of dual-branch cross attention on pedestrian detection performance, we conduct both quantitative and qualitative analyses. The results demonstrate that pedestrian detection performance is optimal only when both branches of the cross attention are used to refine the image features, as shown in Table 7 and Figure 9. In conclusion, the above experiments confirm the effectiveness of the dual-branch cross attention mechanism.

5 Conclusion

To address the limitation of existing methods that primarily focus on improving fused image quality through network design—while overlooking the potential benefits of enhanced image quality for pedestrian detection—we propose a multimodal large language model (MLLM)-driven infrared and visible image fusion method. This method not only aims to improve the quality of the fused images but also emphasizes enhancing their performance

in pedestrian detection tasks. By leveraging a multimodal large language model, we analyze the fused images based on user-provided questions related to improving pedestrian detection performance and generate suggestions for enhancing image quality. To fully utilize the guidance provided by the MLLM, we design a Text-Driven Feature Harmonization (Text-DFH) module, which refines the features output by the fusion network according to the textual suggestions. This ensures improved fusion quality while maintaining strong performance in pedestrian detection. In addition, the proposed method also demonstrates significant application potential in the field of nuclear medical imaging. However, under extreme weather conditions such as rain, fog, and snow, the fusion performance of the current method may degrade. Moreover, when such methods are applied to other types of source images [87–90], their performance may degrade. In future work, we plan to extend this research to develop an infrared and visible image fusion framework tailored for extreme weather scenarios, striving to maintain robust downstream task performance even in challenging environments.

Data availability statement

The original contributions presented in the study are included in the article/supplementary material, further inquiries can be directed to the corresponding author.

Author contributions

KW: Project administration, Writing – original draft, Writing – review and editing, Investigation, Conceptualization, Methodology. DH: Formal Analysis, Writing – review and editing, Data curation, Validation. YaC: Visualization, Supervision, Writing – review and editing, Resources. YkC: Funding acquisition, Project administration, Supervision, Writing – review and editing, Writing – original draft. YL: Validation, Writing – review and editing, Visualization, Formal Analysis. ZJ: Writing – review and editing, Investigation, Data curation, Resources. FC: Formal Analysis,

Writing – review and editing, Data curation. WL: Writing – review and editing, Resources, Visualization.

Funding

The author(s) declare that financial support was received for the research and/or publication of this article. Science and Technology Project of China Southern Power Grid Co., Ltd. (No. YNKJXM20240052).

Conflict of interest

Authors KW, DH, YaC, YkC, YL, ZJ, FC, and WL were employed by Yunnan Power Grid Co., Ltd.

References

- Li H, Wu X-J, Kittler J. Rfn-nest: an end-to-end residual fusion network for infrared and visible images. *Inf Fusion* (2021) 73:72–86. doi:10.1016/j.inffus.2021.02.023
- Ma J, Yu W, Liang P, Li C, Jiang J. Fusiongan: a generative adversarial network for infrared and visible image fusion. *Inf Fusion* (2019) 48:11–26. doi:10.1016/j.inffus.2018.09.004
- Tang L, Yuan J, Zhang H, Jiang X, Ma J. Piafusion: a progressive infrared and visible image fusion network based on illumination aware. *Inf Fusion* (2022) 83–84:79–92. doi:10.1016/j.inffus.2022.03.007
- Xu M, Tang L, Zhang H, Ma J. Infrared and visible image fusion via parallel scene and texture learning. *Pattern Recognition* (2022) 132:108929. doi:10.1016/j.patcog.2022.108929
- Du K, Li H, Zhang Y, Yu Z. Chitnet: a complementary to harmonious information transfer network for infrared and visible image fusion. *IEEE Trans Instrumentation Meas* (2025) 74:1–17. doi:10.1109/TIM.2025.3527523
- Shi Y, Liu Y, Cheng J, Wang ZJ, Chen X. Vdmufusion: a versatile diffusion model-based unsupervised framework for image fusion. *IEEE Trans Image Process* (2025) 34:441–54. doi:10.1109/tip.2024.3512365
- Lv G, Sima C, Gao Y, Dong A, Ma G, Cheng J. Sigfusion: semantic information-guided infrared and visible image fusion. *IEEE Trans Instrumentation Meas* (2024) 73:1–18. doi:10.1109/tim.2024.3457951
- Fu H, Wang S, Duan P, Xiao C, Dian R, Li S, et al. Lraf-net: long-range attention fusion network for visible–infrared object detection. *IEEE Trans Neural Networks Learn Syst* (2024) 35:13232–45. doi:10.1109/tnnls.2023.3266452
- Li Y, Pang Y, Cao J, Shen J, Shao L. Improving single shot object detection with feature scale unmixing. *IEEE Trans Image Process* (2021) 30:2708–21. doi:10.1109/tip.2020.3048630
- Zhao Z-Q, Zheng P, Xu S-T, Wu X. Object detection with deep learning: a review. *IEEE Trans Neural Networks Learn Syst* (2019) 30:3212–32. doi:10.1109/tnnls.2018.2876865
- Liu Y, Zeng J, Tao X, Fang G. Rethinking self-supervised semantic segmentation: achieving end-to-end segmentation. *IEEE Trans Pattern Anal Machine Intelligence* (2024) 46:10036–46. doi:10.1109/tpami.2024.3432326
- Wu L, Fang L, He X, He M, Ma J, Zhong Z. Querying labeled for unlabeled: cross-image semantic consistency guided semi-supervised semantic segmentation. *IEEE Trans Pattern Anal Machine Intelligence* (2023) 45:8827–44. doi:10.1109/TPAMI.2022.3233584
- Zhao S, Zhang Q. A feature divide-and-conquer network for rgb-t semantic segmentation. *IEEE Trans Circuits Syst Video Technology* (2023) 33:2892–905. doi:10.1109/tcsvt.2022.3229359
- Zhao S, Liu Y, Jiao Q, Zhang Q, Han J. Mitigating modality discrepancies for rgb-t semantic segmentation. *IEEE Trans Neural Networks Learn Syst* (2024) 35:9380–94. doi:10.1109/tnnls.2022.3233089
- Li H, Wu X. Densefuse: a fusion approach to infrared and visible images. *IEEE Trans Image Process* (2019) 28:2614–23. doi:10.1109/tip.2018.2887342
- Ma J, Tang L, Fan F, Huang J, Mei X, Ma Y. Swinfusion: cross-domain long-range learning for general image fusion via swin transformer. *IEEE/CAA J Automatica Sinica* (2022) 9:1200–17. doi:10.1109/jas.2022.105686
- Tang W, He F, Liu Y. Ydtr: infrared and visible image fusion via y-shape dynamic transformer. *IEEE Trans Multimedia* (2023) 25:5413–28. doi:10.1109/tmm.2022.3192661
- Li H, Qiu H, Yu Z, Zhang Y. Infrared and visible image fusion scheme based on nsct and low-level visual features. *Infrared Phys and Technology* (2016) 76:174–84. doi:10.1016/j.infrared.2016.02.005
- Li H, Wang Y, Yang Z, Wang R, Li X, Tao D. Discriminative dictionary learning-based multiple component decomposition for detail-preserving noisy image fusion. *IEEE Trans Instrumentation Meas* (2020) 69:1082–102. doi:10.1109/tim.2019.2912239
- Hou R, Zhou D, Nie R, Liu D, Xiong L, Guo Y, et al. Vif-net: an unsupervised framework for infrared and visible image fusion. *IEEE Trans Comput Imaging* (2020) 6:640–51. doi:10.1109/tci.2020.2965304
- Ma J, Zhang H, Shao Z, Liang P, Xu H. Ganmcc: a generative adversarial network with multiclassification constraints for infrared and visible image fusion. *IEEE Trans Instrumentation Meas* (2021) 70:1–14. doi:10.1109/tim.2020.3038013
- Xu H, Wang X, Ma J. Drf: disentangled representation for visible and infrared image fusion. *IEEE Trans Instrumentation Meas* (2021) 70:1–13. doi:10.1109/tim.2021.3056645
- Zhang Y, Yang M, Li N, Yu Z. Analysis-synthesis dictionary pair learning and patch saliency measure for image fusion. *Signal Process.* (2020) 167:107327. doi:10.1016/j.sigpro.2019.107327
- Tang L, Zhang H, Xu H, Ma J. 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* (2023) 99:101870. doi:10.1016/j.inffus.2023.101870
- Liu J, Liu Z, Wu G, Ma L, Liu R, Zhong W, et al. Multi-interactive feature learning and a full-time multi-modality benchmark for image fusion and segmentation. In: *Proceedings of the IEEE/CVF international conference on computer vision (ICCV)* (2023). p. 8115–24.
- Zhang H, Zuo X, Jiang J, Guo C, Ma J. Mrfs: mutually reinforcing image fusion and segmentation. In: *2024 IEEE/CVF conference on computer vision and pattern recognition (CVPR)* (2024). p. 26964–73.
- Wang D, Liu J, Liu R, Fan X. An interactively reinforced paradigm for joint infrared-visible image fusion and saliency object detection. *Inf Fusion* (2023) 98:101828. doi:10.1016/j.inffus.2023.101828
- Yang Z, Zhang Y, Li H, Liu Y. Instruction-driven fusion of infrared-visible images: tailoring for diverse downstream tasks. *Inf Fusion* (2025) 121:103148. doi:10.1016/j.inffus.2025.103148
- Yi X, Xu H, Zhang H, Tang L, Ma J. Text-if: leveraging semantic text guidance for degradation-aware and interactive image fusion. In: *2024 IEEE/CVF conference on computer vision and pattern recognition (CVPR)* (2024). p. 27016–25.
- Wang H, Zhang H, Yi X, Xiang X, Fang L, Ma J. Terf: text-driven and region-aware flexible visible and infrared image fusion. In: *Proceedings of the 32nd ACM international conference on multimedia* (2024). p. 935–44.
- Achiam J, Adler S, Agarwal S, Ahmad L, Akkaya I, Aleman FL, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).
- Team G, Anil R, Borgeaud S, Alayrac J-B, Yu J, Soricut R, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805* (2023).

Generative AI statement

The author(s) declare that Generative AI was used in the creation of this manuscript. AI was only used to polish the paper.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

33. Wu S, Fei H, Qu L, Ji W, Chua T-S. Next-gpt: any-to-any multimodal llm. In: *Forty-first international conference on machine learning* (2024).
34. Liu Y, Zhang K, Li Y, Yan Z, Gao C, Chen R, et al. Sora: a review on background, technology, limitations, and opportunities of large vision models. *arXiv preprint arXiv:2402* (2024).
35. Li K, He Y, Wang Y, Li Y, Wang W, Luo P, et al. Videochat: chat-centric video understanding. *arXiv preprint arXiv:2305.06355* (2023).
36. Lin B, Ye Y, Zhu B, Cui J, Ning M, Jin P, et al. Video-llava: learning united visual representation by alignment before projection. *arXiv preprint arXiv:2311* (2023).
37. Yang Z, Li Y, Tang X, Xie M. Mgfusion: a multimodal large language model-guided information perception for infrared and visible image fusion. *Front Neurobot* (2024) 18:1521603. doi:10.3389/fnbot.2024.1521603
38. Radford A, Kim JW, Hallacy C, Ramesh A, Goh G, Agarwal S, et al. Learning transferable visual models from natural language supervision. In: *International conference on machine learning (PMLR)* (2021). p. 8748–63.
39. Touvron H, Lavril T, Izacard G, Martinet X, Lachaux M-A, Lacroix T, et al. Llama: open and efficient foundation language models. *arXiv preprint arXiv:2302* (2023).
40. Xu H, Ma J, Jiang J, Guo X, Ling H. U2fusion: a unified unsupervised image fusion network. *IEEE Trans Pattern Anal Machine Intelligence* (2022) 44:502–18. doi:10.1109/tpami.2020.3012548
41. Zhang Y, Liu Y, Sun P, Yan H, Zhao X, Zhang L. Ifcnn: a general image fusion framework based on convolutional neural network. *Inf Fusion* (2020) 54:99–118. doi:10.1016/j.inffus.2019.07.011
42. Zhang H, Ma J. Sdnet: a versatile squeeze-and-decomposition network for real-time image fusion. *Int J Computer Vis* (2021) 129:2761–85. doi:10.1007/s11263-021-01501-8
43. Li H, Yang Z, Zhang Y, Jia W, Yu Z, Liu Y. Mulfs-cap: multimodal fusion-supervised cross-modality alignment perception for unregistered infrared-visible image fusion. *IEEE Trans Pattern Anal Machine Intelligence* (2025) 47:3673–90. doi:10.1109/TPAMI.2025.3535617
44. Li H, Zhao J, Li J, Yu Z, Lu G. Feature dynamic alignment and refinement for infrared-visible image fusion: translation robust fusion. *Inf Fusion* (2023) 95:26–41. doi:10.1016/j.inffus.2023.02.011
45. Xiao W, Zhang Y, Wang H, Li F, Jin H. Heterogeneous knowledge distillation for simultaneous infrared-visible image fusion and super-resolution. *IEEE Trans Instrumentation Meas* (2022) 71:1–15. doi:10.1109/tim.2022.3149101
46. Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, et al. Attention is all you need. *Adv Neural Inf Process Syst* (2017) 30.
47. Dosovitskiy A, Beyer L, Kolesnikov A, Weissenborn D, Zhai X, Unterthiner T, et al. An image is worth 16x16 words: transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).
48. Li J, Zhu J, Li C, Chen X, Yang B. Cgtf: convolution-guided transformer for infrared and visible image fusion. *IEEE Trans Instrumentation Meas* (2022) 71:1–14. doi:10.1109/tim.2022.3175055
49. Tang W, He F, Liu Y, Duan Y, Si T. Datfuse: infrared and visible image fusion via dual attention transformer. *IEEE Trans Circuits Syst Video Technology* (2023) 33:3159–72. doi:10.1109/tcsvt.2023.3234340
50. Zhao Z, Bai H, Zhang J, Zhang Y, Xu S, Lin Z, et al. Cddfuse: correlation-driven dual-branch feature decomposition for multi-modality image fusion. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)* (2023). p. 5906–16.
51. Zhao Z, Bai H, Zhang J, Zhang Y, Zhang K, Xu S, et al. Equivariant multi-modality image fusion. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2024). p. 25912–21.
52. Ma J, Liang P, Yu W, Chen C, Guo X, Wu J, et al. Infrared and visible image fusion via detail preserving adversarial learning. *Inf Fusion* (2020) 54:85–98. doi:10.1016/j.inffus.2019.07.005
53. Ma J, Xu H, Jiang J, Mei X, Zhang X. Ddcgan: a dual-discriminator conditional generative adversarial network for multi-resolution image fusion. *IEEE Trans Image Process* (2020) 29:4980–95. doi:10.1109/tip.2020.2977573
54. Li J, Huo H, Li C, Wang R, Feng Q. Attentionfgan: infrared and visible image fusion using attention-based generative adversarial networks. *IEEE Trans Multimedia* (2021) 23:1383–96. doi:10.1109/tmm.2020.2997127
55. Zhou H, Wu W, Zhang Y, Ma J, Ling H. Semantic-supervised infrared and visible image fusion via a dual-discriminator generative adversarial network. *IEEE Trans Multimedia* (2021) 25:635–48. doi:10.1109/tmm.2021.3129609
56. Zhang H, Yuan J, Tian X, Ma J. Gan-fm: infrared and visible image fusion using gan with full-scale skip connection and dual markovian discriminators. *IEEE Trans Comput Imaging* (2021) 7:1134–47. doi:10.1109/tci.2021.3119954
57. Cai Z, Vasconcelos N. Cascade r-cnn: high quality object detection and instance segmentation. *IEEE Trans pattern anal machine intelligence* (2019) 43:1483–98. doi:10.1109/tpami.2019.2956516
58. Girshick R, Donahue J, Darrell T, Malik J. Rich feature hierarchies for accurate object detection and semantic segmentation. In: *Proceedings of the IEEE conference on computer vision and pattern recognition* (2014). p. 580–7.
59. Ren S, He K, Girshick R, Sun J. Faster r-cnn: towards real-time object detection with region proposal networks. *Adv Neural Inf Process Syst* (2015) 28.
60. Redmon J, Divvala S, Girshick R, Farhadi A. You only look once: unified, real-time object detection. In: *Proceedings of the IEEE conference on computer vision and pattern recognition* (2016). p. 779–88.
61. Liu W, Anguelov D, Erhan D, Szegedy C, Reed S, Fu C-Y, et al. Ssd: single shot multibox detector. In: *Computer vision—ECCV 2016: 14th European conference, Amsterdam, The Netherlands, October 11–14, 2016, proceedings, Part I* 14. Springer (2016). p. 21–37.
62. Carion N, Massa F, Synnaeve G, Usunier N, Kirillov A, Zagoruyko S. End-to-end object detection with transformers. In: *European conference on computer vision*. Springer (2020). p. 213–29.
63. Wu P, Zhai W, Cao Y. Background activation suppression for weakly supervised object localization. In: *2022 IEEE/CVF conference on computer vision and pattern recognition (CVPR)*. IEEE (2022). p. 14228–37.
64. Xu J, Hou J, Zhang Y, Feng R, Zhao R-W, Zhang T, et al. Cream: weakly supervised object localization via class re-activation mapping. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2022). p. 9437–46.
65. Zhang S, Yu Z, Liu L, Wang X, Zhou A, Chen K. Group r-cnn for weakly semi-supervised object detection with points. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2022). p. 9417–26.
66. Liu H, Li C, Wu Q, Lee YJ. Visual instruction tuning. *Adv Neural Inf Process Syst* (2023) 36:34892–916.
67. Jia X, Zhu C, Li M, Tang W, Zhou W. Llvip: a visible-infrared paired dataset for low-light vision. In: *Proceedings of the IEEE/CVF international conference on computer vision workshops (ICCVW)* (2021). p. 3496–504.
68. Tang L, Huang H, Zhang Y, Qi G, Yu Z. Structure-embedded ghosting artifact suppression network for high dynamic range image reconstruction. *Knowledge-Based Syst* (2023) 263:110278. doi:10.1016/j.knsys.2023.110278
69. Liu Y, Shi Y, Mu F, Cheng J, Chen X. Glioma segmentation-oriented multi-modal mr image fusion with adversarial learning. *IEEE/CAA J Automatica Sinica* (2022) 9:1528–31. doi:10.1109/jas.2022.105770
70. Xie M, Wang J, Zhang Y. A unified framework for damaged image fusion and completion based on low-rank and sparse decomposition. *Signal Processing: Image Commun* (2021) 29:116400. doi:10.1016/j.image.2021.116400
71. Liu J, Fan X, Huang Z, Wu G, Liu R, Zhong W, et al. Target-aware dual adversarial learning and a multi-scenario multi-modality benchmark to fuse infrared and visible for object detection. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)* (2022). p. 5802–11.
72. Liu Y, Qi Z, Cheng J, Chen X. Rethinking the effectiveness of objective evaluation metrics in multi-focus image fusion: a statistic-based approach. *IEEE Trans Pattern Anal Machine Intelligence* (2024) 46:5806–19. doi:10.1109/tpami.2024.3367905
73. Xydeas CS, Petrovic V, et al. Objective image fusion performance measure. *Electronics Lett* (2000) 36:308–9.
74. Chen H, Varshney PK. A human perception inspired quality metric for image fusion based on regional information. *Inf Fusion* (2007) 8:193–207. doi:10.1016/j.inffus.2005.10.001
75. Wang Z, Bovik A, Sheikh H, Simoncelli E. Image quality assessment: from error visibility to structural similarity. *IEEE Trans Image Process* (2004) 13:600–12. doi:10.1109/tip.2003.819861
76. Liu J, Wu G, Liu Z, Wang D, Jiang Z, Ma L, et al. Infrared and visible image fusion: from data compatibility to task adaption. *IEEE Trans Pattern Anal Machine Intelligence* (2024) 1–20. doi:10.1109/TPAMI.2024.3521416
77. Zhang X, Ye P, Xiao G. Vfbb: a visible and infrared image fusion benchmark. In: *2020 IEEE/CVF conference on computer vision and pattern recognition workshops (CVPRW)* (2020). p. 468–78.
78. Zhao Z, Xu S, Zhang J, Liang C, Zhang C, Liu J. Efficient and model-based infrared and visible image fusion via algorithm unrolling. *IEEE Trans Circuits Syst Video Technology* (2021) 32:1186–96. doi:10.1109/TCSVT.2021.3075745
79. Li H, Liu J, Zhang Y, Liu Y. A deep learning framework for infrared and visible image fusion without strict registration. *Int J Computer Vis* (2024) 132:1625–44. doi:10.1007/s11263-023-01948-x
80. Li H, Xu T, Wu X-J, Lu J, Kittler J. Lrrnet: a novel representation learning guided fusion network for infrared and visible images. *IEEE Trans Pattern Anal Machine Intelligence* (2023) 45:11040–52. doi:10.1109/tpami.2023.3268209
81. Li H, Cen Y, Liu Y, Chen X, Yu Z. Different input resolutions and arbitrary output resolution: a meta learning-based deep framework for infrared and visible image fusion. *IEEE Trans Image Process* (2021) 30:4070–83. doi:10.1109/tip.2021.3069339
82. Liu R, Liu Z, Liu J, Fan X, Luo Z. A task-guided, implicitly-searched and meta-initialized deep model for image fusion. *IEEE Trans Pattern Anal Machine Intelligence* (2024) 46:6594–609. doi:10.1109/tpami.2024.3382308

83. Menze BH, Jakab A, Bauer S, Kalpathy-Cramer J, Farahani K, Kirby J, et al. The multimodal brain tumor image segmentation benchmark (brats). *IEEE Trans Med Imaging* (2015) 34:1993–2024. doi:10.1109/tmi.2014.2377694
84. Mu P, Wu G, Liu J, Zhang Y, Fan X, Liu R. Learning to search a lightweight generalized network for medical image fusion. *IEEE Trans Circuits Syst Video Technology* (2024) 34:5921–34. doi:10.1109/tcsvt.2023.3342808
85. Zhao Z, Bai H, Zhang J, Zhang Y, Zhang K, Xu S, et al. Equivariant multi-modality image fusion. In: *2024 IEEE/CVF conference on computer vision and pattern recognition (CVPR)* (2024). p. 25912–21.
86. Zhang H, Zuo X, Zhou H, Lu T, Ma J. A robust mutual-reinforcing framework for 3d multi-modal medical image fusion based on visual-semantic consistency. *Proc AAAI Conf Artif Intelligence* (2024) 38:7087–95. doi:10.1609/aaai.v38i7.28536
87. Zhang Y, Yang X, Li H, Xie M, Yu Z. Dcpnet: a dual-task collaborative promotion network for pansharpening. *IEEE Trans Geosci Remote Sensing* (2024) 62:1–16. doi:10.1109/tgrs.2024.3377635
88. Li H, Yang Z, Zhang Y, Tao D, Yu Z. Single-image hdr reconstruction assisted ghost suppression and detail preservation network for multi-exposure hdr imaging. *IEEE Trans Comput Imaging* (2024) 10:429–45. doi:10.1109/tci.2024.3369396
89. Li H, Wang D, Huang Y, Zhang Y, Yu Z. Generation and recombination for multifocus image fusion with free number of inputs. *IEEE Trans Circuits Syst Video Technology* (2024) 34:6009–23. doi:10.1109/TCSVT.2023.3344222
90. Liu Y, Yu C, Cheng J, Wang ZJ, Chen X. Mm-net: a mixformer-based multi-scale network for anatomical and functional image fusion. *IEEE Trans Image Process* (2024) 33:2197–212. doi:10.1109/tip.2024.3374072



OPEN ACCESS

EDITED BY

Zhiqin Zhu,
Chongqing University of Posts and
Telecommunications, China

REVIEWED BY

Guangcheng Wang,
Nantong University, China
Venu Allapakam,
VIT University, India

*CORRESPONDENCE

Lu Tang,
✉ xztanglu@xzhmu.edu.cn

RECEIVED 06 March 2025

ACCEPTED 13 May 2025

PUBLISHED 26 May 2025

CITATION

Tian C, Zhang J and Tang L (2025) Perceptual
objective evaluation for multimodal medical
image fusion.

Front. Phys. 13:1588508.

doi: 10.3389/fphy.2025.1588508

COPYRIGHT

© 2025 Tian, Zhang and Tang. This is an
open-access article distributed under the
terms of the [Creative Commons Attribution
License \(CC BY\)](https://creativecommons.org/licenses/by/4.0/). The use, distribution or
reproduction in other forums is permitted,
provided the original author(s) and the
copyright owner(s) are credited and that the
original publication in this journal is cited, in
accordance with accepted academic practice.
No use, distribution or reproduction is
permitted which does not comply with
these terms.

Perceptual objective evaluation for multimodal medical image fusion

Chuangeng Tian¹, Juyuan Zhang¹ and Lu Tang^{2*}

¹School of Information Engineering (School of Big Data), Xuzhou University of Technology, Xuzhou,
China, ²School of Medical Imaging, Xuzhou Medical University, Xuzhou, China

Multimodal medical Image fusion (MMIF) has received widespread attention due to its promising application in clinical diagnostics and treatment. Due to the inherent limitations of fusion algorithms, the quality of obtained medical fused images (MFI) varies significantly. An objective evaluation of MMIF can quantify the visual quality differences in fused images and facilitate the rapid development of advanced MMIF techniques, thereby enhancing fused image quality. However, rare research has been dedicated to the MMIF objective evaluation. In this study, we present a multi-scale aware attention network for MMIF quality evaluation. Specifically, we employ a Multi-scale Transform structure that simultaneously processes these multi-scale images using an ImageNet pre-trained ResNet34. Subsequently, we incorporate an online class activation mapping mechanism to focus visual attention on the lesion region, enhancing representative discrepancy features closely associated with MFI quality. Finally, we aggregate these enhanced features and map them to the quality difference. Due to the lack of dataset for the objective evaluation task, we collect 129 pairs of source images from public datasets, namely, the Whole Brain Atlas, and construct a MMIF quality database containing 1,290 medical fused images generated using MMIF algorithms. Each fused image was annotated with a subjective quality score by experienced radiologists. Experimental results demonstrate that our method produces a satisfactory consistent with subjective perception, superior to the state-of-the-art quality evaluation methods. The source images dataset is publicly available at: <http://www.med.harvard.edu/AANLIB/home.html>.

KEYWORDS

multimodal medical image fusion, objective evaluation, multi-scale transform, class activation mapping mechanism, region of interest

1 Introduction

Multimodal medical image fusion (MMIF) is increasingly common in clinical diagnostics. MMIF algorithms aim to generate high-quality fused images from multimodal input images [1–3]. However, most existing MMIF algorithms struggle to achieve optimal fusion due to inherent model limitations. Even worse, instead of promoting, fused image quality declined during the fusion process, even increasing the risk of misdiagnosis. Figure 1 illustrates fusion results from different MMIF algorithms, where the first four images exhibit lower quality compared to the last one, with the first image being the worst. As observed, low-quality fused images fail to convey the critical information of the original images, contradicting the very

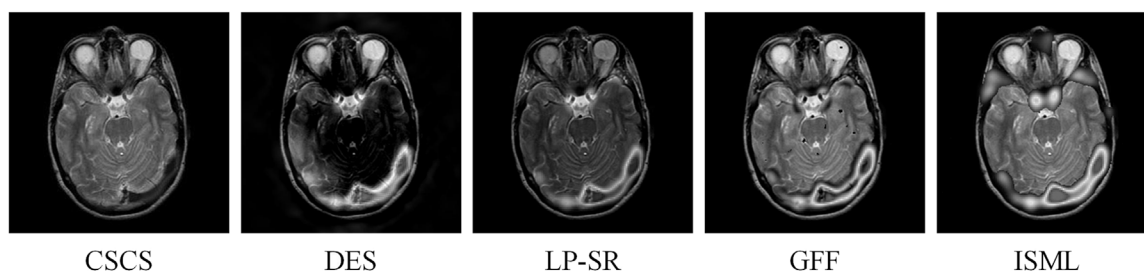


FIGURE 1
A case of fused images via different Multimodal medical image fusion (MMIF) algorithms.

purpose of image fusion. Conversely, high-quality fused images provide clinicians with more reliable information, enhancing diagnostic confidence and decision-making. Hence, it is natural to consider how to achieve a fairer evaluation of these fused images.

In previous work, researchers generally compare the fusion results using both subjective and objective assessments [4–8]. Subjective quality evaluation refers to the visual judgment of image quality by human observers based on perceptual impressions, typically using scoring or ranking methods to quantify visual performance [9]. While this approach closely reflects clinical perception, it is labor-intensive and not scalable for large volumes of medical data. To address this limitation, objective quality assessment methods have been extensively developed to automatically evaluate fused images through computational models and algorithms [10–16]. These methods avoid human bias and enable large-scale assessment by quantifying image quality using well-defined criteria. Generally, objective evaluation methods can be classified into full-reference, reduced-reference, and no-reference approaches [17–19]. Since no ground-truth fused images exist, the no-reference approach is the most suitable for this task. This approach is not only more theoretically realistic but also exhibits higher applicability in clinical settings, as physicians are the ultimate beneficiaries of quality evaluation, the results of image quality assessment can vary depending on the scenario (e.g., the presence or absence of lesion regions in the image), leading to potential instability. No reference evaluation algorithms are roughly divided into hand-crafted metrics and deep learning-based metrics. For instance, Yang et al. [11] gauged structural similarity information of fused images. Qu et al. [15] used mutual information to measure fused images. Tang et al. [17] adopted non-subsampled contourlet transform (NSCT) and pulse coupled neural network (PCNN) for medical fusion image evaluation. However, these studies are limited in their ability to effectively capture hand-crafted features. To alleviate this limitation, deep learning-based metrics have been reported for MMIF quality assessment. Tian et al. [20] exploited a generative adversarial network (GAN) to implement objective evaluation of MMFI. However, such models often face criticism for being “black-box” approaches, making it difficult to gain sufficient trust from radiologists.

In this study, we construct a medical image fusion quality dataset and utilize it to evaluate the performance of the proposed MS-ANN model for MMIF quality assessment. We first conduct multi-scale transform to capture different scale information of fused images.

Meanwhile, input these multi-scale images to fine-tuned ImageNet pre-trained ResNet34. Then, we utilize an online class activation mapping mechanism (CAM) to capture visualization attention to the lesion regions, such operation is highly related to radiologists making decisions. Finally, by aggregating the multi-scale streams to complement each other, we obtain richer, enhanced discrepancy features that are subsequently mapped to the quality differences of the fused images.

The key contributions of the proposed MS-AAN are summarized as follows.

- (1) Given the limited research on objective evaluation for MMIF, we propose a no-reference fused image quality assessment method based on a multi-scale aware attention network, termed MS-AAN. MS-AAN not only automatically predicts the quality of fused images but also enhances model interpretability.
- (2) To characterize quality discrepancies in fused images, we capture and aggregate multi-scale features by utilizing multi-scale transformer and ImageNet pre-trained ResNet34. Such multi-scale streams complement each other and can obtain plentiful details of quality discrepancy-related cues.
- (3) To locate lesion clues and enhance feature representation, we propose a CAM attention network, which can pay attention to the lesion regions via generating localization heat maps. It is highly related to radiologists making decisions. In this way, our MS-AAN earns the trust of radiologists.

2 Related work

2.1 Objective evaluation of multimodal medical image fusion

Multimodal medical image fusion (MMIF) plays an important role in clinical diagnostics and treatment. For radiologists, high-quality fused images can enhance diagnostic confidence and aid in follow-up treatment planning. Plenty of MMIF quality evaluation algorithms have been reported. For instance, Xydeas et al. [10] used gradient information from source images to evaluate fused images. Yang et al. [11] gauged structure similarity information of fused images. Li et al. [12] adopted edge information from the source image to the fused image for objective assessment. Zhao et al. [13] proposed phase congruency to evaluate fused images. Zheng et al.

[14] designed perceptual evaluation via a ratio of spatial frequency error. Qu et al. [15] used mutual information to measure fused images. Liu et al. [16] adopted entropy for fused image objective assessment. Tang et al. [17] adopted non-subsampled contourlet transform and pulse coupled neural network for medical fusion image evaluation. However, these handcrafted methods often lack the ability to effectively capture complex representation features. As a result, deep learning-based metrics for MMIF evaluation have attracted much attention. Tian et al. [20] introduced a generative adversarial network to implement MMIF evaluation. Wang et al. [21] proposed a no-reference image quality assessment framework that incorporates an adaptive graph attention module to enhance both local and contextual information. Liu et al. [9] developed a CNN-based multi-focus image fusion quality assessment model using hierarchical semantic features to better capture focus-level details. Additionally, Yue et al. [18] introduced a pyramid-based framework for assessing the quality of retinal images, which improves robustness to various types of distortions commonly found in clinical data. However, such studies often face challenges in addressing the “black-box” nature of the model. This limits the ability to sufficient trust from radiologists. Despite the growing interest in MMIF evaluation, few studies have focused on objective evaluation, and there is a lack of high-quality fused images. As a result, no reference metric demonstrates significant practical value for this task.

2.2 Multi-scale aware network

In recent years, multi-scale transform has achieved progress in the field of multimodal medical image fusion [22, 23], especially non-subsampled contourlet transform (NSCT) has displayed tremendous results [24, 25]. Specifically, Huang et al. [25] proposed SPECT and CT image fusion based on NSCT and PCNN. Yin et al. [24] used NSCT and PCNN for medical image fusion. Tang et al. [17] proposed a medical fusion image evaluation method based on NSCT and PCNN. Therefore, the combination of NSCT and PCNN has been proven to be a highly effective strategy for MMIF and MMIF quality evaluation. Inspired by this, can we replace PCNN with deep learning? Recent advancements in pre-trained CNNs on ImageNet have demonstrated their ability to extract richer features [26–28]. Motivated by the above fact, we employ a simple yet effective approach by combining NSCT with a pre-trained CNN to capture richer multi-scale feature representations.

2.3 CAM attention mechanism

Recent years have witnessed that the CAM is an effective tool for model interpretability. Zhou et al. used CAM to locate class-relevant objects [29]. Subsequently, gradient-weighted CAM was further extended to obtain better localization [30]. Ouyang et al. adopted gradient-weighted CAM to learn chest X-ray abnormality localization [31]. Tang et al. utilized an online CAM mechanism to concentrate on thyroid nodule localization, improving the model interpretability [32]. Thus, in this paper, we further extend the CAM attention mechanism to guide the network in focusing on lesion

regions, enhancing the representative discriminative features, which ensures alignment with radiologists’ decision-making.

3 Methods

The proposed MS-ANN model is designed to comprehensively capture perceptual quality information from multimodal fused medical images. Its architecture comprises three main components: a multi-scale transform module, an ImageNet pre-trained ResNet34 backbone, and a CAM attention mechanism, as illustrated in Figure 2. First, we construct a multi-scale stream network with NSCT by down-sampling the input fused images to generate representations at four different scales. Each scale is processed by four ResNet34 backbone, which is selected for its efficiency and strong feature representation ability. Using a pre-trained model also facilitates robust learning with limited data. To enhance model interpretability and ensure the network emphasizes diagnostically relevant regions, we incorporate a CAM-based attention mechanism after feature extraction. Finally, the attention-refined features from all scales are concatenated and mapped to a quality score through fully connected layers.

3.1 Multi-scale aware neural network

We adopt the NSCT to perform multi-scale and multi-directional decomposition on the medical fused image. NSCT is a shift-invariant extension of the contourlet transform that enables rich representation of image features across different scales and directions, which is particularly beneficial for medical image analysis. Specifically, the medical fused image F is transformed into multiple sub-band $\{F_{m,\alpha}\}$ at each level $m \in [1, 4]$ and direction. This decomposition allows the network to capture structural details at various resolutions, which is formulated as Equation 1:

$$F_{m,\alpha} = MST(F) \quad (1)$$

Where $MST(\cdot)$ repents the MST functions. Following this transformation, we use an ImageNet pre-trained ResNet34 as the backbone to extract high-level semantic features from the decomposed components. Particularly, these multiple sub-bands are input to ResNet34, and we use Rectified Linear Unit (ReLU) as the activation function, which is formulated as Equation 2:

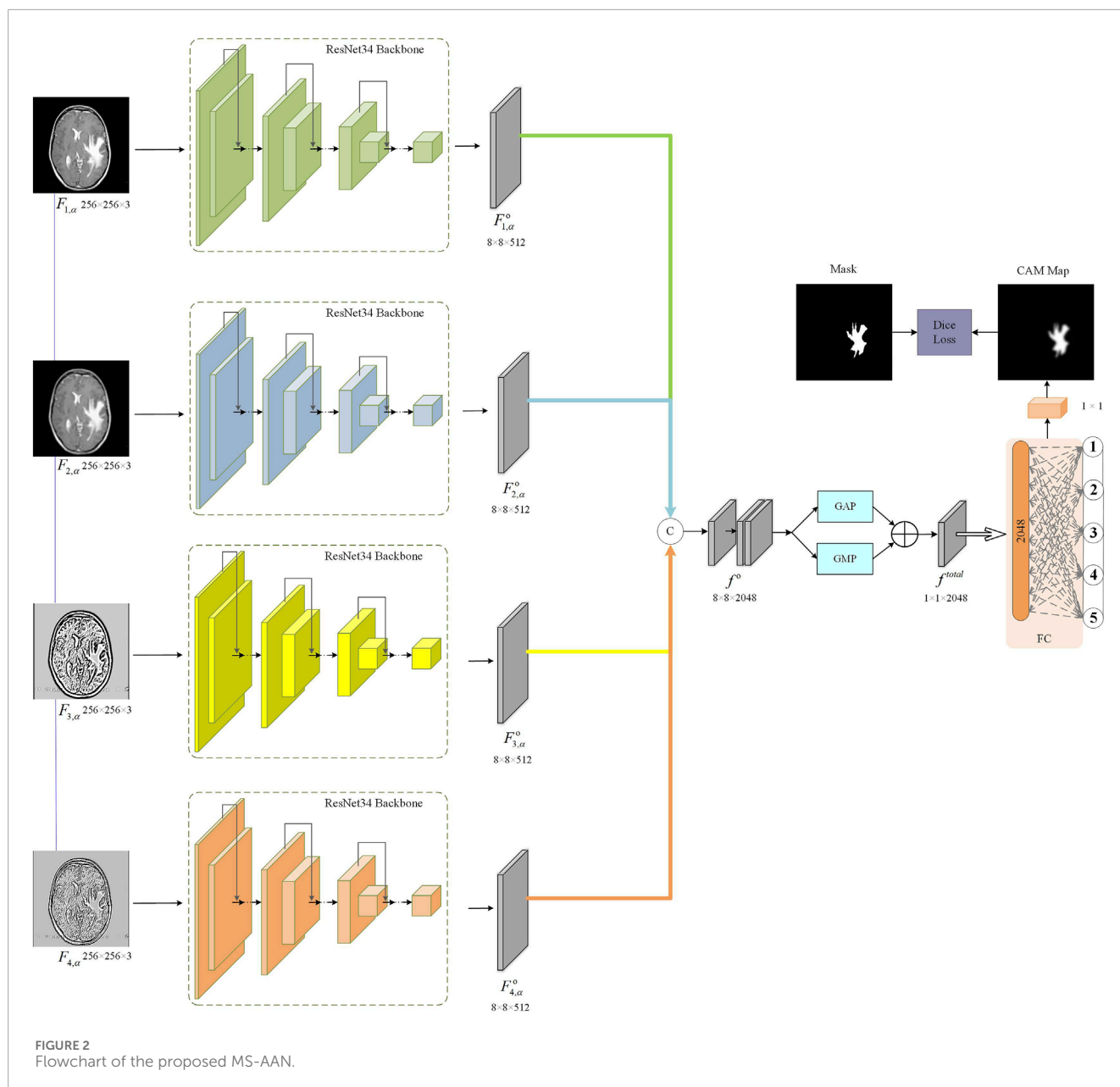
$$F_{m,\alpha}^o = ReLU(conv(F_{m,\alpha}, K)) = ReLU\left(\sum_{n=1}^t F_{m,\alpha}^n \odot K + A\right) \quad (2)$$

Where $F_{m,\alpha}^o$ stands for output features. K represents a kernel of convolutional layer. $F_{m,\alpha}^n$ is n_{th} channel of $F_{m,\alpha}$ with totally t channels, A and \odot represent the bias and convolution operation, respectively.

3.2 Aggregation of multi-scale feature

Considering the advantages of multi-scale transform, we aggregate the output features of multi-scale streams for MMIF quality evaluation. Firstly, we perform concatenate operations on four multi-scale stages, as shown in Equation 3:

$$f^o = F_{1,\alpha}^o \oplus F_{2,\alpha}^o \oplus F_{3,\alpha}^o \oplus F_{4,\alpha}^o \quad (3)$$



where \oplus stands for concatenate operation. Then, we compute global average pooling (GAP), as shown in Equation 4:

$$I_G = \frac{1}{W_{I_c} * H_{I_c}} \sum_{j=1}^{W_{I_c} * H_{I_c}} I_c^j \quad (4)$$

Where I_c^j denotes the pixel value of j -th in I_c , I_c stands for output of the last layer. W_{I_c} and H_{I_c} represent the width and height of I_c , respectively. The enhancement feature transfers to a convolution layer, and we conduct GAP and global max pooling (GMP). Finally, a simple addition operation is carried out to aggregate GAP and GMP, which is formulated as Equation 5:

$$f^{total} = GAP(f^o) + GMP(f^o) \quad (5)$$

3.3 CAM attention mechanism

To capture quality discrepancy features of lesion region from the whole medical fused images, we introduce the CAM attention mechanism. Specifically, we generate the attention feature map M by applying a nonlinear activation function to the final aggregated feature map f^{total} , which is described in Equation 5. This representation integrates multi-scale semantic information and is more suitable for highlighting perceptually important regions. The resulting attention map has a spatial resolution of $1/16$ relative to the input image $\{F_{m,\alpha}\}$ and guides the network to focus on diagnostically relevant areas during quality assessment. Then, conducting a normalization on M to $[0, 1]$. After that, performing the sigmoid operation for soft masking, named $S(M)$, is formulated as Equation 6:

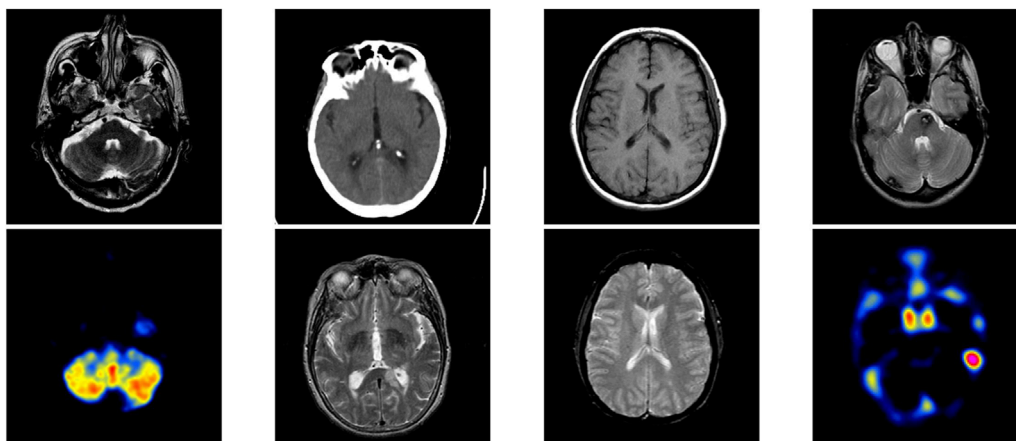


FIGURE 3
Some examples of source images.

TABLE 1 Comparison performance of MS-ANN with other six metrics.

Metric	PLCC	SRCC	KRCC	RMSE
MPRI	0.3031	0.3167	0.2375	0.2611
TE	0.1797	0.1946	0.1407	0.3909
MI	0.2270	0.1738	0.1071	0.3712
OEEP	0.3064	0.3367	0.2342	0.2810
NSCT-PCNN	0.6252	0.6420	0.4166	0.2480
RSFE	0.4054	0.2275	0.1700	0.2663
AGA	0.6956	0.6871	0.5721	0.3669
SBA	0.5861	0.6012	0.4156	0.4266
PNQC	0.8106	0.8016	0.7681	0.2119
Proposed MS-ANN	0.9131	0.9061	0.8560	0.1166

Bold values represent the best results.

$$S(M) = \frac{1}{1 + \exp(-\mu(M - \beta))} \quad (6)$$

Where μ and β stand for hyper-parameters. Dice loss is used as the attention loss function, denoted as L_a , and is defined as shown in Equation 7:

$$L_a = \text{Dice}(S(M), G) \quad (7)$$

Where G is the ground truth of the lesion mask. Finally, in the fully connected layer, we conduct Cross Entropy loss for quality classification, dubbed L_c , as shown in Equation 8:

$$L_c = - \sum [f \log(\hat{f}_x) - (1 - f)(1 - \log \hat{f}_x)] \quad (8)$$

Where f stands for class label, $\hat{f}_x = [\hat{f}_1, \hat{f}_2, \hat{f}_3, \hat{f}_4, \hat{f}_5]$, $x = 1, 2, 3, 4, 5$, which denote the five classes quality results of medical fused images.

3.4 Total loss function

As observe in Figure 2, total loss function of our MS-ANN, comprise of attention L_a and classification L_c , which is denoted as shown in Equation 9:

$$L_t = L_a + \gamma L_c \quad (9)$$

4 Experiments

4.1 Dataset

In this study, we perform medical fused data for appraising the developed MS-ANN in MMIF quality assessment. Specifically, we collect 129 pairs of source images from public datasets, i.e., Whole Brain Atlas, which include CT and MR, MR-T1 and MR-T2, MR-T2 and PET, MR-T2 and SPECT, as shown in Figure 3. The selected images span a wide range of anatomical structures and clinical conditions (e.g., tumors, lesions, and degenerative changes), ensuring that the dataset is both diverse and representative of real-world clinical fusion scenarios. We then apply ten representative state-of-the-art MMIF algorithms [16, 24, 33–40], resulting in a total of 1,290 fused images. This dataset construction process is consistent with our previous work, where more technical details of the fusion methods can be found [20, 41]. For subjective quality assessment, each fused image is annotated with a Mean Opinion Score (MOS) ranging from 1 (lowest quality) to 5 (highest quality), as independently rated by two experienced radiologists. To ensure the reliability and consistency of the subjective assessment, a senior radiologist further reviewed and validated the assigned scores.

To rigorously evaluate the effectiveness of the proposed MS-ANN, we adopt four widely recognized quantitative assessment metrics [42, 43]: Pearson's Linear Correlation Coefficient (PLCC), Spearman's Rank Correlation Coefficient (SRCC), Kendall's Rank Correlation Coefficient (KRCC), and Root Mean Square Error (RMSE). These metrics are designed to measure the alignment

TABLE 2 Ablation studies on the proposed MS-ANN.

Model	Pre	Multi-scale	CAM	PLCC	SRCC	KRCC	RMSE
Baseline				0.7971	0.8022	0.7199	0.2936
Proposed ResNet34	★			0.8633	0.8571	0.7761	0.1696
Proposed ResNet34	★	★		0.8916	0.8811	0.8256	0.1301
Proposed ResNet34	★	★	★	0.9131	0.9061	0.8560	0.1166

Bold values represent the best results.

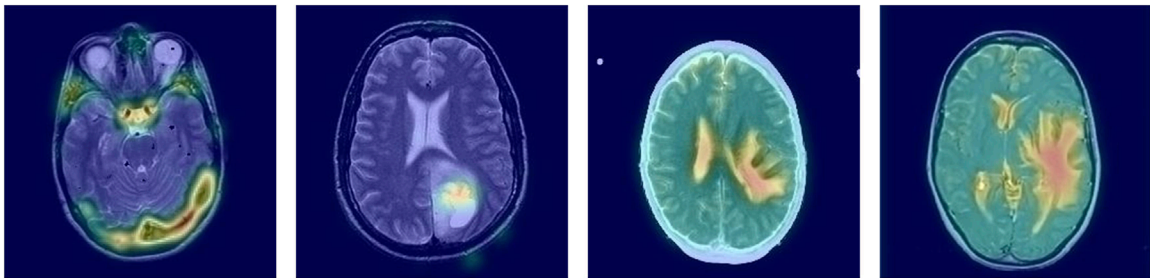


FIGURE 4
Generated attention maps of our methods on four medical fused images.

TABLE 3 The result of the external validation.

Model	PLCC	SRCC	KRCC	RMSE
Our	0.8591	0.8388	0.7916	0.1721

between the predicted quality scores generated by the model and the ground-truth MOS provided by expert radiologists. Specifically, PLCC, SRCC, and KRCC are used to evaluate the consistency between the predicted quality scores and the ground-truth MOS, with higher values indicating better consistency with human perception. RMSE measures the absolute prediction error, where lower values represent better performance. These metrics are widely used in the field and ensure comparability with previous IQA studies [9, 18, 19, 21, 44].

4.2 Performance comparison

To validate the effectiveness of the proposed MS-ANN, we compare it with six mainstream methods, including multiple pseudo reference images-based quality metric (MPRI) [44], Tsallis entropy-based quality metric (TE) [45], mutual information-based quality metric (MI) [46], the objective evaluation of fusion performance (OEF) [10], the ratio of spatial frequency error-based quality metric (RSFE) [14], the NSCT-PCNN-based quality metric (NSCT-PCNN) [17], the adaptive graph attention (AGA) for blind image quality assessment method [21], statistically based approach (SBA) for multi-focus image fusion quality assessment [9], and pyramid networks with quality-aware contrast loss (PNQC) for retinal image

quality assessment [18]. Among these metrics, higher values of MPRI, TE, MI, OEF, NSCT-PCNN, AGA, SBA, and PNQC indicate better quality, whereas lower values of RSFE denote better quality. We compute the PLCC, SRCC, KRCC and RMSE values of six mainstream methods and MS-ANN, as shown in Table 1. The highest scores are highlighted in bold. Based on Table 1, our MS-ANN achieves the best performance, significantly outperforming the six competing models. Specifically, compared to the second-ranked RIQA, our proposed method improves PLCC from 0.8106 to 0.9131, SRCC from 0.8016 to 0.9061, KRCC from 0.7681 to 0.8560, while declining RMSE 0.2119 from to 0.1166.

4.3 Ablation study

We conduct ablation studies to discuss the contribution of each important part of the MS-ANN. We first train each component independently on the medical fused dataset and then jointly optimize all components of MS-ANN. The results are presented in Table 2. First, the baseline model refers to ResNet34 without ImageNet pre-training, achieving a PLCC of 0.7971, SRCC of 0.8022, KRCC of 0.7199, and RMSE of 0.2936. Second, we apply a pre-training strategy to enhance the ability to capture features. As shown in the second row of Table 2, performance significantly improves, with PLCC increasing from 0.7971 to 0.8633, SRCC from 0.8022 to 0.8571, and KRCC from 0.7199 to 0.7761, while RMSE decreases from 0.2936 to 0.1696. These results demonstrate that the ImageNet pre-trained model outperforms the baseline model without pre-training. This improvement may be attributed to the effective use of pre-trained knowledge, which helps mitigate the challenge of limited

training data. Third, we further introduce NSCT to capture more multi-scale features. With the addition of multi-scale transform, the results show noticeable improvements when comparing baseline + Pre and baseline + Pre + multi-scale: PLCC increases by 2.83% (0.8633 vs. 0.8916), SRCC by 2.40% (0.8571 vs. 0.8811), and KRCC by 4.95% (0.7761 vs. 0.8256), while RMSE decreases by 3.95% (0.1696 vs. 0.1301). Moreover, we integrate the CAM mechanism to guide the model's attention toward lesion regions, thereby enhancing both feature representation and interpretability. As shown in Table 2, the proposed MS-ANN (Baseline + Pre + multi-scale + CAM) achieves superior performance compared to the variant without CAM (Baseline + Pre + multi-scale). Specifically, PLCC increases from 0.8916 to 0.9131, SRCC from 0.8811 to 0.9061, KRCC from 0.8256 to 0.8560, and RMSE decreases from 0.1301 to 0.1166. These improvements demonstrate that CAM significantly enhances the model's ability to capture quality-related features. More importantly, the lesion-focused attention maps provide intuitive visual explanations, which can assist radiologists in verifying model predictions and build greater confidence in clinical use. As shown in Figure 4, the CAM-based heatmaps illustrate the model's ability to concentrate on diagnostically relevant regions, offering visual support for the model's quantitative superiority.

4.4 External validation

To further validate the generalization ability of our MS-ANN, we conduct an external independent evaluation using the multimodal medical image fusion database [17]. It is important to note that the performance metrics reported in Table 3 differ from those in Table 2 because they are obtained under different evaluation settings. Specifically, Table 2 reports results from ablation studies conducted on the training dataset to analyze the contribution of each model component, whereas Table 3 presents results from a separate external dataset. As shown in Table 3, our model achieves promising performance, with a PLCC of 0.8591, SRCC of 0.8388, KRCC of 0.7916, and RMSE of 0.1721. These results demonstrate the robustness and effectiveness of MS-ANN in assessing multimodal medical image fusion quality across different datasets.

5 Conclusion

In this paper, we develop a quality evaluation metric for multimodal medical image fusion, called no reference multi-scale aware attention network (MS-ANN). Specifically, we first apply a multi-scale transform to extract different scale information from fused images and feed these transformed images into an ImageNet pre-trained ResNet34. This multi-scale strategy enables complementary feature extraction, capturing rich details relevant to quality assessment. Then, we propose a CAM attention network, which captures visualization attention to the lesion regions to facilitate model interpretability. Finally, we employ a concatenation operation to refine quality discrepancy features and map them to the quality differences in multimodal fusion images. However, the dataset used in this study exhibits an imbalance between MRI-PET and MRI-SPECT image pairs, with MRI-SPECT images being more prevalent. Moreover, the diversity of medical conditions

and anatomical regions is somewhat limited, which may affect the model's generalization to other clinical settings or imaging modalities. In future work, we aim to address these limitations by expanding the dataset to cover a broader range of organs and clinical conditions, thereby improving the robustness and generalization capability of the proposed MS-ANN model. Additionally, while our study adopts widely accepted statistical metrics to evaluate image quality prediction, it is important to recognize the potential influence of MMIF quality on downstream clinical tasks such as diagnosis accuracy or treatment decisions. High-quality fused images can provide clearer lesion boundaries, improved structural detail, and more reliable functional information, which are crucial in radiological assessment and therapy planning. In future work, we intend to design user studies or integrate radiologist-in-the-loop evaluations to measure the actual diagnostic utility of images rated by our model. Such assessments would offer a more comprehensive validation of the model's clinical value and help bridge the gap between objective image quality assessment and practical medical outcomes. Despite these limitations, the proposed MS-ANN shows strong consistency with subjective perception, offering potential to facilitate clinical diagnosis and guide the development of advanced multimodal medical image fusion techniques.

Data availability statement

The original contributions presented in the study are included in the article/supplementary material, further inquiries can be directed to the corresponding author.

Author contributions

CT: Conceptualization, Methodology, Supervision, Validation, Writing – original draft. JZ: Methodology, Validation, Writing – original draft. LT: Conceptualization, Project administration, Supervision, Validation, Writing – original draft, Writing – review and editing.

Funding

The author(s) declare that financial support was received for the research and/or publication of this article. This work was supported by the Xu Zhou Science and technology Program, China (KC22466) and National Natural Science Foundation of China (82001912).

Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Generative AI statement

The author(s) declare that no Generative AI was used in the creation of this manuscript.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated

organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

References

- Birkfellner W, Figl M, Furtado H, Renner A, Hatamikia S, Hummel J. Multi-modality imaging: a software fusion and image-guided therapy perspective. *Front Phys* (2018) 6:66. doi:10.3389/fphy.2018.00066
- Azam MA, Khan KB, Salahuddin S, Rehman E, Khan SA, Khan MA. A review on multimodal medical image fusion: compendious analysis of medical modalities, multimodal databases, fusion techniques and quality metrics. *Comput Biol Med* (2022) 144:105253. doi:10.1016/j.combiomed.2022.105253
- Zhou T, Cheng Q, Lu H, Li Q, Zhang X, Qiu S. Deep learning methods for medical image fusion: a review. *Comput Biol Med* (2023) 160:106959. doi:10.1016/j.combiomed.2023.106959
- Cheng S, Liu R, He Y, Fan X, Luo Z. Blind image deblurring via hybrid deep priors modeling. *Neurocomputing* (2020) 387:334–45. doi:10.1016/j.neucom.2020.01.004
- Shao W-Z, Lin Y-Z, Liu Y-Y, Wang L-Q, Ge Q, Bao B-K, et al. Gradient-based discriminative modeling for blind image deblurring. *Neurocomputing* (2020) 413:305–27. doi:10.1016/j.neucom.2020.06.093
- Wang Z, Bovik AC, Sheikh HR, Simoncelli EP. Image quality assessment: from error visibility to structural similarity. *IEEE Trans Image Process* (2004) 13:600–12. doi:10.1109/TIP.2003.819861
- Shen L, Chen X, Pan Z, Fan K, Li F, Lei J. No-reference stereoscopic image quality assessment based on global and local content characteristics. *Neurocomputing* (2021) 424:132–42. doi:10.1016/j.neucom.2020.10.024
- Ma K, Liu W, Zhang K, Duanmu Z, Wang Z, Zuo W. End-to-End blind image quality assessment using deep neural networks. *IEEE Trans Image Process* (2018) 27:1202–13. doi:10.1109/TIP.2017.2774045
- Liu Y, Qi Z, Cheng J, Chen X. Rethinking the effectiveness of objective evaluation metrics in multi-focus image fusion: a statistic-based approach. *IEEE Trans Pattern Anal Mach Intell* (2024) 46:5806–19. doi:10.1109/TPAMI.2024.3367905
- Ydeas CS, Petrović V. Objective image fusion performance measure. *Electron Lett* (2000) 36:308–9. doi:10.1049/el:20000267
- Yang C, Zhang J-Q, Wang X-R, Liu X. A novel similarity based quality metric for image fusion. *Inf Fusion* (2008) 9:156–60. doi:10.1016/j.inffus.2006.09.001
- Li S, Kwok JT, Wang Y. Combination of images with diverse focuses using the spatial frequency. *Inf Fusion* (2001) 2:169–76. doi:10.1016/S1566-2535(01)00038-0
- Zhao J, Laganieri R, Liu Z. Performance assessment of combinative pixel-level image fusion based on an absolute feature measurement. *Int J Innovat Comput Inf Control* (2006) 3. doi:10.1109/ICICIC.2006.296
- Zheng Y, Essock EA, Hansen BC, Haun AM. A new metric based on extended spatial frequency and its application to DWT based fusion algorithms. *Inf Fusion* (2007) 8:177–92. doi:10.1016/j.inffus.2005.04.003
- Qu G, Zhang D, Yan P. Information measure for performance of image fusion. *Electron Lett* (2002) 38:313–5. doi:10.1049/el:20020212
- Liu Y, Liu S, Wang Z. A general framework for image fusion based on multi-scale transform and sparse representation. *Inf Fusion* (2015) 24:147–64. doi:10.1016/j.inffus.2014.09.004
- Tang L, Tian C, Li L, Hu B, Yu W, Xu K. Perceptual quality assessment for multimodal medical image fusion. *Signal Process Image Commun* (2020) 85:115852. doi:10.1016/j.image.2020.115852
- Yue G, Zhang S, Zhou T, Jiang B, Liu W, Wang T. Pyramid network with quality-aware contrastive loss for retinal image quality assessment. *IEEE Trans Med Imaging* (2025) 44:1416–31. doi:10.1109/TMI.2024.3501405
- Guo Y, Hu M, Min X, Wang Y, Dai M, Zhai G, et al. Blind image quality assessment for pathological microscopic image under screen and immersion scenarios. *IEEE Trans Med Imaging* (2023) 42:3295–306. doi:10.1109/TMI.2023.3282387
- Tian C, Zhang L. G2NPAN: GAN-guided nuance perceptual attention network for multimodal medical fusion image quality assessment. *Front Neurosci* (2024) 18:1415679. doi:10.3389/fnins.2024.1415679
- Wang H, Liu J, Tan H, Lou J, Liu X, Zhou W, et al. Blind image quality assessment via adaptive graph attention. *IEEE Trans Circuits Syst Video Technol* (2024) 34:10299–309. doi:10.1109/TCSVT.2024.3405789
- Duan H, Wang W, Xing L, Xie B, Zhang Q, Zhang Y. Identifying geological structures in the Pamir region using non-subsampled shearlet transform and gravity gradient tensor. *Geophys J Int* (2025) 240:2125–43. doi:10.1093/gji/ggaf036
- Ma J, Chen Y, Chen L, Tang Z. Dual-attention pyramid transformer network for no-reference image quality assessment. *Expert Syst Appl* (2024) 257:125008. doi:10.1016/j.eswa.2024.125008
- Yin M, Liu X, Liu Y, Chen X. Medical image fusion with parameter-adaptive pulse coupled neural network in nonsubsampling shearlet transform domain. *IEEE Trans Instrum Meas* (2019) 68:49–64. doi:10.1109/TIM.2018.2838778
- Huang C, Tian G, Lan Y, Peng Y, Ng EYK, Hao Y, et al. A new pulse coupled neural network (PCNN) for brain medical image fusion empowered by shuffled frog leaping algorithm. *Front Neurosci* (2019) 13:210. doi:10.3389/fnins.2019.00210
- Norouzi M, Hosseini SH, Khoshnevisan M, Moshiri B. Applications of pre-trained CNN models and data fusion techniques in Unity3D for connected vehicles. *Appl Intell* (2025) 55:390. doi:10.1007/s10489-024-06213-3
- Swamy MR, P V, Rajendran V. Deep learning approaches for online signature authentication: a comparative study of pre-trained CNN models. *Eng Res Express* (2025) 7:015230. doi:10.1088/2631-8695/ada86d
- Arnia F, Saddami K, Roslilar R, Muharrar R, Munadi K. Towards accurate diabetic foot ulcer image classification: leveraging CNN pre-trained features and extreme learning machine. *Smart Health* (2024) 33:100502. doi:10.1016/j.smhl.2024.100502
- Zhou B, Khosla A, Lapedriza A, Oliva A, Torralba A. Learning deep features for discriminative localization. In: *2016 IEEE conference on computer vision and pattern recognition (CVPR)*. Las Vegas, NV: IEEE (2016). p. 2921–9. doi:10.1109/CVPR.2016.319
- Selvaraju RR, Cogswell M, Das A, Vedantam R, Parikh D, Batra D. Grad-CAM: visual explanations from deep networks via gradient-based localization. Venice: IEEE: 2017 IEEE International Conference on Computer Vision ICCV (2017). p. 618–26. doi:10.1109/ICCV.2017.74
- Ouyang X, Karanam S, Wu Z, Chen T, Huo J, Zhou XS, et al. Learning hierarchical attention for weakly-supervised chest X-ray abnormality localization and diagnosis. *IEEE Trans Med Imaging* (2021) 40:2698–710. doi:10.1109/TMI.2020.3042773
- Tang L, Tian C, Yang H, Cui Z, Hui Y, Xu K, et al. TS-DSANN: texture and shape focused dual-stream attention neural network for benign-malignant diagnosis of thyroid nodules in ultrasound images. *Med Image Anal* (2023) 89:102905. doi:10.1016/j.media.2023.102905
- Min X, Zhai G, Gu K, Yang X, Guan X. Objective quality evaluation of dehazed images. *IEEE Trans Intell Transport Syst* (2019) 20:2879–92. doi:10.1109/TITS.2018.2868771
- Liu Y, Chen X, Ward RK, Jane Wang Z. Image fusion with convolutional sparse representation. *IEEE Signal Process Lett* (2016) 23:1882–6. doi:10.1109/LSP.2016.2618776
- Das S, Kundu MK. NSCT-based multimodal medical image fusion using pulse-coupled neural network and modified spatial frequency. *Med Biol Eng Comput* (2012) 50:1105–14. doi:10.1007/s11517-012-0943-3
- Li S, Kang X, Hu J. Image fusion with guided filtering. *IEEE Trans Image Process* (2013) 22:2864–75. doi:10.1109/TIP.2013.2244222
- Shen R, Cheng I, Basu A. Cross-scale coefficient selection for volumetric medical image fusion. *IEEE Trans Biomed Eng* (2013) 60:1069–79. doi:10.1109/TBME.2012.2211017
- Du J, Li W, Xiao B, Nawaz Q. Union Laplacian pyramid with multiple features for medical image fusion. *Neurocomputing* (2016) 194:326–39. doi:10.1016/j.neucom.2016.02.047
- Tang L, Tian C, Xu K. Exploiting quality-guided adaptive optimization for fusing multimodal medical images. *IEEE Access* (2019) 7:96048–59. doi:10.1109/ACCESS.2019.2926833
- Das S, Kundu MK. A neuro-fuzzy approach for medical image fusion. *IEEE Trans Biomed Eng* (2013) 60:3347–53. doi:10.1109/TBME.2013.2282461
- Tang L, Hui Y, Yang H, Zhao Y, Tian C. Medical image fusion quality assessment based on conditional generative adversarial network. *Front Neurosci* (2022) 16:986153. doi:10.3389/fnins.2022.986153
- Hu B, Wang S, Gao X, Li L, Gan J, Nie X. Reduced-reference image deblurring quality assessment based on multi-scale feature enhancement and aggregation. *Neurocomputing* (2023) 547:126378. doi:10.1016/j.neucom.2023.126378

43. Sim K, Yang J, Lu W, Gao X. Blind stereoscopic image quality evaluator based on binocular semantic and quality channels. *IEEE Trans Multimedia* (2022) 24:1389–98. doi:10.1109/TMM.2021.3064240
44. Min X, Zhai G, Gu K, Liu Y, Yang X. Blind image quality estimation via distortion aggravation. *IEEE Trans Broadcast* (2018) 64:508–17. doi:10.1109/TBC.2018.2816783
45. Sholehkerdar A, Tavakoli J, Liu Z. In-depth analysis of Tsallis entropy-based measures for image fusion quality assessment. *Opt Eng* (2019) 58:1. doi:10.1117/1.OE.58.3.033102
46. Hossny M, Nahavandi S, Creighton D. Comments on 'Information measure for performance of image fusion'. *Electron Lett* (2008) 44:1066–7. doi:10.1049/el:20081754



OPEN ACCESS

EDITED BY

Zhiqin Zhu,
Chongqing University of Posts and
Telecommunications, China

REVIEWED BY

Fan Li,
Kunming University of Science and
Technology, China
Haicheng Bai,
Yunnan Normal University, China

*CORRESPONDENCE

Zheng Liu,
✉ 490956823@qq.com

RECEIVED 25 March 2025

ACCEPTED 19 May 2025

PUBLISHED 06 June 2025

CITATION

Hu D, Wang K, Zhang C, Liu Z, Che Y, Dong S
and Kong C (2025) Target-aware unregistered
infrared and visible image fusion.
Front. Phys. 13:1599968.
doi: 10.3389/fphy.2025.1599968

COPYRIGHT

© 2025 Hu, Wang, Zhang, Liu, Che, Dong and
Kong. This is an open-access article
distributed under the terms of the [Creative
Commons Attribution License \(CC BY\)](#). The
use, distribution or reproduction in other
forums is permitted, provided the original
author(s) and the copyright owner(s) are
credited and that the original publication in
this journal is cited, in accordance with
accepted academic practice. No use,
distribution or reproduction is permitted
which does not comply with these terms.

Target-aware unregistered infrared and visible image fusion

Dengshu Hu, Ke Wang, Cuijin Zhang, Zheng Liu*, Yukui Che,
Shoubing Dong and Chuirui Kong

Qujing Power Supply Bureau, Yunnan Power Grid Co., Ltd., Qujing, China

Introduction: Infrared (IR) and visible (VI) image fusion can provide richer texture details for subsequent object detection tasks. Conversely, object detection can offer semantic information about targets, which in turn helps improve the quality of the fused images. As a result, joint learning approaches that integrate infrared-visible image fusion and object detection have attracted increasing attention.

Methods: However, existing methods typically assume that the input source images are perfectly aligned spatially—an assumption that does not hold in real-world applications. To address this issue, we propose a novel method that enables mutual enhancement between infrared-visible image fusion and object detection, specifically designed to handle misaligned source images. The core idea is to use the object detection loss, propagated via backpropagation, to guide the training of the fusion network, while a specially designed loss function mitigates the modality gap between infrared and visible images.

Results: Comprehensive experiments on three public datasets demonstrate the effectiveness of our approach.

Discussion: In addition, our approach can be used with other radiation frequencies where different modalities require image fusion like, for example, radio-frequency, x- and gamma rays used in medical imaging.

KEYWORDS

infrared and visible image fusion, object detection, feature alignment, target-aware, unregistered

1 Introduction

Images captured by a single sensor often fail to provide a comprehensive description of a scene. For example, infrared (IR) sensors can capture thermal radiation emitted by objects and highlight salient targets, but they lack the ability to represent fine texture details and are more susceptible to noise. On the other hand, visible-light (VI) sensors capture visual information with clear texture details but are easily affected by lighting conditions and occlusions. If the information from both infrared and visible images can be integrated into a single, information-rich fused image, the scene representation can be significantly enhanced. As a result, infrared and visible image fusion has been widely applied as a low-level preprocessing task in various high-level vision applications, such as object detection [1], tracking [2], person re-identification [3], and semantic segmentation [4]. An example in [Figure 1](#) visually illustrates the application of fused images in object detection. It can be observed that detection results obtained from individual sensor images are less accurate than those derived from fused images.

Due to its practical value, infrared and visible image fusion has garnered substantial attention in the research community. Over the past decades, numerous image fusion

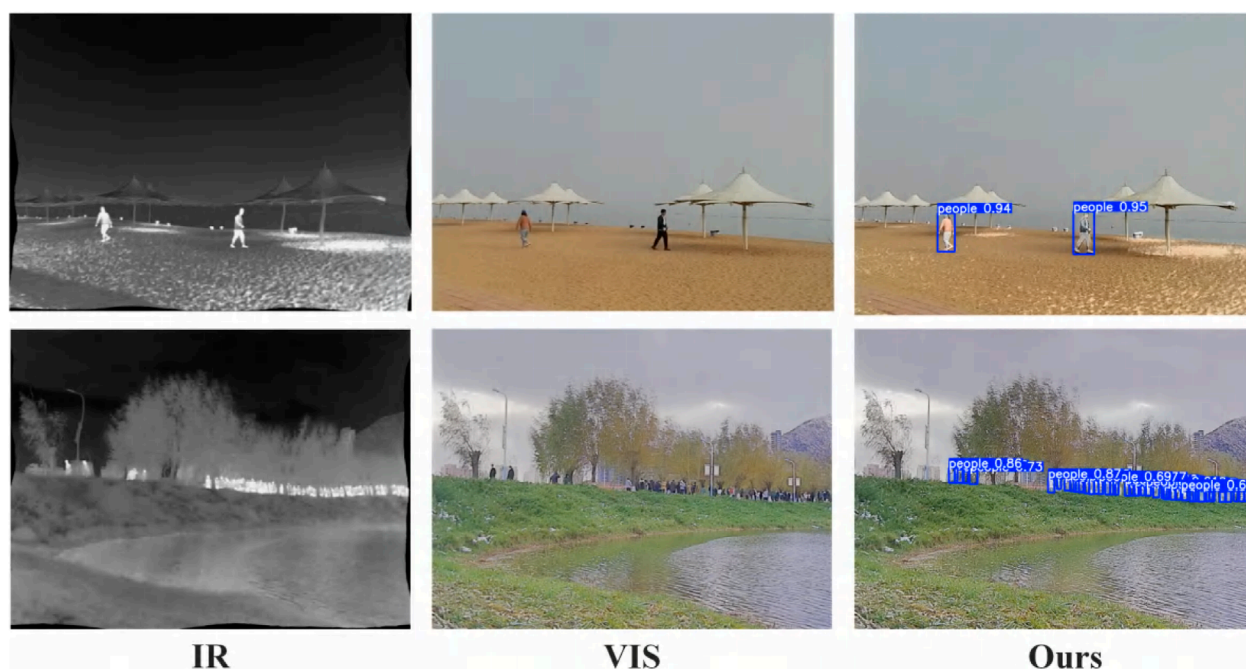


FIGURE 1
Object detection results of the proposed method on the M³FD dataset.

techniques have been proposed, including both traditional and deep learning-based methods. Traditional methods typically fall into two categories: multi-scale transform-based methods [5–7] and sparse representation-based methods [8–12]. Deep learning-based approaches include methods based on autoencoders (AE) [9, 13, 14], convolutional neural networks (CNNs) [15–18], and generative adversarial networks (GANs) [19, 20].

Although recent deep learning-based fusion algorithms can generate visually pleasing results, several critical challenges remain unsolved. On one hand, most existing fusion algorithms focus on optimizing visual quality and evaluation metrics, but rarely consider whether the fused results benefit downstream task performance. On the other hand, even recent methods that incorporate high-level vision tasks into the fusion process—such as TarDAL [21], which proposes a dual-level optimization model using a task-aware dual adversarial learning network to simultaneously address fusion and object detection; SeAFusion [22], which constrains the fusion process with semantic loss to retain richer semantic information; and DetFusion [23], which guides multimodal fusion using target-related features learned by the object detection network—still assume that the source images are perfectly aligned spatially. This assumption does not hold in real-world applications.

In this study, we propose a framework named Target-Aware Unregistered Infrared and Visible Image Fusion Network, designed to achieve robust performance in both misaligned image fusion and high-level vision tasks. Specifically, we introduce an object detection network to predict detection results on the fused image and construct a detection loss. This loss is then backpropagated to guide the training of the fusion network, encouraging the fused image to retain more information useful for object detection. Additionally, to effectively align unregistered images, we design

a modality consistency loss to reduce the domain gap between infrared and visible images.

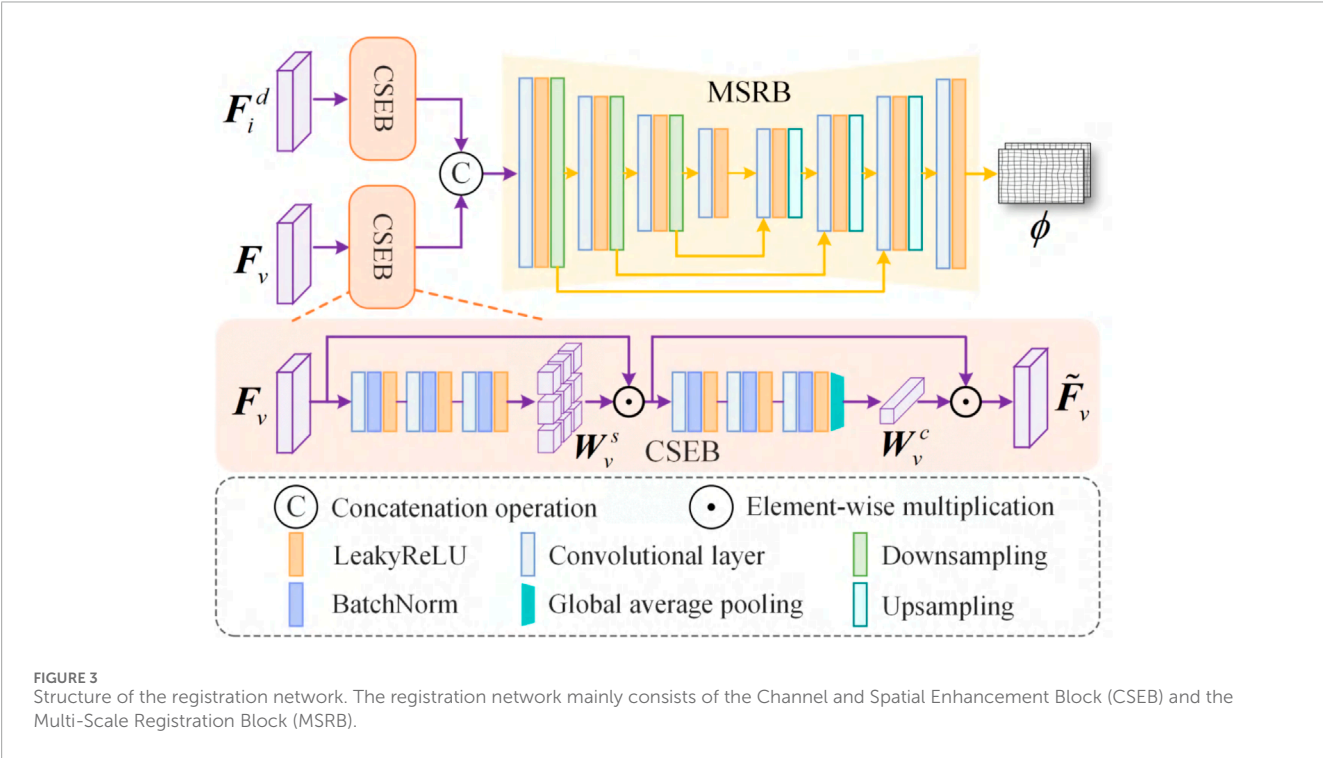
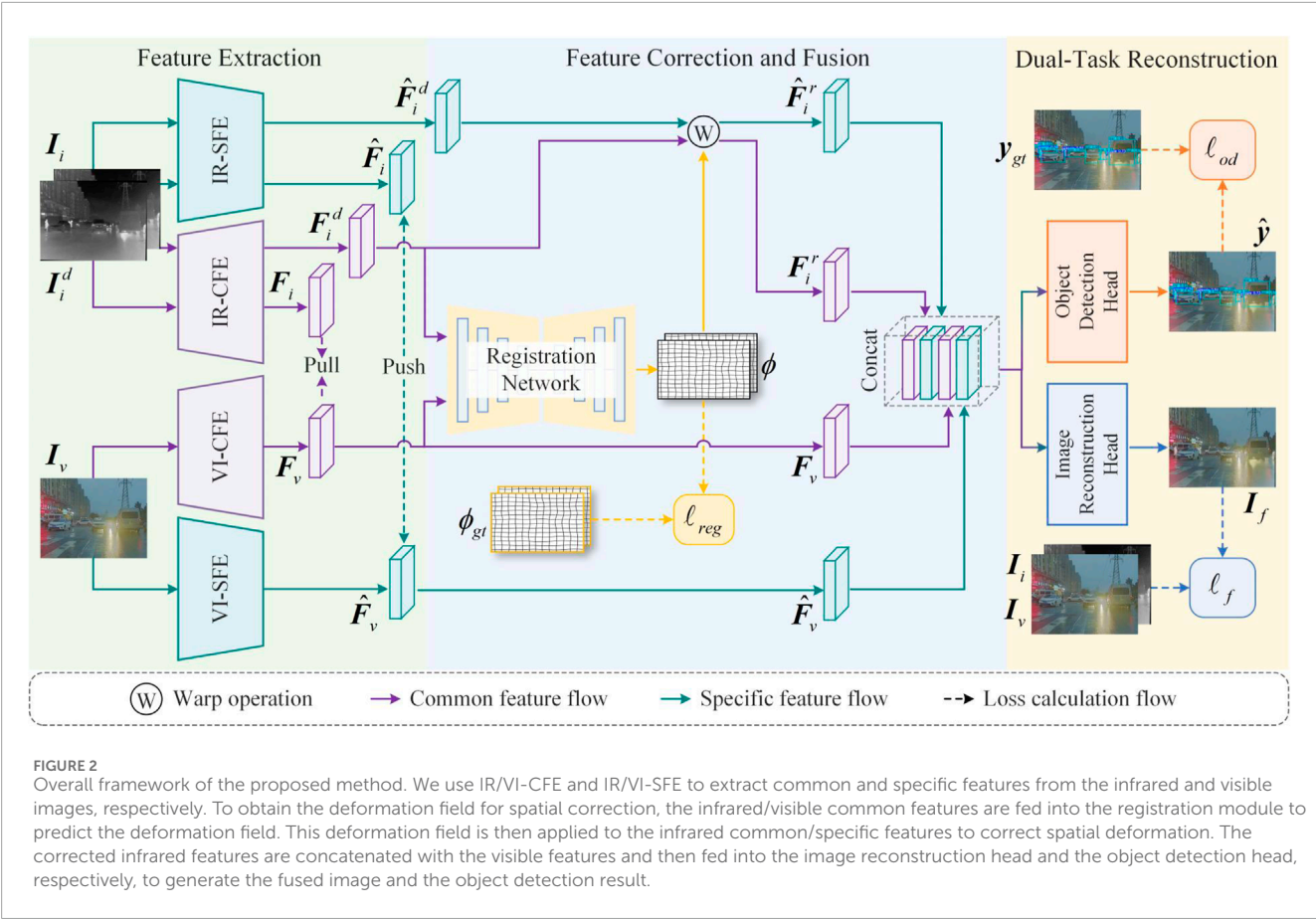
In summary, our main contributions are as follows:

- (1) We are the first to unify unregistered image fusion and object detection within a single framework, breaking the limitations of object detection in real-world applications.
- (2) We propose a modality consistency loss that effectively eliminates the domain discrepancy between infrared and visible images, improving image registration accuracy.
- (3) Our method demonstrates excellent performance in image alignment, fusion, and object detection across multiple datasets. And our method can be used with other radiation frequencies where different modalities require image fusion like, for example, radio-frequency, x- and gamma rays used in medical imaging.

The rest of this paper is organized as follows. Section 2 briefly reviews related work on high-level vision task-driven image fusion and unregistered infrared-visible image fusion. Section 3 describes the proposed method in detail. Section 4 presents and discusses the experimental results. Section 5 concludes the paper.

2 Related work

In this section, we first provide a brief overview of high-level vision task-driven infrared and visible image fusion methods, and then review existing approaches for unregistered infrared and visible image fusion.



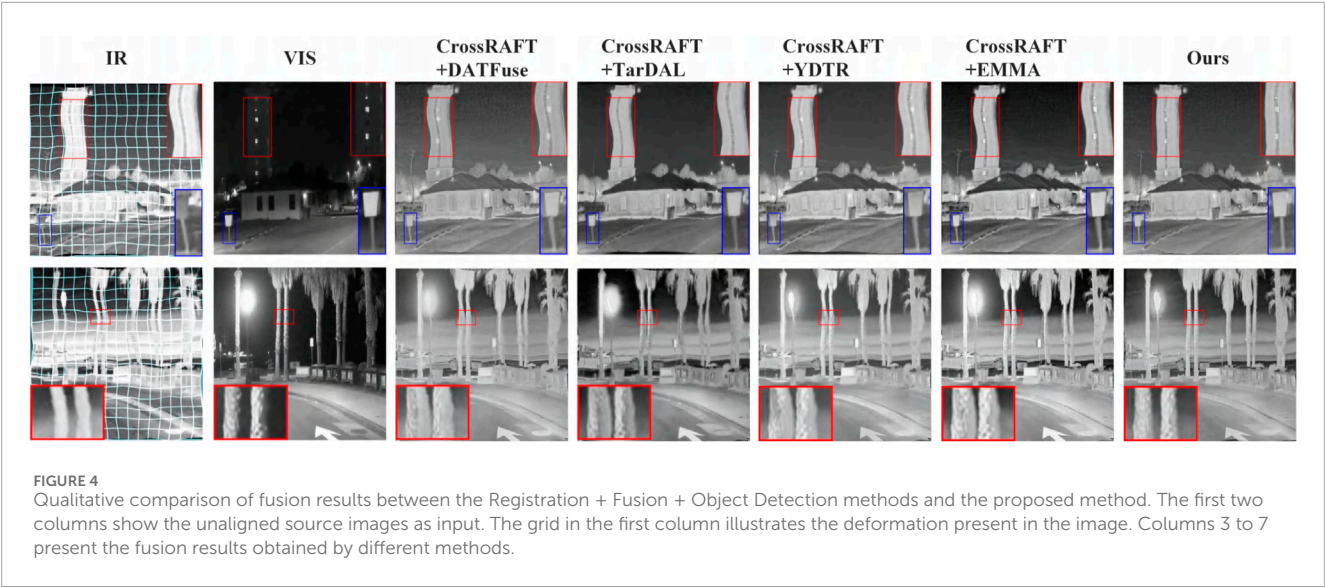


TABLE 1 Quantitative comparison of fusion results between the Registration + Fusion + Object Detection methods and the proposed method.

Methods	$Q_{CC}\uparrow$	$Q_{AB/F}\uparrow$	$Q_{CV}\downarrow$	$Q_{SSIM}\uparrow$
DATFuse	0.8303	0.3246	1425.2631	1.2189
TarDAL	0.8317	0.3313	1396.1484	1.2205
YDTR	0.8246	0.3179	1383.2556	1.2133
EMMA	0.8255	0.3341	1399.4075	1.2236
Ours	0.8325	0.3420	1375.5238	1.2271

Bolded values indicate the best performance.

2.1 High-level vision task-driven infrared and visible image fusion

High-level vision task-driven fusion methods typically incorporate a semantic segmentation [24–27] or object detection network [23, 28] after the fusion network, using the loss functions from these downstream tasks to constrain the fusion results and improve the quality of the fused image. However, introducing high-level vision tasks at the fused image level only provides indirect guidance for the feature extraction network to learn features relevant to the downstream tasks.

To provide direct task-level guidance at the feature level and further enhance fusion performance, PSFusion [29] injects semantic features extracted from a segmentation task directly into the fusion network. SegMiF [25] feeds the fused result into a semantic segmentation network to extract semantic features, which are then interacted with the multimodal image features from the encoder to enhance the fusion result. MRFS [26] interacts and fuses the source image features before feeding them into a semantic segmentation head to enforce semantic supervision, thereby improving the global scene perception of the fusion network. MetaFusion [28] sends the fused result into an object detection network to extract features,

which are then combined with the source image features and passed into a meta-feature generator to guide feature extraction in the fusion branch.

Although these methods improve fusion performance to some extent by leveraging downstream high-level tasks, they all assume that the input images are perfectly aligned in spatial position—a condition rarely met in real-world applications. In practice, such methods rely on additional image registration algorithms to achieve accurate alignment before performing fusion. This not only makes the fusion quality highly dependent on the registration accuracy but also significantly increases the complexity of the overall network design.

2.2 Unregistered infrared and visible image fusion

To address the problem of unregistered infrared and visible image fusion, most existing approaches combine registration and fusion algorithms, i.e., first aligning the input misaligned image pairs and then performing fusion. However, due to the large modality gap between infrared and visible images, ignoring the adverse

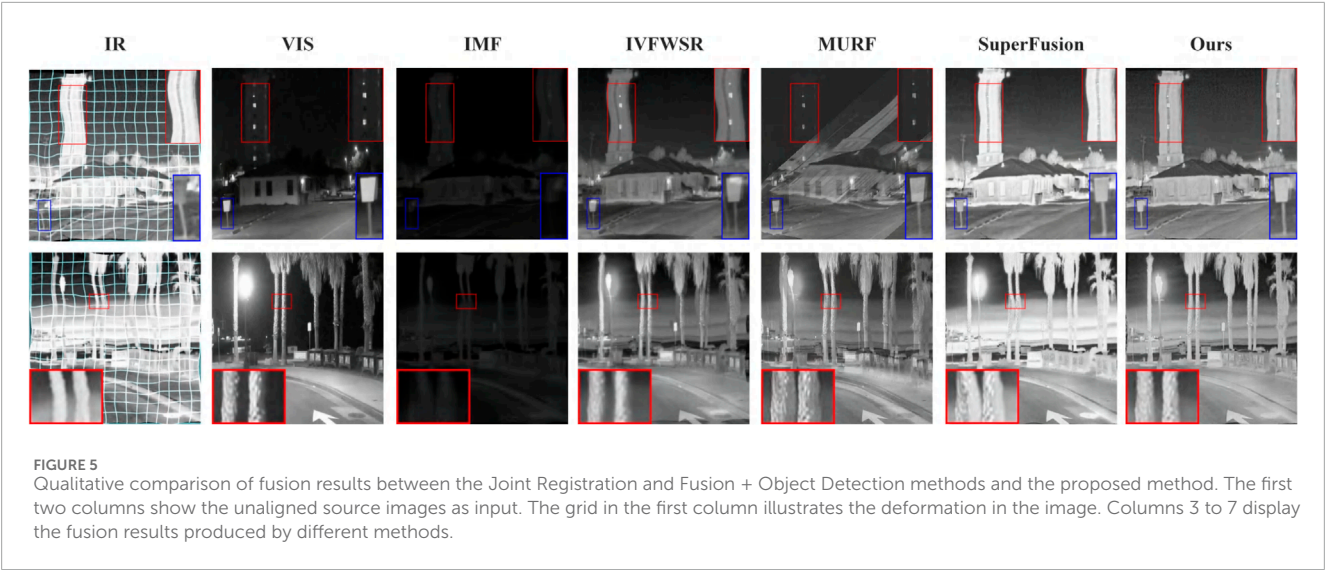


TABLE 2 Quantitative comparison of fusion results between the Joint Registration and Fusion + Object Detection methods and the proposed method.

Methods	$Q_{CC}\uparrow$	$Q_{AB/F}\uparrow$	$Q_{CV}\downarrow$	$Q_{SSIM}\uparrow$
IMF	0.8221	0.3119	1477.6932	1.2058
IVFWSR	0.8269	0.3208	1586.8251	1.2115
MURF	0.8315	0.3254	1456.3259	1.2140
SuperFusion	0.8320	0.3396	1399.4521	1.2207
Ours	0.8325	0.3420	1375.5238	1.2271

Bolded values indicate the best performance.

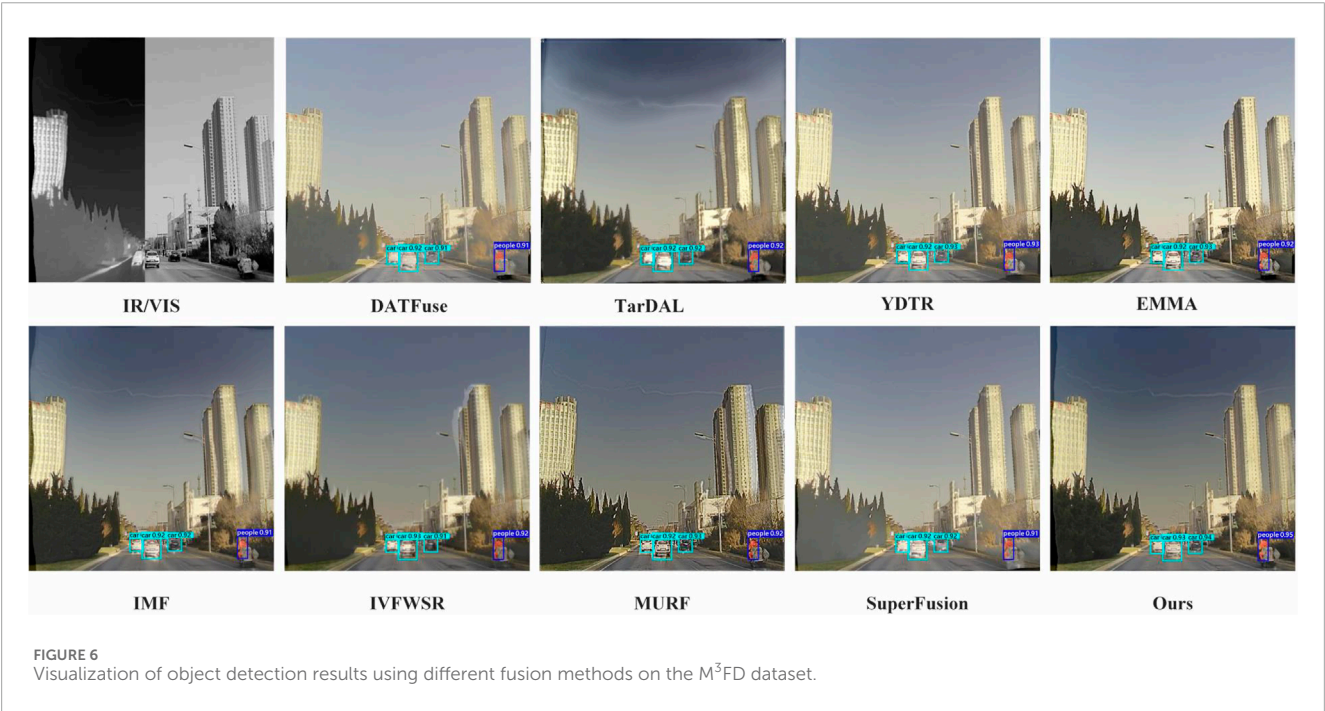


TABLE 3 Quantitative object detection results of different fusion methods on the M³FD dataset.

Methods	mAP _{50→90} ↑
DATFuse	53.10
TarDAL	53.20
YDTR	53.80
EMMA	54.20
IMF	52.40
IVFWSR	52.60
MURF	52.20
SuperFusion	53.80
Ours	54.50

Bolded values indicate the best performance.

impact of modality discrepancy on registration can greatly degrade fusion quality. For instance, ReCoNet [30] adopts this strategy but produces suboptimal fusion results due to this issue. UMF-CMGR [31] and IMF [32] consider the effect of modality differences on registration results. They propose to convert visible images into pseudo-infrared images via an image generation network and then perform mono-modal registration between the pseudo-infrared and misaligned infrared images. However, the quality of the generated image has a direct impact on the final performance of these methods. Moreover, these methods treat registration and fusion as two independent tasks, failing to establish a unified framework where both tasks can benefit each other.

To address this, RFNet [33] and MURF [34] treat image fusion as a downstream task of registration and improve registration performance by enhancing the sparsity of the gradient in the fused result. However, to tackle the modality discrepancy issue during registration, both methods aim to transform the multimodal registration into a mono-modal one. Specifically, RFNet uses an image generation model to produce a pseudo-image with the same modality as the misaligned one before performing mono-modal registration, while MURF leverages contrastive learning to extract modality-invariant features from the input image pair for registration. Similarly, Super-Fusion [35] extracts modality-invariant features using shared-parameter encoders and consistency constraints on the fused result for registration.

Nevertheless, the information carried by modality-invariant features in infrared-visible pairs is often far less rich than the complementary information present in the image pair. As a result, it is difficult to achieve satisfactory cross-modal registration using only modality-invariant features. In addition, the above methods all follow a two-stage approach (registration + fusion). This two-stage strategy greatly limits deployment in practical applications due to computational constraints. Although RFVIF [36], IVFWSR [37] and MulFS-CAP [38] attempt to achieve registration and fusion within a single-stage framework, the types of deformations they can handle remain limited. Unlike the methods mentioned above,

our approach considers multiple challenges simultaneously: the impact of modality discrepancy on cross-modal registration, the deployment limitations of two-stage processing, and the feature requirements of downstream high-level vision tasks for both registration and fusion.

3 Methods

3.1 Overview

As shown in Figure 2, the proposed method consists of three core components: feature extraction, feature alignment and fusion, and dual-task reconstruction. The feature extraction component is designed to obtain both modality-specific and modality-common features from the source images. The feature alignment and fusion component is used to predict a deformation field, which is then used to spatially align the infrared-specific and common features. These aligned features are then fused with the corresponding visible image's specific and common features. In the dual-task reconstruction stage, the fused features are fed into the object detection head and the image reconstruction head, respectively, to generate both the object detection result map and the fused image.

3.2 Feature extraction

The main objective of feature extraction is to extract both the common and specific features of infrared and visible images, in order to facilitate subsequent cross-modal registration and feature fusion. This process consists of four modules: the IR-Specific Feature Extraction (IR-SFE) module, the VI-Specific Feature Extraction (VI-SFE) module, the IR-Common Feature Extraction (IR-CFE) module, and the VI-Common Feature Extraction (VI-CFE) module. Among them, the IR/VI-SFE modules are used to extract modality-specific features from the infrared/visible images, while the IR/VI-CFE modules are used to extract their common features. Assume that each sample in the training dataset contains three images: a pixel-wise strictly aligned infrared image I_i , a visible image I_v , and a deformed infrared image I_i^d . We feed I_i and I_i^d into the IR-CFE and IR-SFE, respectively, to obtain the infrared common feature F_i , the deformed infrared common feature F_i^d , the infrared specific feature \hat{F}_i , and the deformed infrared specific feature \hat{F}_i^d . At the same time, we feed I_v into the VI-CFE and VI-SFE to obtain the visible common feature F_v and the visible specific feature \hat{F}_v .

In the cross-modal registration process, it is usually necessary to rely on the common information between cross-modal images to establish pixel-wise correspondences. To reduce the modality gap between infrared and visible images and thus establish more accurate pixel-wise correspondences, we introduce a modality consistency loss ℓ_c :

$$\ell_c = \frac{1}{HWC} \|F_i - F_v\|_1, \quad (1)$$

Here, H , W , and C denote the height, width, and number of channels of the feature maps, respectively, and $\|\cdot\|_1$ represents the l1-norm. In addition, considering that the goal of image fusion is to integrate as much complementary information as possible

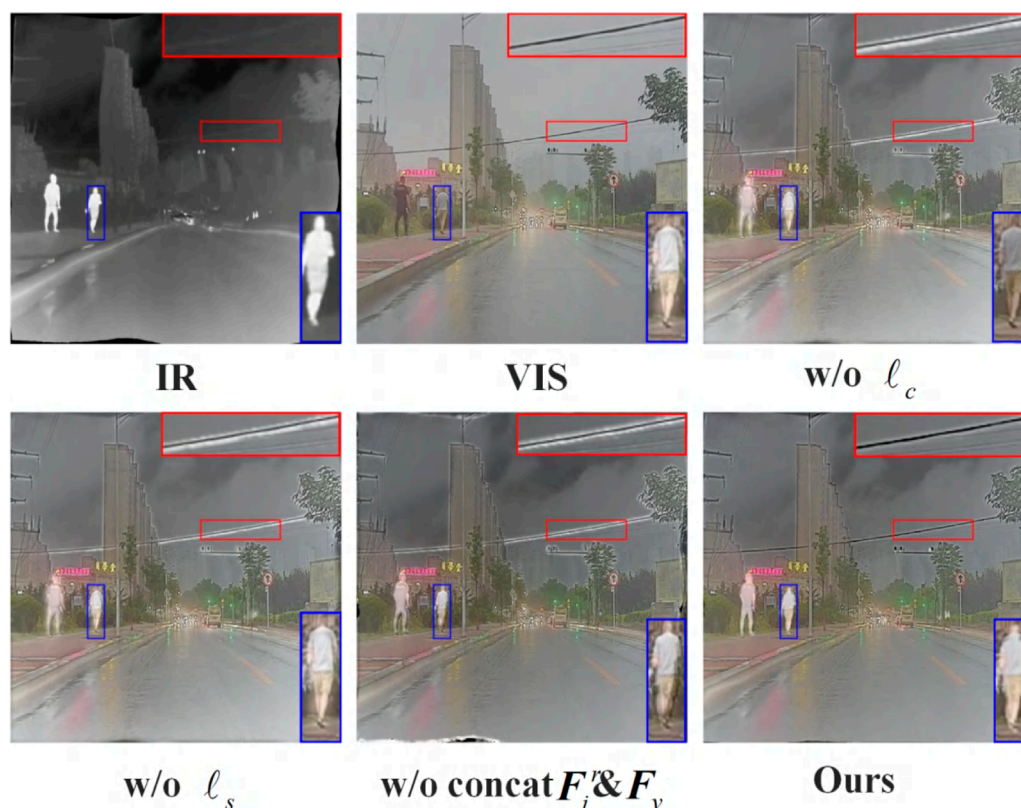


FIGURE 7
Ablation study of the core designs.

TABLE 4 Quantitative results of the ablation study on the core designs.

Methods	$Q_{CC} \uparrow$	$Q_{AB/F} \uparrow$	$Q_{CV} \downarrow$	$Q_{SSIM} \uparrow$
w/o ℓ_c	0.8304	0.3469	1339.9062	1.2204
w/o ℓ_s	0.8313	0.3439	1336.8814	1.2256
w/o Concat F_i' and F_v	0.8274	0.3451	1369.4537	1.2114
Ours	0.8325	0.3420	1375.5238	1.2271

Bolded values indicate the best performance.

from cross-modal source images into a single image, we introduce the modality complementary information loss ℓ_s to further enrich the complementary information from the source images in the fused image:

$$\ell_s = -\frac{1}{HWC} \|\hat{F}_i - \hat{F}_v\|_1. \quad (2)$$

3.3 Feature alignment and fusion

Feature alignment corrects the deformation in infrared features by predicting a deformation field, thereby achieving spatial alignment between infrared and visible features. This process is mainly implemented by the registration network. Subsequently, the

aligned infrared features are fused with the visible features to obtain the fused features. As shown in Figure 3, the registration network is composed of a Channel and Spatial Enhancement Block (CSEB) and a Multi-Scale Registration Block (MSRB). The CSEB is mainly used to enhance the information beneficial to registration at both the channel and spatial levels, thereby improving the accuracy of the predicted deformation field. The CSEB consists of six feature extraction layers and a Global Average Pooling (GAP) layer. Each feature extraction layer is composed of a convolutional layer with a kernel size of 3×3 , stride 1, followed by Batch Normalization (BatchNorm) and a LeakyReLU activation function. The MSRB is used to predict the deformation field to correct the deformed infrared features and ensure spatial alignment between the infrared and visible features. The MSRB adopts a U-Net-like architecture.

TABLE 5 Quantitative analysis results of the hyperparameter study.

γ	λ_1	λ_2	λ_3	$Q_{CC}\uparrow$	$Q_{AB/F}\uparrow$	$Q_{CV}\downarrow$	$Q_{SSIM}\uparrow$
2	10	5	1	0.8235	0.3352	1450.3498	1.2222
2	10	5	10	0.8198	0.3389	1683.8772	1.2195
2	10	1	5	0.8123	0.3321	1502.6641	1.2088
2	10	10	5	0.8260	0.3334	1465.2293	1.2247
2	1	5	5	0.8011	0.3195	1450.3288	1.1954
2	20	5	5	0.8059	0.3248	1529.1245	1.1996
1	10	5	5	0.8144	0.3340	1499.3888	1.2111
5	10	5	5	0.8080	0.3302	1775.1124	1.2020
2	10	5	5	0.8325	0.3420	1375.5238	1.2271

Bolded values indicate the best performance.

TABLE 6 Computational efficiency comparison of four SOTA Joint Registration and Fusion methods, the value is tested on GPU.

Methods	FLOPs(G)	Size(M)	Time(s)
IMF	1724.08	13.30	0.82
IVFWSR	859.43	14.09	0.33
MURF	120.72	1.76	1.18
SuperFusion	65.43	0.14	0.27
Ours	60.12	0.97	0.40

Bolded values indicate the best performance.

We input the deformed infrared common feature F_i^d and the visible common feature F_v into two CSEBs with unshared parameters, obtaining the enhanced features \tilde{F}_i^d and \tilde{F}_v , respectively. Taking the enhancement process of F_v as an example, F_v is fed into three feature extraction layers to generate the spatial enhancement weights W_v^s . To enhance registration-relevant information at the spatial level, we perform element-wise multiplication between W_v^s and F_v :

$$F_v^s = W_v^s \odot F_v, \quad (3)$$

Here, F_v^s denotes the feature enhanced at the spatial level, and \odot represents the element-wise multiplication operation. We feed F_v^s into three feature extraction layers and a global average pooling (GAP) layer to obtain feature W_v^c for channel-level enhancement. Then, W_v^c is element-wise multiplied with F_v^s to produce the enhanced feature \tilde{F}_v , which has been refined at both the spatial and channel levels:

$$\tilde{F}_v = W_v^c \odot F_v^s, \quad (4)$$

Similarly, we obtain the deformed infrared common feature \tilde{F}_i^d enhanced at both the spatial and channel levels. We concatenate \tilde{F}_i^d and \tilde{F}_v along the channel dimension and feed the resulting

feature into the MSRB to predict the deformation field ϕ . To ensure the accuracy of the predicted deformation field, we introduce a registration loss ℓ_{reg} :

$$\ell_{reg} = \frac{1}{2HW} \|\phi - \phi_{gt}\|_1, \quad (5)$$

Here, ϕ_{gt} is the label of ϕ .

We use ϕ to correct F_i^d and \tilde{F}_i^d respectively, resulting in the corrected infrared common feature F_i^r and infrared-specific feature \tilde{F}_i^r :

$$\begin{aligned} F_i^r &= \phi \circ F_i^d, \\ \tilde{F}_i^r &= \phi \circ \tilde{F}_i^d, \end{aligned} \quad (6)$$

Here, \circ denotes the Warp operation, which resamples the deformed feature maps based on ϕ to correct the deformations within them. During the fusion process, to minimize information loss, we concatenate F_i^r , \tilde{F}_i^r , F_v , and \tilde{F}_v along the channel dimension to obtain the fused feature F_f :

$$F_f = [F_i^r, \tilde{F}_i^r, F_v, \tilde{F}_v], \quad (7)$$

Here, $[\cdot]$ represents the operation of concatenation along the channel dimension.

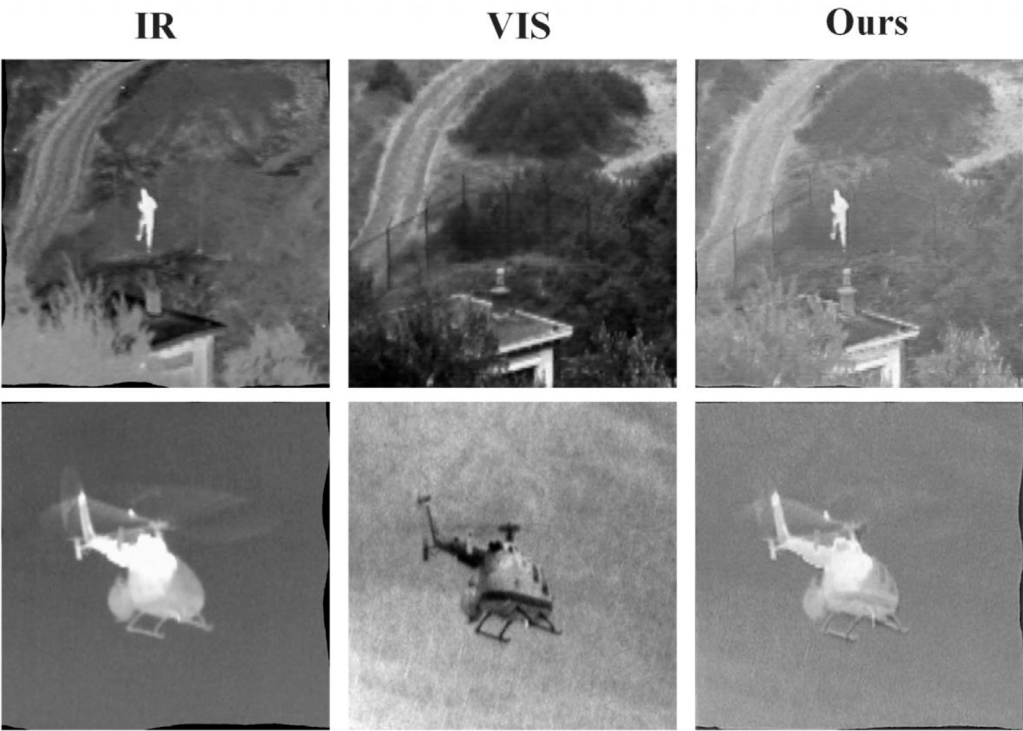


FIGURE 8
Fusion results of our method on different scenarios.



FIGURE 9
Failure cases of our method on the real-world dataset CVC-14.

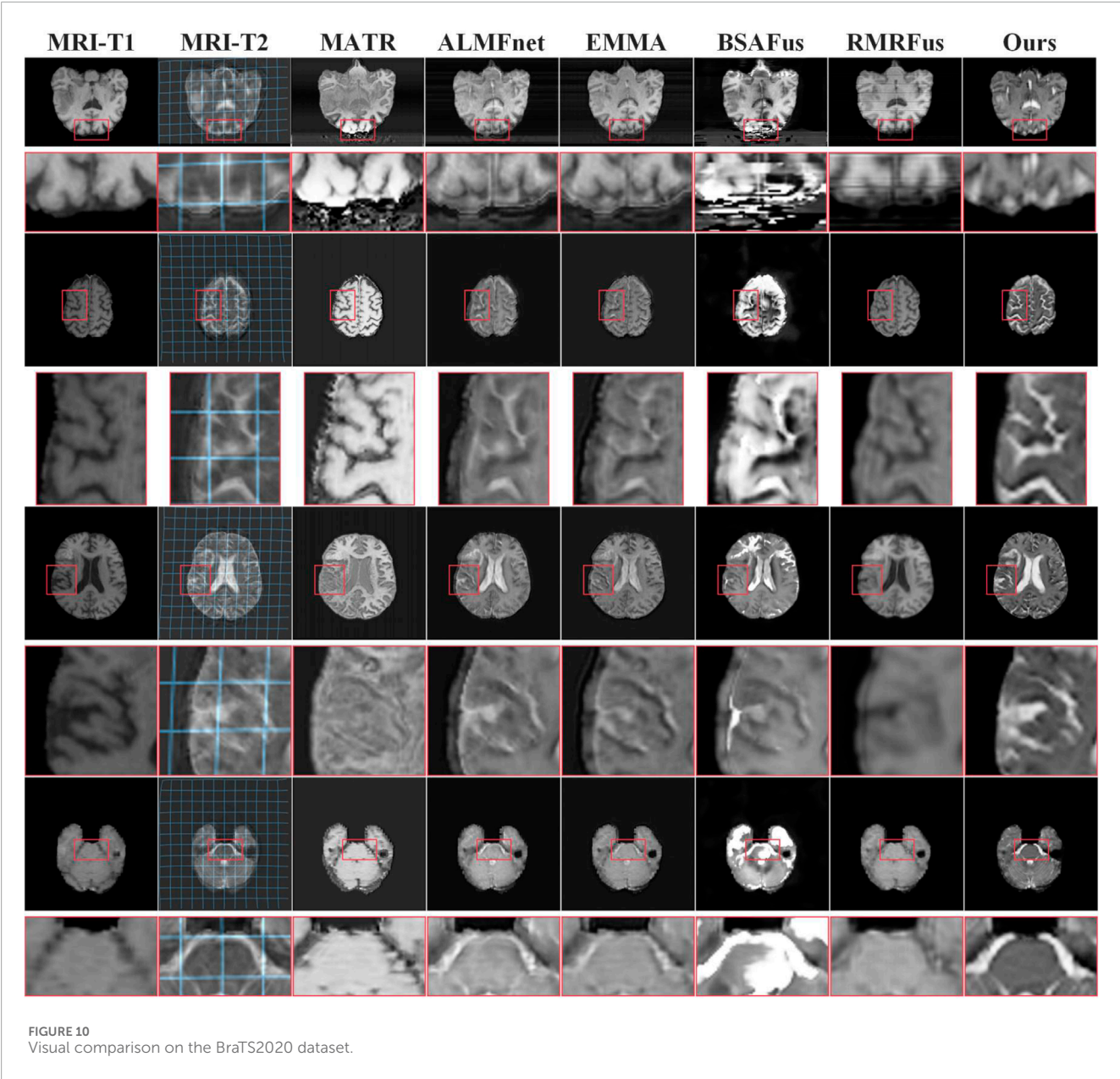


TABLE 7 Quantitative analysis results on the BraTS2020 dataset.

Methods	$Q_{CC}\uparrow$	$Q_{AB/F}\uparrow$	$Q_{CV}\downarrow$	$Q_{SSIM}\uparrow$
MATR	0.7889	0.2901	1345.4510	1.2299
ALMFnet	0.7749	0.2888	1606.5911	1.2155
EMMA	0.7906	0.2853	1568.7139	1.2220
BSAFus	0.7812	0.3001	1436.1287	1.2318
RMRFus	0.7784	0.2992	1409.9831	1.2007
Ours	0.7934	0.3063	1399.5234	1.2454

Bolded values indicate the best performance.

3.4 Dual-task reconstruction

In the dual-task reconstruction, the fused feature is fed into both the object detection head and the image reconstruction head to respectively generate the object detection result map and the fused image. The dual-task reconstruction primarily consists of the object detection head and the image reconstruction head. We adopt YOLOv5 [39] as the object detection head. The image reconstruction head is composed of three feature extraction layers, where the LeakyReLU activation function in the final layer is replaced with a Tanh activation function. The fused feature F_f is input into both the object detection head and the image reconstruction head to obtain the object detection result map \hat{y} and the fused image I_f , respectively. To ensure high-quality object detection results, we introduce the object detection loss ℓ_{ob} to constrain the network:

$$\ell_{ob} = c_{yolov5}(y, y_{gt}), \quad (8)$$

Here, $c_{yolov5}(\cdot)$ refers to the loss function used during the training of YOLOv5. In addition, to encourage the fused image to retain as much shared and complementary information from both infrared and visible images as possible, we introduce luminance loss ℓ_b and gradient loss ℓ_g , and construct the fusion loss ℓ_f accordingly:

$$\ell_f = \ell_b + \gamma \ell_g, \quad (9)$$

Here, γ denotes the balancing hyperparameter. The gradient loss ℓ_g is defined as:

$$\ell_g = \frac{1}{HW} \|\nabla I_f - \max(\nabla I_i, \nabla I_v)\|_1, \quad (10)$$

Here, ∇ denotes the Sobel operator. The luminance loss ℓ_b is defined as:

$$\ell_b = \frac{1}{HW} \|I_f - \max(I_i, I_v)\|_1. \quad (11)$$

Finally, we define the total loss ℓ_t as follows:

$$\ell_t = \ell_c + \ell_s + \lambda_1 \ell_{reg} + \lambda_2 \ell_f + \lambda_3 \ell_{ob}, \quad (12)$$

Here, $\lambda_n (n = 1, 2, 3)$ denotes the balancing hyperparameter.

4 Experiments

4.1 Experimental setup

4.1.1 Datasets and implementation details

4.1.1.1 Datasets

Following standard experimental practices in the image fusion field [40–43], we trained our model on 152 pairs of infrared and visible images with a resolution of 512×512 from the RoadScene Xu et al. [44, 45] dataset. For testing, we used 18 pairs of images from RoadScene and 17 pairs from M³FD [21]. The misaligned infrared images were generated by randomly applying a combination of rigid and non-rigid deformations to the originally well-aligned infrared images. This type of mixed deformation is applied randomly to the original aligned images in each epoch to augment the training data.

4.1.1.2 Implementation details

The proposed method was implemented using the PyTorch framework and trained on a single NVIDIA GeForce RTX 3090 GPU. The model was trained for 150 epochs with a batch size of 8, a learning rate of $1e-3$, and the Adam optimizer was used to update the model parameters. The four hyperparameters in the loss function were set to $\gamma = 2$, $\lambda_1 = 10$, $\lambda_2 = 5$, and $\lambda_3 = 5$.

4.1.2 Evaluation metrics

We selected four commonly used image quality evaluation metrics to objectively assess the quality of the fusion results, including correlation coefficient (Q_{CC}) [46], gradient-based fusion performance ($Q_{AB/F}$) [47], Chen-Varshney metric (Q_{CV}) [48], and structural similarity (Q_{SSIM}) [49]. Metric Q_{CC} evaluates the linear correlation between the fused image and the source images, reflecting their similarity. Metric $Q_{AB/F}$ assesses the amount of edge information transferred from the source images to the fused image. Metric Q_{CV} takes into account both edge information and human visual perception. Metric Q_{SSIM} quantifies information loss and distortion in the fused image by comparing it with the source images. Among these metrics, a lower value of indicates better fusion quality, while higher values of the other metrics indicate better performance. In addition, we adopted metric $mAP_{50 \rightarrow 90}$ [50] as the evaluation metric for the object detection task, where a higher $mAP_{50 \rightarrow 90}$ value indicates better detection performance.

4.2 Comparison with state-of-the-art methods

In our experiments, we first compare the proposed method with two categories of fusion approaches for unaligned infrared and visible images based on their fusion results. We then compare the subsequent object detection results obtained using these two categories of methods. The first category involves registering the images to be fused, followed by image fusion and then object detection. We refer to this category as Registration + Fusion + Object Detection. The second category performs joint training of registration and fusion to directly handle unaligned images, followed by object detection. We refer to this as Joint Registration and Fusion + Object Detection.

4.2.1 Comparison with registration + fusion + object detection methods

For the Registration + Fusion + Object Detection methods, we follow the standard processing pipeline used in prior work. We first adopt the high-performing registration method CrossRAFT [51] to align the images to be fused. Then, we apply four advanced infrared and visible image fusion methods to the aligned results, including DATFuse [52], TarDAL [21], YDTR [53], and EMMA [54]. Figure 4 shows the visual results of different methods. As seen from the fusion results, our proposed method not only demonstrates stronger capability in preserving structures and textures but also effectively avoids distortions and artifacts caused by feature misalignment. In addition, we performed objective evaluations of the results from different methods. As shown in Table 1, our method achieves the best performance across all four evaluation metrics.

4.2.2 Comparison with joint registration and fusion + object detection methods

In recent years, joint registration and fusion methods have attracted significant attention. To demonstrate the superiority of our approach over these methods, we compared its performance with four joint registration and fusion methods: IMF, IVFWSR, MURE, and SuperFusion. Figure 5 presents a qualitative comparison of the fusion results produced by different methods. It can be observed that our method exhibits clear advantages in terms of feature alignment, contrast preservation, and detail retention. In addition, we conducted quantitative experiments to visually compare the performance differences. As shown in Table 2, our method achieves the best performance across all four evaluation metrics.

4.2.3 Performance evaluation on infrared and visible image object detection

We evaluated the object detection performance of the two aforementioned categories of methods, as well as the proposed method, on the M³FD dataset. Figure 6 shows the visualized results of object detection. In comparison, our proposed method achieves superior performance. Table 3 presents the quantitative results. The fused outputs generated by our method help the detection network achieve the highest object detection accuracy. This further demonstrates the superior fusion capability of our approach for object detection tasks.

4.3 Ablation study

The core of the proposed method lies in the losses designed to eliminate modality differences, namely, losses ℓ_c and ℓ_s . In this section, we conduct ablation studies on these key components to verify their effectiveness. All experiments are conducted on the M³FD dataset. From the ablation results, it can be observed that removing losses ℓ_c and ℓ_s leads to a decline in the model's ability to correct local deformations, as shown in Figure 7. In addition, when the shared information is excluded during fusion and only complementary information is used for concatenation, the visual quality of the fused image does not deteriorate significantly, but the objective evaluation results in Table 4 show a noticeable drop in performance.

4.4 Analysis of hyperparameters

In our proposed method, four main hyperparameters are defined: $\lambda_1, \lambda_2, \lambda_3$, which balances different losses, i.e., ℓ_{reg} , ℓ_f , and ℓ_{ob} , and γ , which balances luminance loss ℓ_b and gradient loss ℓ_g . During model training, $\lambda_1, \lambda_2, \lambda_3, \gamma$ are set to 10, 5, 5, two respectively.

Next, we analyze the impact of variations in these hyperparameters on model performance. To analyze the impact of $\lambda_1, \lambda_2, \lambda_3$ on fusion performance, we perform a search over $\lambda_1, \lambda_2, \lambda_3$ values in the ranges of 1–20, 1 to 10, and 1 to 10. The quantitative evaluation results for both fusion and downstream object detection are presented in Table 5. As shown in Table 5, the model achieves optimal performance on fusion when $\lambda_1 = 10, \lambda_2 = 5$, and $\lambda_3 = 5$.

To verify the effectiveness of the hyperparameter γ , we fix λ_1, λ_2 and λ_3 to 10, 5, 5 and analyze the model performance as γ varies

from 1 to 5. As shown in Table 5, the model achieves the best fusion performance when γ is set to 2. Therefore, we set the hyperparameter γ to 2.

4.5 Analysis of computational complexity

As shown in Table 6, a complexity evaluation is introduced to evaluate the efficiency of our method from three aspects, i.e., FLOPs, training parameters and runtime. Wherein, for FLOPs calculation, the size of the input images is standardized to 512×512 pixels. The inference time is calculated as the average time taken to process 18 scene images from RoadScene's test dataset. From Table 6, our model performs the best in FLOPs, implying that our method has fast calculation speed and is application-friendly. The average inference time for our model to fuse two source images is 0.40 s, only a bit longer than the SOTA method, demonstrating that our model's inference speed is relatively fast and acceptable. Besides, the parameter size of our model is only 0.97M, which can be easily deployed in practical applications. This indicates the efficiency of our method, which can serve practical vision tasks well with better visual performance.

4.6 Analysis of generalization ability

To validate the generalization ability of our method, we conduct experiments under other scenarios. Fusion results are shown in Figure 8. From the qualitative results we can see that our proposed model performs perfectly under other scenarios.

4.7 Analysis of limitation

The proposed method enables mutual enhancement between infrared-visible image fusion and object detection, specifically designed to handle misaligned source images, achieving better experimental results compared to other methods. However, our approach still has certain limitations. Specifically, since our model is trained on the generated unaligned dataset, where the deformations in real-world images cannot be fully included, failure cases appear under real-world scenarios. As shown in Figure 9, our method fails to handle deformations under real-world scenarios. Improving the robustness of our method is vital for future research.

4.8 Further discussion

To validate the effectiveness of the proposed method in the field of medical imaging, we conduct a comparative study on the publicly available BraTS2020 Menze et al. [55] dataset. Specifically, we first employ the state-of-the-art medical image registration method CorrMLP Meng et al. [56] to align the deformed MRI-T2 images to the reference MRI-T1 images, and subsequently apply several advanced fusion methods (including MATR Tang et al. [57], ALMFnet Mu et al. [58], EMMA Zhao et al. [54], BSAFus Li et al. [47], and RMRfus Zhang et al. [59]) for image fusion. As shown in Figure 10, the fusion images generated by the

proposed method exhibit superior image quality and effectively correct artifacts and spatial deformations. In contrast, existing "registration + fusion" methods often introduce noticeable artifacts when handling unregistered medical images, significantly degrading the visual quality of the fused images. Furthermore, as reported in Table 7, the quantitative analysis results further demonstrate the significant advantages of the proposed method in terms of fusion performance.

5 Conclusion

This paper proposes a mutual promotion algorithm for infrared and visible image fusion and object detection, tailored for unaligned image scenarios. Considering the significant modality differences between infrared and visible images, we design specific loss functions to reduce such differences, thereby easing the difficulty of cross-modality image registration and improving its accuracy. In addition, we adopt a mutually beneficial learning strategy that enables the fusion task and the downstream object detection task to enhance each other, leading to improved quality in both the fused images and detection results. Extensive qualitative and quantitative experiments demonstrate the superiority of our method over existing state-of-the-art approaches. In addition, our approach can be used with other radiation frequencies where different modalities require image fusion like, for example, radio-frequency, x- and gamma rays used in medical imaging.

Data availability statement

The original contributions presented in the study are included in the article/supplementary material, further inquiries can be directed to the corresponding author.

Author contributions

DH: Conceptualization, Methodology, Writing – review and editing, Writing – original draft, Investigation. KW: Writing – review and editing, Project administration, Data curation. CZ: Validation,

Writing – review and editing, Formal Analysis. ZL: Methodology, Supervision, Writing – original draft, Funding acquisition, Writing – review and editing. YC: Formal Analysis, Visualization, Project administration, Writing – review and editing. SD: Resources, Data curation, Validation, Writing – review and editing. CK: Resources, Writing – review and editing, Formal Analysis.

Funding

The author(s) declare that financial support was received for the research and/or publication of this article.

Conflict of interest

Authors DH, KW, CZ, ZL, YC, SD, and CK were employed by Yunnan Power Grid Co., Ltd.

The authors declare that this study received funding from the Science and Technology Project of China Southern Power Grid Co., Ltd. (No. YNKJXM20240052). The funder had the following involvement in the study: study design, collection, analysis, interpretation of data, the writing of this article, and the decision to submit it for publication.

Generative AI statement

The author(s) declare that Generative AI was used in the creation of this manuscript. AI was only used to polish the paper.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

References

1. Cao Y, Guan D, Huang W, Yang J, Cao Y, Qiao Y. Pedestrian detection with unsupervised multispectral feature learning using deep neural networks. *information fusion* (2019) 46:206–17. doi:10.1016/j.inffus.2018.06.005
2. Li C, Zhu C, Huang Y, Tang J, Wang L. Cross-modal ranking with soft consistency and noisy labels for robust rgb-t tracking. In: *Proceedings of the European conference on computer vision (ECCV)* (2018). p. 808–23.
3. Lin X, Li J, Ma Z, Li H, Li S, Xu K, et al. Learning modal-invariant and temporal-memory for video-based visible-infrared person re-identification. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)* (2022). p. 20973–82.
4. Ha Q, Watanabe K, Karasawa T, Ushiku Y, Harada T. Mfnet: towards real-time semantic segmentation for autonomous vehicles with multi-spectral scenes. In: *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE (2017). p. 5108–15.
5. Li S, Kang X, Fang L, Hu J, Yin H. Pixel-level image fusion: a survey of the state of the art. *information Fusion* (2017) 33:100–12. doi:10.1016/j.inffus.2016.05.004
6. Li H, Qi X, Xie W. Fast infrared and visible image fusion with structural decomposition. *Knowledge-Based Syst* (2020) 204:106182. doi:10.1016/j.knosys.2020.106182
7. Li H, Qiu H, Yu Z, Zhang Y. Infrared and visible image fusion scheme based on nsct and low-level visual features. *Infrared Phys and Technology* (2016) 76:174–84. doi:10.1016/j.infrared.2016.02.005
8. Zhang Q, Liu Y, Blum RS, Han J, Tao D. Sparse representation based multi-sensor image fusion for multi-focus and multi-modality images: a review. *Inf fusion* (2018) 40:57–75. doi:10.1016/j.inffus.2017.05.006
9. Xie M, Wang J, Zhang Y. A unified framework for damaged image fusion and completion based on low-rank and sparse decomposition. *Signal Processing: Image Commun* (2021) 98:116400. doi:10.1016/j.image.2021.116400
10. Li H, Wang Y, Yang Z, Wang R, Li X, Tao D. Discriminative dictionary learning-based multiple component decomposition for detail-preserving noisy image fusion. *IEEE Trans Instrumentation Meas* (2020) 69:1082–102. doi:10.1109/tim.2019.2912239

11. Xiao W, Zhang Y, Wang H, Li F, Jin H. Heterogeneous knowledge distillation for simultaneous infrared-visible image fusion and super-resolution. *IEEE Trans Instrumentation Meas* (2022) 71:1–15. doi:10.1109/tim.2022.3149101
12. Zhang Y, Yang M, Li N, Yu Z. Analysis-synthesis dictionary pair learning and patch saliency measure for image fusion. *Signal Process.* (2020) 167:107327. doi:10.1016/j.sigpro.2019.107327
13. Li H, Wu X-J, Kittler J. Rfn-nest: an end-to-end residual fusion network for infrared and visible images. *Inf Fusion* (2021) 73:72–86. doi:10.1016/j.inffus.2021.02.023
14. Li H, Wu X. Densefuse: a fusion approach to infrared and visible images. *IEEE Trans Image Process* (2019) 28:2614–23. doi:10.1109/tip.2018.2887342
15. Shi Y, Liu Y, Cheng J, Wang ZJ, Chen X. Vdmufusion: a versatile diffusion model-based unsupervised framework for image fusion. *IEEE Trans Image Process* (2025) 34:441–54. doi:10.1109/tip.2024.3512365
16. Ma J, Tang L, Xu M, Zhang H, Xiao G. Stdffusionnet: an infrared and visible image fusion network based on salient target detection. *IEEE Trans Instrumentation Meas* (2021) 70:1–13. doi:10.1109/TIM.2021.3075747
17. Du K, Li H, Zhang Y, Yu Z. Chitnet: a complementary to harmonious information transfer network for infrared and visible image fusion. *IEEE Trans Instrumentation Meas* (2025) 74:1–17. doi:10.1109/TIM.2025.3527523
18. Yang Z, Li Y, Tang X, Xie M. Mgfusion: a multimodal large language model-guided information perception for infrared and visible image fusion. *Front Neurobotics* (2024) 18:1521603. doi:10.3389/fnbot.2024.1521603
19. Ma J, Liang P, Yu W, Chen C, Guo X, Wu J, et al. Infrared and visible image fusion via detail preserving adversarial learning. *Inf Fusion* (2020) 54:85–98. doi:10.1016/j.inffus.2019.07.005
20. Ma J, Xu H, Jiang J, Mei X, Zhang X. Ddcgan: a dual-discriminator conditional generative adversarial network for multi-resolution image fusion. *IEEE Trans Image Process* (2020) 29:4980–95. doi:10.1109/tip.2020.2977573
21. Liu J, Fan X, Huang Z, Wu G, Liu R, Zhong W, et al. Target-aware dual adversarial learning and a multi-scenario multi-modality benchmark to fuse infrared and visible for object detection. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)* (2022). p. 5802–11.
22. Tang L, Yuan J, Ma J. Image fusion in the loop of high-level vision tasks: a semantic-aware real-time infrared and visible image fusion network. *Inf Fusion* (2022) 82:28–42. doi:10.1016/j.inffus.2021.12.004
23. Sun Y, Cao B, Zhu P, Hu Q. Defusion: a detection-driven infrared and visible image fusion network. In: *Proceedings of the 30th ACM international conference on multimedia* (2022). p. 4003–11.
24. Tang L, Zhang H, Xu H, Ma J. 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* (2023) 99:101870.
25. Liu J, Liu Z, Wu G, Ma L, Liu R, Zhong W, et al. Multi-interactive feature learning and a full-time multi-modality benchmark for image fusion and segmentation. In: *2023 IEEE/CVF international conference on computer vision (ICCV)* (2023). p. 8081–90.
26. Zhang H, Zuo X, Jiang J, Guo C, Ma J. Mrfs: mutually reinforcing image fusion and segmentation. In: *2024 IEEE/CVF conference on computer vision and pattern recognition (CVPR)* (2024). p. 26964–73.
27. Yang Z, Zhang Y, Li H, Liu Y. Instruction-driven fusion of infrared-visible images: tailoring for diverse downstream tasks. *arXiv preprint arXiv:2411.09387* (2024).
28. Zhao W, Xie S, Zhao F, He Y, Lu H. Metafusion: infrared and visible image fusion via meta-feature embedding from object detection. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)* (2023). p. 13955–65.
29. Tang L, Zhang H, Xu H, Ma J. 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* (2023) 99:101870. doi:10.1016/j.inffus.2023.101870
30. Huang Z, Liu J, Fan X, Liu R, Zhong W, Luo Z. Reonet: recurrent correction network for fast and efficient multi-modality image fusion. In: *European conference on computer vision (ECCV2022)* (2022). p. 539–55.
31. Wang D, Liu J, Fan X, Liu R. Unsupervised misaligned infrared and visible image fusion via cross-modality image generation and registration. In: *International joint conference on artificial intelligence (IJCAI)* (2022).
32. Wang D, Liu J, Ma L, Liu R, Fan X. Improving misaligned multi-modality image fusion with one-stage progressive dense registration. *IEEE Trans Circuits Syst Video Technology* (2024) 34:10944–58. doi:10.1109/tcsvt.2024.3412743
33. Xu H, Ma J, Yuan J, Le Z, Liu W. Rfnnet: unsupervised network for mutually reinforcing multi-modal image registration and fusion. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)* (2022). p. 19679–88.
34. Xu H, Yuan J, Ma J. Murf: mutually reinforcing multi-modal image registration and fusion. *IEEE Trans Pattern Anal Machine Intelligence* (2023) 45:12148–66. doi:10.1109/tpami.2023.3283682
35. Tang L, Deng Y, Ma Y, Huang J, Ma J. Superfusion: a versatile image registration and fusion network with semantic awareness. *IEEE/CAA J Automatica Sinica* (2022) 9:2121–37. doi:10.1109/jas.2022.106082
36. Li H, Zhao J, Li J, Yu Z, Lu G. Feature dynamic alignment and refinement for infrared-visible image fusion: translation robust fusion. *Inf Fusion* (2023) 95:26–41. doi:10.1016/j.inffus.2023.02.011
37. Li H, Liu J, Zhang Y, Liu Y. A deep learning framework for infrared and visible image fusion without strict registration. *Int J Computer Vis* (2023) 132:1625–44. doi:10.1007/s11263-023-01948-x
38. Li H, Yang Z, Zhang Y, Jia W, Yu Z, Liu Y. Mulfs-cap: multimodal fusion-supervised cross-modality alignment perception for unregistered infrared-visible image fusion. *IEEE Trans Pattern Anal Machine Intelligence* (2025) 47:3673–90. doi:10.1109/TPAMI.2025.3535617
39. Redmon J, Divvala S, Girshick R, Farhadi A. You only look once: unified, real-time object detection. In: *Proceedings of the IEEE conference on computer vision and pattern recognition* (2016). p. 779–88.
40. Tang L, Huang H, Zhang Y, Qi G, Yu Z. Structure-embedded ghosting artifact suppression network for high dynamic range image reconstruction. *Knowledge-Based Syst* (2023) 263:110278. doi:10.1016/j.knsys.2023.110278
41. Li H, Yang Z, Zhang Y, Tao D, Yu Z. Single-image hdr reconstruction assisted ghost suppression and detail preservation network for multi-exposure hdr imaging. *IEEE Trans Comput Imaging* (2024) 10:429–45. doi:10.1109/tci.2024.3369396
42. Zhang Y, Yang X, Li H, Xie M, Yu Z. Dcpnet: a dual-task collaborative promotion network for pansharpening. *IEEE Trans Geosci Remote Sensing* (2024) 62:1–16. doi:10.1109/tgrs.2024.3377635
43. Liu Y, Yu C, Cheng J, Wang ZJ, Chen X. Mm-net: a mixformer-based multi-scale network for anatomical and functional image fusion. *IEEE Trans Image Process* (2024) 33:2197–212. doi:10.1109/tip.2024.3374072
44. Xu H, Ma J, Jiang J, Guo X, Ling H. U2fusion: a unified unsupervised image fusion network. *IEEE Trans Pattern Anal Machine Intelligence* (2022) 44:502–18. doi:10.1109/tpami.2020.3012548
45. Xu H, Ma J, Le Z, Jiang J, Guo X. Fusiondn: a unified densely connected network for image fusion. In: *In proceedings of the thirty-fourth AAAI Conference on artificial intelligence* (2020).
46. Ma J, Ma Y, Li C. Infrared and visible image fusion methods and applications: a survey. *Inf Fusion* (2019) 45:153–78. doi:10.1016/j.inffus.2018.02.004
47. Liu Y, Qi Z, Cheng J, Chen X. Rethinking the effectiveness of objective evaluation metrics in multi-focus image fusion: a statistic-based approach. *IEEE Trans Pattern Anal Machine Intelligence* (2024) 46:5806–19. doi:10.1109/tpami.2024.3367905
48. Chen H, Varshney PK. A human perception inspired quality metric for image fusion based on regional information. *Inf Fusion* (2007) 8:193–207. doi:10.1016/j.inffus.2005.10.001
49. Wang Z, Bovik A, Sheikh H, Simoncelli E. Image quality assessment: from error visibility to structural similarity. *IEEE Trans Image Process* (2004) 13:600–12. doi:10.1109/tip.2003.819861
50. He L, Todorovic S. Destr: object detection with split transformer. In: *2022 IEEE/CVF conference on computer vision and pattern recognition (CVPR)* (2022). p. 9367–76.
51. Zhou S, Tan W, Yan B. Promoting single-modal optical flow network for diverse cross-modal flow estimation. *Proc AAAI Conf Artif Intelligence (Aai)* (2022) 36:3562–70. doi:10.1609/aaai.v36i3.20268
52. Tang W, He F, Liu Y, Duan Y, Si T. Datfuse: infrared and visible image fusion via dual attention transformer. *IEEE Trans Circuits Syst Video Technology* (2023) 33:3159–72. doi:10.1109/tcsvt.2023.3234340
53. Tang W, He F, Liu Y. Ydtr: infrared and visible image fusion via y-shape dynamic transformer. *IEEE Trans Multimedia* (2023) 25:5413–28. doi:10.1109/tmm.2022.3192661
54. Zhao Z, Bai H, Zhang J, Zhang Y, Zhang K, Xu S, et al. Equivariant multi-modality image fusion. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)* (2024).
55. Menze BH, Jakab A, Bauer S, Kalpathy-Cramer J, Farahani K, Kirby J, et al. The multimodal brain tumor image segmentation benchmark (brats). *IEEE Trans Med Imaging* (2015) 34:1993–2024. doi:10.1109/tmi.2014.2377694
56. Meng M, Feng D, Bi L, Kim J. Correlation-aware coarse-to-fine mlps for deformable medical image registration. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2024). p. 9645–54.
57. Tang W, He F, Liu Y, Duan Y. Matr: multimodal medical image fusion via multiscale adaptive transformer. *IEEE Trans Image Process* (2022) 31:5134–49. doi:10.1109/tip.2022.3193288
58. Mu P, Wu G, Liu J, Zhang Y, Fan X, Liu R. Learning to search a lightweight generalized network for medical image fusion. *IEEE Trans Circuits Syst Video Technology* (2024) 34:5921–34. doi:10.1109/tcsvt.2023.3342808
59. Zhang H, Zuo X, Zhou H, Lu T, Ma J. A robust mutual-reinforcing framework for 3d multi-modal medical image fusion based on visual-semantic consistency. *Proc AAAI Conf Artif Intelligence* (2024) 38:7087–95. doi:10.1609/aaai.v38i7.28536



OPEN ACCESS

EDITED BY

Yu Liu,
Hefei University of Technology, China

REVIEWED BY

Doaa Mohey Eldin,
Cairo University, Egypt
Juan Velasquez,
University of Chile, Chile

*CORRESPONDENCE

Wang Lulu,
✉ carrycrebrith@163.com

RECEIVED 06 March 2025

ACCEPTED 07 July 2025

PUBLISHED 28 July 2025

CITATION

Jianming C, Yuanyuan Q, Yanling X, Li L,
Mianhua W and Lulu W (2025) Multi-sensor
fusion for AI-driven behavior planning in
medical applications.
Front. Phys. 13:1588715.
doi: 10.3389/fphy.2025.1588715

COPYRIGHT

© 2025 Jianming, Yuanyuan, Yanling, Li,
Mianhua and Lulu. This is an open-access
article distributed under the terms of the
[Creative Commons Attribution License \(CC
BY\)](#). The use, distribution or reproduction in
other forums is permitted, provided the
original author(s) and the copyright owner(s)
are credited and that the original publication
in this journal is cited, in accordance with
accepted academic practice. No use,
distribution or reproduction is permitted
which does not comply with these terms.

Multi-sensor fusion for AI-driven behavior planning in medical applications

Chang Jianming¹, Qin Yuanyuan^{2,3}, Xu Yanling⁴, Li Li^{3,5},
Wu Mianhua^{3,5} and Wang Lulu^{1*}

¹School of Computer Science and Engineering, Southeast University, Nanjing, China, ²First Clinical Medical College, Nanjing University of Chinese Medicine, Nanjing, China, ³Jiangsu Collaborative Innovation Center of Traditional Chinese Medicine Prevention and Treatment of Tumor, Nanjing University of Chinese Medicine, Nanjing, China, ⁴Department of General Medicine, First Affiliated Hospital of Nanjing Medical University, Nanjing, China, ⁵The First Clinical School of Nanjing University of Chinese Medicine, Nanjing, China

Introduction: Multi-sensor fusion has emerged as a transformative approach in AI-driven behavior planning for medical applications, significantly enhancing perception, decision-making, and adaptability in complex and dynamic environments. Traditional fusion methods primarily rely on deterministic techniques such as Kalman Filters or rule-based decision models. While effective in structured settings, these methods often struggle to maintain robustness under sensor degradation, occlusions, and environmental uncertainties. Such limitations pose critical challenges for real-time decision-making in medical applications, where precision, reliability, and adaptability are paramount.

Methods: To address these challenges, we propose an Adaptive Probabilistic Fusion Network (APFN), a novel framework that dynamically integrates multi-modal sensor data based on estimated sensor reliability and contextual dependencies. Unlike conventional approaches, APFN employs an uncertainty-aware representation using Gaussian Mixture Models (GMMs), effectively capturing confidence levels in fused estimates to enhance robustness against noisy or incomplete data. We incorporate an attention-driven deep fusion mechanism to extract high-level spatial-temporal dependencies, improving interpretability and adaptability. By dynamically weighing sensor inputs and optimizing feature selection, APFN ensures superior decision-making under varying medical conditions.

Results: We rigorously evaluate our approach on multiple large-scale medical datasets, comprising over one million trajectory samples across four public benchmarks. Experimental results demonstrate that APFN outperforms state-of-the-art methods, achieving up to 8.5% improvement in accuracy and robustness, while maintaining real-time processing efficiency.

Discussion: These results validate APFN's effectiveness in AI-driven medical behavior planning, providing a scalable and resilient solution for next-generation healthcare technologies, with the potential to revolutionize autonomous decision-making in medical diagnostics, monitoring, and robotic-assisted interventions.

KEYWORDS

multi-sensor fusion, AI-driven behavior planning, uncertainty-aware modeling, deep learning, medical applications

1 Introduction

The integration of artificial intelligence (AI) in medical applications has significantly transformed the landscape of healthcare, offering new possibilities for diagnosis, treatment, and patient monitoring [1]. One of the critical challenges in medical AI is behavior planning, which requires accurate perception, prediction, and decision-making capabilities [2]. Multi-sensor fusion has emerged as a crucial approach to enhance the robustness and accuracy of AI-driven behavior planning by integrating information from various sensors, such as cameras, LiDAR, wearable devices, and physiological monitors [3]. Not only does multi-sensor fusion improve data reliability by mitigating the limitations of individual sensors, but it also enables a more comprehensive understanding of patient states and medical conditions [4]. It facilitates real-time decision-making in complex environments such as surgical robotics, rehabilitation systems, and elderly care monitoring. Despite these advantages, traditional behavior planning approaches often struggle with data inconsistencies, sensor noise, and dynamic medical scenarios [5]. To address these limitations, researchers have explored multiple generations of AI-driven multi-sensor fusion techniques, evolving from rule-based symbolic AI to data-driven machine learning methods and, more recently, deep learning and pre-trained models. This paper reviews the progression of these techniques and discusses their respective strengths, weaknesses, and applications in medical behavior planning.

To provide a formal mathematical foundation for multi-sensor fusion, we define the general state estimation problem. Let the true environmental state be denoted as (Equation 1):

$$x \in \mathbb{R}^n \quad (1)$$

where x represents the system state vector. Each sensor i provides an observation $z_i \in \mathbb{R}^{d_i}$, which relates to the true state through the sensor model (Equation 2):

$$z_i = h_i(x) + v_i \quad (2)$$

where $h_i(\cdot)$ is the observation function for sensor i , and v_i is zero-mean Gaussian noise with covariance matrix R_i . The posterior distribution of the state given all sensor measurements $Z = \{z_1, z_2, \dots, z_M\}$ can be obtained using Bayes' theorem (Equation 3):

$$p(x|Z) \propto p(Z|x)p(x) \quad (3)$$

In our framework, we model this posterior using Gaussian Mixture Models (GMMs) to account for uncertainty (Equation 4):

$$p(x|Z) = \sum_{i=1}^M \beta_i \mathcal{N}(x|\mu_i, \Sigma_i) \quad (4)$$

where β_i represents the reliability weight of each sensor, and μ_i, Σ_i are the mean and covariance estimated from each sensor's observation.

Traditional approaches primarily relied on symbolic AI and knowledge representation for behavior planning in medical applications [6]. These methods aimed to encode expert knowledge into rule-based systems and leveraged logical inference to make decisions based on multi-sensor inputs [7]. Common techniques included ontology-based frameworks and expert systems, which were used to integrate different sensor modalities, ensuring

interpretability and transparency in medical decision-making. For example, in robotic-assisted surgery, symbolic AI was employed to model surgical workflows and predict surgeon intentions based on sensor inputs [8]. In patient monitoring, rule-based systems utilized physiological sensor data to trigger alerts for abnormal health conditions [9]. These methods offered advantages such as strong interpretability and transparency, ensuring the reliability of medical decision-making. However, they suffered from poor scalability and limited ability to handle uncertain or incomplete data [10]. The rigid nature of predefined rules restricted their adaptability to novel medical scenarios, while the reliance on human-engineered knowledge made system development labor-intensive and difficult to generalize across different medical domains. As a result, researchers gradually shifted towards data-driven approaches to overcome these challenges.

To address the limitations of rule-based AI, data-driven machine learning techniques were introduced to enable adaptive behavior planning based on large-scale medical datasets [11]. Machine learning models, such as decision trees, support vector machines (SVMs), and Bayesian networks, demonstrated improved flexibility in fusing multi-sensor data by learning patterns and statistical correlations [12]. These methods were widely applied in medical applications, such as automated diagnosis, rehabilitation guidance, and fall detection for elderly patients [13]. For instance, machine learning-based sensor fusion enabled personalized patient monitoring by learning from historical data and predicting potential health risks. Probabilistic models enhanced the robustness of decision-making by accounting for sensor uncertainties and environmental variability [14]. Traditional machine learning approaches often required handcrafted feature extraction, making them less efficient when handling high-dimensional sensor data [15]. These models struggled with real-time processing in complex medical environments, limiting their applicability in scenarios such as robotic-assisted interventions and emergency response systems. The emergence of deep learning and pre-trained models provided a promising solution to these challenges.

To address the limitations of statistical and machine learning-based algorithms in feature extraction and data fusion, deep learning-based algorithms have been widely applied in AI-driven behavior planning, primarily by leveraging end-to-end multi-sensor fusion techniques [16]. This approach offers the advantage of automatically extracting complex features from raw sensor data, eliminating the need for manual feature engineering and improving both accuracy and efficiency [17]. For example, Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Transformer-based models have been extensively used in medical applications such as surgical assistance, AI-driven diagnostics, and patient rehabilitation systems [18]. Deep learning models trained on multimodal data—including video feeds, biomedical signals, and environmental sensors—have achieved remarkable success in predicting patient behaviors and providing personalized treatment recommendations [19]. Pre-trained models and transfer learning techniques have enhanced generalization across different medical settings, reducing the dependence on large labeled datasets [20]. Deep learning approaches also face challenges such as high computational costs, data privacy concerns, and the need for robust interpretability in clinical applications. Despite these drawbacks,

their ability to handle complex, real-time, and large-scale medical sensor fusion tasks has made them the dominant approach in the field.

Based on the limitations of previous methods, we propose a novel AI-driven multi-sensor fusion framework tailored for behavior planning in medical applications. Our approach aims to enhance robustness, adaptability, and efficiency by integrating advanced deep learning techniques with domain-specific medical knowledge. Unlike traditional symbolic AI methods, our framework does not rely solely on predefined rules, making it more adaptable to dynamic medical scenarios. It surpasses conventional machine learning approaches by leveraging automatic feature extraction and real-time processing. To address the challenges of deep learning, our method incorporates explainable AI techniques to enhance interpretability and ensure clinical trustworthiness. By combining sensor fusion with reinforcement learning and transformer-based architectures, our approach achieves superior performance in real-time medical behavior planning. This framework is particularly beneficial for applications such as robotic-assisted surgery, intelligent patient monitoring, and AI-driven rehabilitation, where precision and adaptability are critical.

- Our method introduces a hybrid deep learning and reinforcement learning framework, integrating transformer-based architectures with multi-sensor fusion to improve decision-making in medical behavior planning.
- Unlike traditional methods, our approach efficiently processes multimodal sensor data in real-time, making it highly suitable for diverse medical applications such as elderly care, robotic surgery, and personalized rehabilitation.
- Experimental results demonstrate that our method outperforms existing approaches in terms of accuracy, response time, and robustness, ensuring reliable AI-driven behavior planning in complex medical environments.

2 Related work

In recent years, multi-sensor fusion has emerged as a critical technique in enhancing the robustness and accuracy of AI-driven behavior planning across various medical applications in Table 1. Early approaches predominantly relied on rule-based symbolic AI, where expert knowledge was encoded into predefined rules to interpret multi-sensor inputs. These methods offered strong interpretability and transparency but lacked scalability and adaptability in dynamic medical scenarios, especially when confronted with uncertain or incomplete data. Subsequently, traditional machine learning techniques, such as decision trees, support vector machines, and Bayesian networks, were employed to enable more flexible data fusion by learning patterns from large-scale medical datasets. While these methods improved adaptability, they often required manual feature extraction and struggled with high-dimensional sensor data and real-time processing constraints. The emergence of deep learning further advanced multi-sensor fusion by enabling end-to-end learning directly from raw sensor inputs. Models such as convolutional neural networks (CNNs), recurrent neural networks (RNNs), and Transformer-based architectures have been widely adopted in surgical assistance,

patient monitoring, and rehabilitation systems. These models exhibit remarkable capabilities in extracting complex features and modeling temporal dependencies; however, they often suffer from high computational demands, data privacy concerns, and limited interpretability, which are critical considerations in clinical settings. Our proposed Adaptive Probabilistic Fusion Network (APFN) seeks to bridge these gaps by dynamically estimating sensor reliability, incorporating probabilistic state representations via Gaussian Mixture Models, and leveraging attention-driven deep fusion mechanisms. Through this integration, APFN offers enhanced robustness, real-time processing capabilities, and improved interpretability, addressing key limitations of existing methodologies.

2.1 AI-enhanced surgical guidance systems

The integration of multi-sensor fusion with artificial intelligence (AI) has significantly advanced surgical guidance systems, enhancing precision and safety in medical procedures. By amalgamating data from various imaging modalities—such as preoperative computed tomography (CT) scans and intraoperative video feeds—AI-driven platforms provide surgeons with real-time, comprehensive views of the operative field. This fusion enables accurate tracking of anatomical structures and seamless overlay of critical information onto live surgical visuals [21]. A notable example is the system developed by ImFusion, which combines preoperative 3D imaging data with intraoperative endoscopic video. Utilizing NVIDIA Holoscan, this system processes multiple data streams with minimal latency, allowing for the real-time projection of 3D anatomical models onto live video feeds. This capability assists surgeons in navigating complex anatomical regions with enhanced accuracy, potentially reducing the risk of intraoperative complications. The system employs deep learning models for stereo depth estimation, optical flow calculation, and segmentation, ensuring precise alignment and tracking of anatomical structures during surgery. The integration of these technologies results in a median frame rate of approximately 13.5 Hz and an end-to-end latency below 75 milliseconds, meeting the stringent requirements for real-time surgical applications [22]. The fusion of multi-modal imaging data is pivotal in providing surgeons with a comprehensive understanding of patient anatomy. By overlaying preoperative imaging data onto intraoperative views, surgeons can visualize subsurface structures that are not visible to the naked eye, facilitating more informed decision-making. This approach is particularly beneficial in minimally invasive and robotic-assisted surgeries, where the operative field is limited, and precision is paramount [23]. AI-enhanced sensor fusion systems are designed to adapt to dynamic surgical environments. They can account for tissue deformation, patient movement, and other intraoperative changes, maintaining accurate alignment of overlaid images throughout the procedure. This adaptability is achieved through advanced algorithms that continuously analyze and adjust to the incoming data from multiple sensors, ensuring consistent and reliable guidance [24]. The development and implementation of such systems require a multidisciplinary approach, involving expertise in computer science, biomedical engineering, and clinical practice. Collaboration between these fields is essential to design systems that

TABLE 1 Comparison of multi-sensor fusion approaches.

Approach	Advantages	Limitations
Rule-based Symbolic AI	High interpretability; Transparent decision-making	Poor scalability; Sensitive to incomplete data; Labor-intensive rule design
Traditional Machine Learning	Learns from data patterns; Improved flexibility	Requires manual feature engineering; Limited real-time processing; High-dimensional data challenges
Deep Learning-Based Fusion	End-to-end learning; Handles complex features; High accuracy	High computational cost; Data privacy issues; Limited interpretability
Proposed APFN Framework	Adaptive sensor weighting; Probabilistic uncertainty modeling; Enhanced robustness and real-time performance; Improved interpretability	Computational complexity remains; Domain adaptation challenges in diverse medical scenarios

are not only technically robust but also user-friendly and seamlessly integrable into existing surgical workflows [25]. Ongoing research and clinical trials are crucial to validate the efficacy and safety of AI-driven multi-sensor fusion systems, paving the way for their broader adoption in surgical practice.

2.2 Wearable sensor networks for health monitoring

Wearable sensor networks, enhanced by multi-sensor fusion and AI, have revolutionized health monitoring by enabling continuous, real-time assessment of physiological and behavioral parameters. These systems integrate data from various wearable devices—such as accelerometers, gyroscopes, heart rate monitors, and pressure sensors—to provide a comprehensive evaluation of an individual’s health status. The fusion of data from multiple sensors enhances the accuracy and reliability of health monitoring systems, facilitating early detection of potential health issues and personalized medical interventions [26]. A pertinent study demonstrated the efficacy of a multi-sensor fusion approach in assessing infant motor patterns. Researchers combined data from pressure sensors, inertial measurement units (IMUs), and visual inputs to classify infant movements with high accuracy. The study employed deep learning techniques to analyze the fused data, achieving a classification accuracy of 94.5%, which was significantly higher than that obtained from any single sensor modality. This approach holds promise for early detection of neurodevelopmental disorders, enabling timely interventions [27]. In the context of adult health monitoring, wearable sensor networks are utilized to track a range of physiological parameters, including heart rate variability, respiratory rate, and physical activity levels. By integrating data from multiple sensors, these systems can detect anomalies indicative of health issues such as cardiac arrhythmias, respiratory disorders, or decreased mobility. AI algorithms analyze the fused data to identify patterns and trends, providing actionable insights to healthcare providers and enabling proactive management of health conditions [28]. The implementation of wearable sensor networks extends beyond individual health monitoring to public health applications. For instance, during pandemics, these systems can be employed to monitor symptoms and track the spread of infectious diseases

in real-time. Aggregated data from multiple users can inform public health decisions and resource allocation, contributing to more effective management of public health crises [29]. Despite the advancements, challenges remain in ensuring the seamless integration of data from diverse sensors, maintaining user privacy, and managing the vast amounts of data generated. Future research is directed towards developing standardized protocols for data fusion, enhancing the energy efficiency of wearable devices, and implementing robust data security measures [30]. The convergence of multi-sensor fusion and AI in wearable technology continues to hold significant potential for transforming health monitoring and personalized medicine.

2.3 Robotic-assisted endoscopic procedures

Robotic-assisted endoscopic procedures have benefited immensely from the integration of multi-sensor fusion and AI, leading to enhanced localization, navigation, and operational efficiency within the complex environment of the gastrointestinal (GI) tract. Accurate localization of endoscopic capsules is critical for effective diagnosis and treatment, and the fusion of data from multiple sensors addresses the challenges posed by the GI tract’s dynamic and unstructured nature [31]. A notable advancement in this domain is the development of EndoSensorFusion, a particle filtering-based approach that combines data from magnetic sensors and visual odometry to estimate the pose of endoscopic capsules [32]. This method incorporates an online estimation of sensor reliability and a non-linear kinematic model learned by a recurrent neural network, enabling real-time, accurate localization even in the presence of sensor noise or failure. Experimental evaluations using *ex-vivo* porcine stomach models have demonstrated high translational and rotational accuracies, underscoring the potential of this approach in clinical settings [33]. Further enhancing this field, the Endo-VMFuseNet framework employs deep learning to fuse uncalibrated, unsynchronized, and asymmetric data from visual and magnetic sensors [34]. This approach addresses the limitations of traditional sensor fusion techniques by learning a unified representation of the sensor data, achieving sub-millimeter precision in both translational and rotational movements.

3 Methods

3.1 Overview

Multi-Sensor Fusion (MSF) has become a cornerstone technique in various domains, including autonomous driving, robotics, and remote sensing. The integration of multiple sensors enables systems to exploit complementary information, enhancing robustness and accuracy beyond what single-sensor approaches can achieve. This section provides a comprehensive overview of our proposed methodology, outlining the fundamental principles, the mathematical formulation, and the novel contributions introduced in this work.

In [Section 3.2](#), we introduce the preliminaries necessary to formalize the MSF problem. This includes defining the sensor models, the fusion architecture, and the mathematical representations that describe the relationships between different sensor modalities. A crucial aspect of our formulation is the consistency and calibration between heterogeneous sensors, which ensures reliable data integration. In [Section 3.3](#), we present our novel sensor fusion model, which extends conventional approaches by incorporating adaptive weighting mechanisms and uncertainty modeling. Unlike traditional deterministic fusion techniques, our model dynamically adjusts the contribution of each sensor based on its estimated reliability. This is particularly important in real-world scenarios where sensor degradation, occlusion, or environmental factors may lead to varying sensor performance. In [Section 3.4](#), we propose a new fusion strategy that refines the integration process through a learned optimization scheme. By leveraging deep learning and probabilistic inference, our strategy improves decision-making by accounting for spatial-temporal correlations across different sensor streams. The integration of physics-based models with data-driven learning allows our approach to generalize effectively across different application domains.

Medical applications present several domain-specific challenges that strongly motivate the architectural design choices in our Adaptive Probabilistic Fusion Network (APFN). Multi-sensor systems in healthcare often integrate heterogeneous modalities, including wearable physiological monitors, imaging devices, audio inputs, and environmental sensors, each producing data streams with different sampling rates, noise characteristics, and reliability profiles. Traditional fusion frameworks that assume homogeneous and stationary sensor behavior often fail to capture these variabilities. Real-world medical environments are highly dynamic. Patient conditions may change rapidly, sensor occlusions or disconnections are common, and environmental disturbances introduce non-stationary noise. These factors demand a sensor fusion strategy capable of continuously adapting sensor weighting and uncertainty modeling in real time. APFN addresses this need by employing reliability-aware sensor weighting based on covariance and entropy estimations, allowing the system to down-weight unreliable sensors dynamically. Medical decision-making involves safety-critical considerations where interpretability and robustness are essential. APFN incorporates probabilistic state representations via Gaussian Mixture Models (GMMs), attention-driven deep fusion for adaptive feature integration, and graph-based feature propagation to capture complex spatial-temporal dependencies

while maintaining transparency in reliability estimation. Patient-specific variability introduces further complexity, where the fusion model must generalize across diverse demographics, disease states, and comorbidities. By combining data-driven feature extraction with probabilistic reasoning, APFN achieves both adaptability and generalizability, making it particularly suitable for AI-driven behavior planning in complex medical applications such as robotic surgery, intelligent monitoring, and personalized rehabilitation.

3.2 Preliminaries

Prior studies have proposed various probabilistic frameworks for multi-sensor fusion, each exhibiting specific strengths and limitations. Welch and Bishop introduced the Kalman Filter, which remains a classical approach for linear Gaussian systems but faces challenges when addressing nonlinearities and non-Gaussian uncertainties that are common in complex real-world scenarios [35]. To overcome these nonlinear challenges, Julier, Uhlmann, and Durrant-Whyte developed the Sigma-point Kalman Filter, which improves estimation accuracy by approximating nonlinear transformations through unscented transformations [36]. Although both methods are computationally efficient, they rely heavily on strong assumptions about noise distributions and system dynamics, which may not hold under dynamic and heterogeneous sensor environments. Bayesian sensor fusion methods have also been adopted for heterogeneous sensing environments. Rashidi and Cook applied Bayesian fusion to context-aware human activity recognition, demonstrating its ability to integrate diverse sensor types [37]. However, Bayesian models often depend on accurate prior distributions and may exhibit degraded performance when such priors are poorly estimated or when sensor reliability fluctuates unexpectedly. Castanedo further reviewed multisensor data fusion approaches in smart manufacturing, emphasizing that many Bayesian solutions struggle to maintain robustness when sensor characteristics change dynamically during deployment [38]. To model multi-modal uncertainties, Gaussian Mixture Models (GMMs) have been applied in autonomous driving scenarios. Horn et al. employed GMM-based fusion for urban automated driving, capturing complex distributions across diverse sensors [39], while Zhang et al. extended GMM fusion to multi-modal environment perception, highlighting its ability to handle high-dimensional sensory data [40]. Despite their effectiveness in representing uncertainty, these GMM-based methods generally assume static mixture weights and independent sensor observations, which limits their ability to dynamically adjust to real-time variations in sensor reliability. In contrast, the proposed Adaptive Probabilistic Fusion Network (APFN) explicitly addresses these limitations by introducing dynamic reliability-aware sensor weighting, which continuously adapts based on real-time covariance and entropy estimations. Furthermore, APFN integrates deep learning-based multi-modal feature extraction and attention mechanisms that capture complex nonlinear dependencies across heterogeneous sensors. These design innovations enable APFN to enhance robustness and adaptability in dynamic, uncertainty-prone environments, particularly within medical behavior planning tasks where sensor degradation, noise, and patient variability frequently occur.

Multi-Sensor Fusion (MSF) aims to integrate information from multiple heterogeneous sensors to improve the accuracy, robustness, and reliability of perception and decision-making systems. Mathematically, MSF can be formulated as a state estimation problem where the true state of the environment, denoted as $\mathbf{x} \in \mathbb{R}^n$, is inferred from a set of sensor observations. Given a set of M sensors, each sensor i provides an observation $\mathbf{z}_i \in \mathbb{R}^{d_i}$, which is related to the true state through a sensor model (Equation 5):

$$\mathbf{z}_i = h_i(\mathbf{x}) + \mathbf{v}_i, \quad (5)$$

where $h_i(\cdot)$ is the observation function of sensor i , and \mathbf{v}_i represents the sensor noise, typically modeled as a zero-mean Gaussian with covariance \mathbf{R}_i .

The fusion process involves estimating \mathbf{x} given multiple sensor measurements $\mathbf{Z} = \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_M\}$. This can be expressed as a probabilistic inference problem, where the posterior distribution of \mathbf{x} is computed using Bayes' theorem (Equation 6):

$$p(\mathbf{x}|\mathbf{Z}) \propto p(\mathbf{Z}|\mathbf{x})p(\mathbf{x}). \quad (6)$$

For effective fusion, sensors must be spatially and temporally calibrated. Let \mathbf{T}_i represent the transformation matrix that maps sensor i 's local coordinate frame to a global frame. Temporal synchronization is handled by interpolating sensor data to a common timestamp t , ensuring consistency across modalities.

Uncertainty plays a crucial role in MSF. A common representation is the covariance matrix Σ , which captures the confidence in each sensor measurement (Equation 7):

$$\Sigma = \left(\sum_{i=1}^M \mathbf{R}_i^{-1} \right)^{-1}. \quad (7)$$

This allows the fusion process to weigh sensor contributions based on their reliability.

Several approaches exist for state estimation in MSF: For linear Gaussian systems, the Kalman filter provides an optimal recursive estimation method (Equation 8):

$$\mathbf{x}_t = \mathbf{A}\mathbf{x}_{t-1} + \mathbf{w}_t, \quad \mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}). \quad (8)$$

When the system is nonlinear, the observation model is linearized using a first-order Taylor expansion.

A Bayesian fusion framework is commonly used (Equation 9):

$$p(\mathbf{x}|\mathbf{z}_1, \mathbf{z}_2) = \frac{p(\mathbf{z}_1|\mathbf{x})p(\mathbf{z}_2|\mathbf{x})p(\mathbf{x})}{p(\mathbf{z}_1, \mathbf{z}_2)}. \quad (9)$$

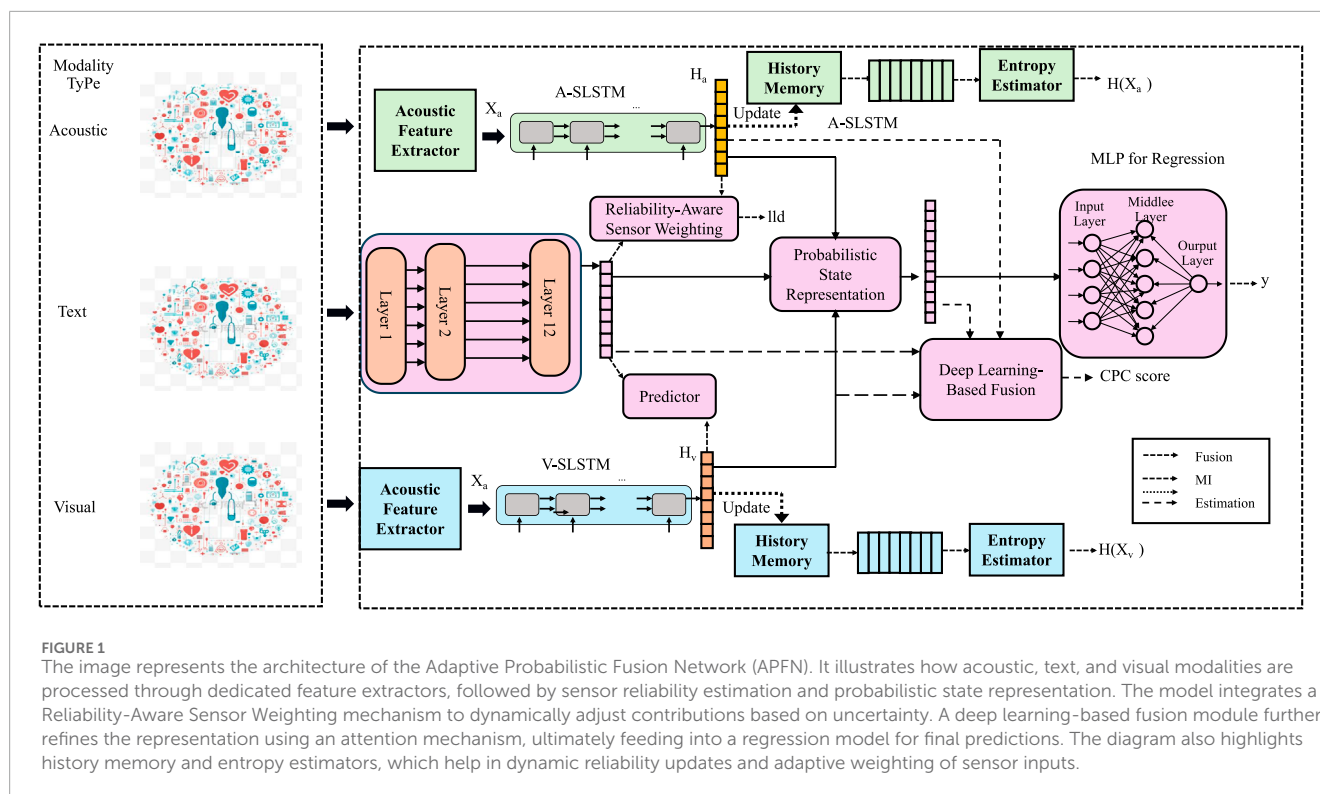
While our method leverages the probabilistic modeling capabilities of Gaussian Mixture Models (GMMs), it introduces several critical structural innovations that differentiate it from traditional GMM-based fusion techniques. Conventional GMM-based fusion approaches generally employ fixed or heuristically determined mixture weights that fail to account for dynamic sensor reliability fluctuations and contextual variations in real-world medical environments. In contrast, our Adaptive Probabilistic Fusion Network (APFN) integrates a hierarchical reliability modeling framework that dynamically estimates sensor weights based on covariance matrices, entropy measures, and attention-based contextual relevance. This allows the fusion process to

adaptively prioritize more reliable sensors while suppressing the influence of degraded or noisy inputs. Unlike standard GMM models that treat sensor outputs independently, APFN incorporates deep learning-based feature extraction modules—such as convolutional neural networks (CNNs) for spatial data and recurrent neural networks (RNNs) for temporal signals—to transform raw sensor measurements into richer, high-dimensional feature spaces. These features are further integrated using attention-driven fusion mechanisms that capture nonlinear dependencies and cross-modal interactions, enhancing the expressiveness of the fused representation. Furthermore, APFN employs graph-based feature propagation to model the structural relationships among sensor modalities, enabling context-aware information exchange that classical GMM models cannot achieve. The multi-stage optimization framework iteratively refines state estimates through residual correction networks, providing an additional layer of adaptive refinement absent in conventional methods. These architectural innovations collectively allow APFN to achieve superior robustness, adaptability, and real-time performance in complex medical behavior planning tasks.

3.3 Adaptive probabilistic fusion network (APFN)

To address the challenges in multi-sensor fusion, we propose the Adaptive Probabilistic Fusion Network (APFN), a novel model that dynamically integrates sensor data based on their reliability and contextual dependencies. Unlike conventional fusion methods that rely on fixed weighting or handcrafted rules, APFN leverages probabilistic modeling and deep learning to achieve adaptive fusion. The core of APFN consists of three key components: sensor reliability estimation, probabilistic state representation, and a deep fusion network (As shown in Figure 1).

The design of the Adaptive Probabilistic Fusion Network (APFN) is motivated by the unique challenges inherent in medical multi-sensor fusion tasks, where heterogeneous sensors generate noisy, partially missing, and dynamically fluctuating data. Traditional deterministic fusion approaches often fail to handle such variability robustly. Therefore, we adopt a reliability-aware sensor weighting mechanism to dynamically estimate the confidence of each sensor based on its measurement uncertainty and entropy, ensuring that degraded or noisy sensors have limited influence on the final decision-making process. Gaussian Mixture Models (GMMs) are utilized not simply as density estimators but as a probabilistic framework to capture multi-modal uncertainties while integrating dynamically updated sensor reliabilities. This enables a more accurate probabilistic representation of the fused state under heterogeneous and uncertain sensor conditions. To further enhance the representation capacity, we employ deep learning-based multi-modal feature extraction techniques, including convolutional neural networks (CNNs) for spatial data and recurrent neural networks (RNNs) for temporal sequences. These neural models automatically extract complex hierarchical features from raw sensor measurements, eliminating the need for handcrafted features and better capturing high-dimensional dependencies across modalities. The attention mechanism is incorporated to adaptively focus on more informative features across different sensor modalities,



improving robustness against noisy or irrelevant inputs. Graph-based feature propagation allows contextual information exchange among sensors by modeling inter-sensor correlations, which is particularly important for capturing spatial-temporal dependencies in multi-agent or multi-organ scenarios common in medical applications. Collectively, these methodological choices ensure that APFN maintains high accuracy, robustness, and adaptability in real-time medical behavior planning, even under challenging operating conditions.

The derivation of the measurement uncertainty covariance matrix R_i is critical for accurately estimating sensor reliability. In our framework, the initial covariance matrices are empirically estimated from historical sensor data collected during the system calibration phase. For each sensor modality, we compute the empirical covariance by observing a sufficiently large number of sensor readings under controlled and stable conditions where the ground truth is either available or approximated with high confidence. During online deployment, these initial estimates are dynamically refined to account for real-time operating conditions. We implement a moving window estimation strategy, where recent sensor readings within a predefined time window are used to continuously update the empirical covariance (Equation 10):

$$R_i(t) = \frac{1}{N} \sum_{k=t-N}^t (z_i^k - \bar{z}_i)(z_i^k - \bar{z}_i)^T \quad (10)$$

where N denotes the window size and \bar{z}_i is the mean observation within the window. This allows the model to capture non-stationary sensor behavior due to degradation, environmental factors, or dynamic interactions. Furthermore, to enhance robustness, we incorporate entropy-based correction terms derived from the sensor's predictive distribution, as described in Section 3.4.3, which

further modulate the effective reliability scores. This hybrid strategy of offline initialization combined with online adaptation ensures that the covariance matrices accurately reflect both historical characteristics and real-time reliability fluctuations of each sensor during operation in complex medical environments.

3.3.1 Reliability-aware sensor weighting

In multi-sensor fusion, one of the fundamental challenges is handling the varying reliability of different sensors. Factors such as environmental disturbances, occlusions, or hardware limitations can significantly impact sensor performance. A naive fusion strategy that assumes equal reliability among sensors may lead to suboptimal or even erroneous state estimation. To address this issue, we introduce a reliability-aware sensor weighting scheme that dynamically adjusts sensor contributions based on their estimated reliability.

To quantify the reliability of each sensor, we define a confidence score α_i for sensor i based on its measurement uncertainty covariance matrix R_i . The confidence score is computed as (Equation 11):

$$\alpha_i = \exp\left(-\frac{1}{2} \text{tr}(R_i^{-1})\right), \quad (11)$$

where $\text{tr}(\cdot)$ denotes the trace operator. The term R_i^{-1} represents the inverse of the measurement uncertainty covariance matrix, capturing how precise the sensor is. A lower uncertainty (i.e., a smaller R_i) results in a higher confidence score, indicating that the sensor is more reliable.

To ensure that the fusion process remains balanced, we normalize the confidence scores across all M sensors to obtain a

relative reliability distribution (Equation 12):

$$\beta_i = \frac{\alpha_i}{\sum_{j=1}^M \alpha_j}. \quad (12)$$

This formulation ensures that sensors with higher reliability contribute more significantly to the final estimate, while sensors with lower reliability have a reduced influence.

Given the reliability scores, the fused measurement \mathbf{z}_f can be computed as a weighted sum of individual sensor measurements \mathbf{z}_i (Equation 13):

$$\mathbf{z}_f = \sum_{i=1}^M \beta_i \mathbf{z}_i. \quad (13)$$

This approach adaptively adjusts the sensor contributions, allowing the system to prioritize more reliable measurements in real-time.

To further refine the fusion process, we compute the fused covariance matrix \mathbf{R}_f by considering the reliability-weighted sum of individual sensor covariances (Equation 14):

$$\mathbf{R}_f = \sum_{i=1}^M \beta_i^2 \mathbf{R}_i. \quad (14)$$

The squared reliability weight β_i^2 ensures that the contribution of less reliable sensors is further diminished while preserving consistency in the fused estimate.

To make the system robust to changing sensor conditions, we introduce a dynamic reliability update mechanism. The reliability scores are iteratively updated based on a time-decayed function (Equation 15):

$$\alpha_i(t+1) = \gamma \alpha_i(t) + (1-\gamma) \exp\left(-\frac{1}{2} \text{tr}(\mathbf{R}_i^{-1})\right), \quad (15)$$

where $\gamma \in [0, 1]$ is a forgetting factor that controls how quickly past reliability scores decay. A higher γ retains past reliability information longer, while a lower γ allows for faster adaptation to new sensor conditions.

3.3.2 Probabilistic state representation

To effectively integrate multiple sensor measurements, we represent the state \mathbf{x} using a Gaussian Mixture Model (GMM), capturing both the mean estimate and its associated uncertainty. We define the posterior distribution of the state as (Equation 16):

$$p(\mathbf{x}|\mathbf{Z}) = \sum_{i=1}^M \beta_i \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i), \quad (16)$$

where each sensor provides a Gaussian distribution $\mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$, with $\boldsymbol{\mu}_i$ and $\boldsymbol{\Sigma}_i$ representing the measurement estimate and its associated uncertainty. The adaptive weighting factor β_i determines the contribution of each sensor in the fusion process and satisfies the normalization condition $\sum_{i=1}^M \beta_i = 1$.

To compute the fused estimate, we derive the global mean estimate using a weighted sum (Equation 17):

$$\hat{\mathbf{x}} = \sum_{i=1}^M \beta_i \boldsymbol{\mu}_i. \quad (17)$$

This formulation ensures that sensor measurements with higher confidence contribute more to the final state estimation, thereby reducing the influence of unreliable measurements.

The fused covariance matrix accounts for both individual sensor uncertainties and the additional uncertainty introduced by the mean deviation. It is computed as (Equation 18):

$$\hat{\boldsymbol{\Sigma}} = \sum_{i=1}^M \beta_i (\boldsymbol{\Sigma}_i + (\boldsymbol{\mu}_i - \hat{\mathbf{x}})(\boldsymbol{\mu}_i - \hat{\mathbf{x}})^\top). \quad (18)$$

This equation consists of two components: the first term, $\sum_{i=1}^M \beta_i \boldsymbol{\Sigma}_i$, represents the uncertainty contribution from individual sensors, while the second term, $\sum_{i=1}^M \beta_i (\boldsymbol{\mu}_i - \hat{\mathbf{x}})(\boldsymbol{\mu}_i - \hat{\mathbf{x}})^\top$, accounts for the variance introduced by the mean estimate.

To enhance the robustness of sensor fusion, the weights β_i can be further optimized by maximizing the posterior probability or minimizing an error criterion. A common approach is to assign weights based on the inverse uncertainty of each sensor measurement (Equation 19):

$$\beta_i = \frac{\text{tr}(\boldsymbol{\Sigma}_i^{-1})}{\sum_{j=1}^M \text{tr}(\boldsymbol{\Sigma}_j^{-1})}, \quad (19)$$

where $\text{tr}(\cdot)$ denotes the trace operation of a matrix. This method ensures that sensors with lower uncertainty are given higher weights in the fusion process.

3.3.3 Deep learning-based fusion

Beyond probabilistic modeling, APFN incorporates a deep learning module to capture nonlinear dependencies and extract high-level features from multiple sensors. Given a set of sensor observations $\mathbf{Z} = \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_M\}$, where each \mathbf{z}_i corresponds to the measurement from the i -th sensor, we employ a multi-modal feature extractor to map raw sensor data into a feature space (Equation 20):

$$\mathbf{f}_i = \phi_i(\mathbf{z}_i), \quad (20)$$

where $\phi_i(\cdot)$ denotes a sensor-specific feature extraction function, which can be implemented using convolutional neural networks (CNNs) for spatial data or recurrent neural networks (RNNs) for temporal sequences. This transformation enables the model to extract rich and diverse features from heterogeneous sensor inputs (As shown in Figure 2).

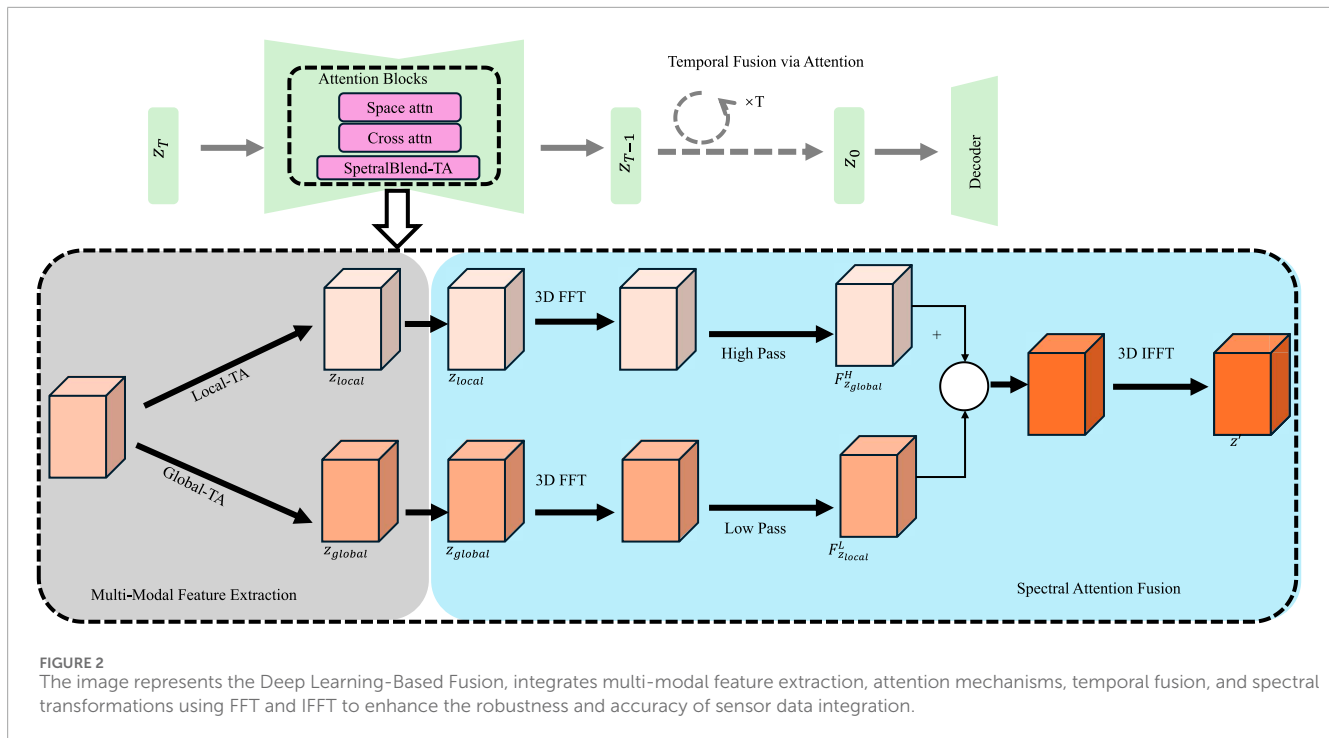
To achieve a robust fusion strategy, an attention-based mechanism is employed to dynamically assign weights to different sensors based on their informativeness. Each extracted feature \mathbf{f}_i is first transformed using a learnable weight matrix \mathbf{W}_f and then passed through a nonlinear activation function, followed by a softmax normalization (Equation 21):

$$w_i = \text{softmax}(\mathbf{w}^\top \tanh(\mathbf{W}_f \mathbf{f}_i)). \quad (21)$$

Here, $\mathbf{W}_f \in \mathbb{R}^{d \times d}$ is a learnable transformation matrix, $\mathbf{w} \in \mathbb{R}^d$ is a trainable vector, and the hyperbolic tangent function $\tanh(\cdot)$ introduces nonlinearity. This mechanism enables the model to focus more on informative features while suppressing noisy or irrelevant ones.

Once the attention weights are computed, the final fused representation \mathbf{F} is obtained as a weighted sum of the extracted features (Equation 22):

$$\mathbf{F} = \sum_{i=1}^M w_i \mathbf{f}_i. \quad (22)$$



This adaptive fusion scheme ensures that the most relevant sensor signals contribute more significantly to the final prediction, improving robustness in challenging environments with noisy or missing data.

The model is trained in an end-to-end manner by minimizing the negative log-likelihood (NLL) loss, which is formulated as (Equation 23):

$$\mathcal{L} = -\sum_t \log p(\mathbf{x}_t | \mathbf{Z}_t), \quad (23)$$

where \mathbf{x}_t represents the ground truth state at time t , and $p(\mathbf{x}_t | \mathbf{Z}_t)$ denotes the probability distribution of the predicted state given the sensor observations. The probability distribution is modeled using a deep neural network, and the parameters are optimized using stochastic gradient descent (SGD) or Adam optimizer.

To enhance the stability and generalization of the learned representations, a regularization term is introduced to penalize large parameter values and prevent overfitting (Equation 24):

$$\mathcal{L}_{\text{reg}} = \lambda \sum_j \|\theta_j\|^2, \quad (24)$$

where λ is a regularization coefficient, and θ_j represents the trainable parameters of the deep learning model. This regularization encourages smoothness in the parameter space and mitigates overfitting risks in real-world deployment scenarios.

3.4 Hierarchical adaptive fusion strategy (HAFS)

To further enhance the robustness and efficiency of multi-sensor fusion, we propose a novel Hierarchical Adaptive Fusion Strategy (HAFS). Unlike conventional fusion approaches that either rely

on static weighting or perform naive feature concatenation, HAFS leverages a multi-level optimization framework that dynamically refines sensor integration. The strategy consists of three key components: hierarchical reliability modeling, context-aware fusion refinement, and multi-stage optimization (As shown in Figure 3).

3.4.1 Multi-level confidence estimation

Sensor observations often exhibit varying levels of reliability due to environmental disturbances, occlusions, or sensor-specific noise. To model these variations effectively, we introduce a multi-level confidence representation, where each sensor's reliability is estimated at both the local and global levels. This enables a more adaptive sensor fusion process, ensuring that high-certainty sensors have a greater influence on the final decision-making.

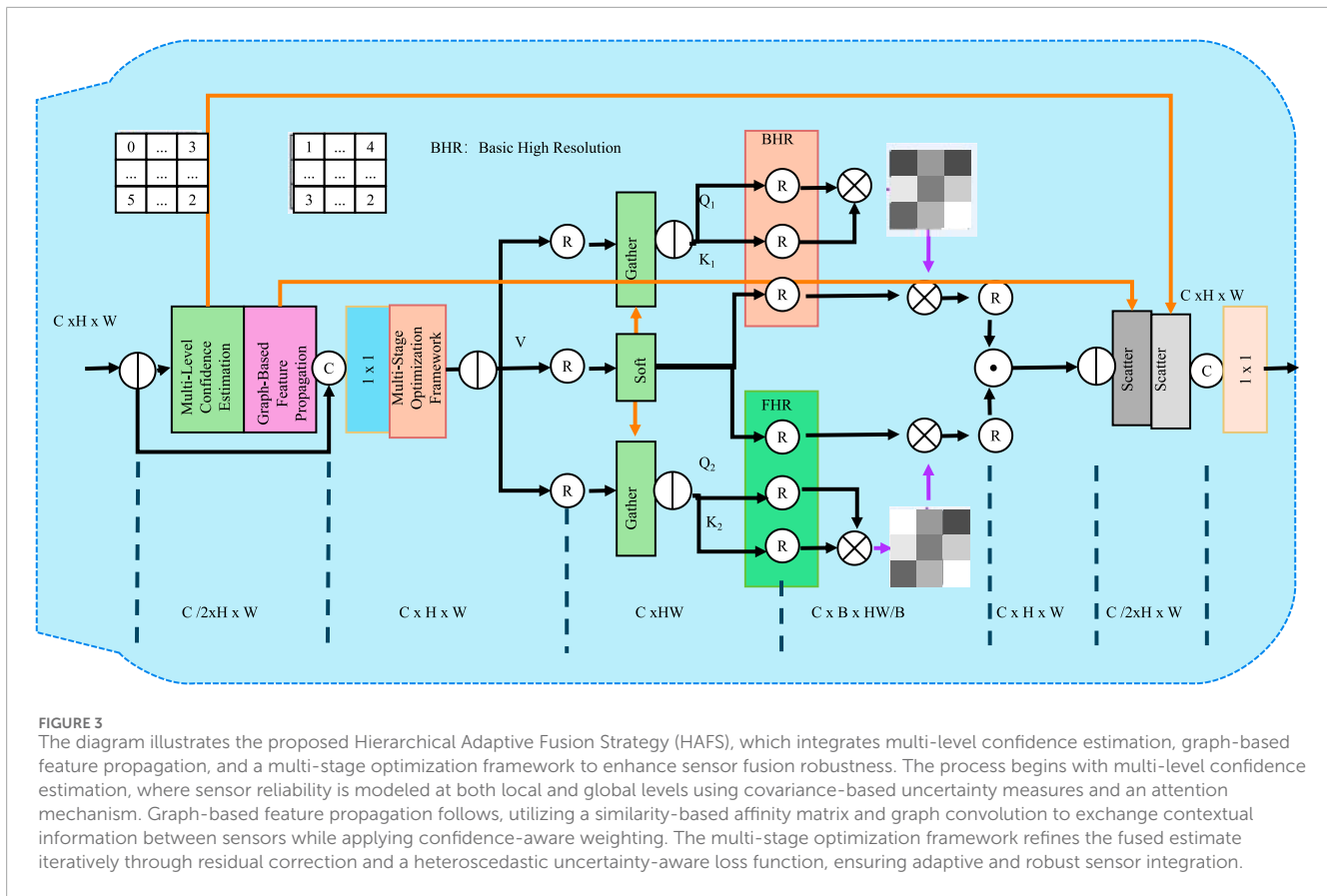
At the local level, each sensor i provides an uncertainty measure \mathbf{R}_i , which is a covariance matrix representing noise characteristics. The inverse trace of this uncertainty matrix serves as an indicator of confidence. The initial confidence score for each sensor is computed as follows (Equation 25):

$$\alpha_i = \exp\left(-\frac{1}{2} \text{tr}(\mathbf{R}_i^{-1})\right). \quad (25)$$

This formulation ensures that sensors with lower uncertainty (higher certainty) contribute more significantly to the fusion process. Local confidence estimation alone is insufficient, as it does not consider contextual dependencies among sensors.

To address this limitation, we introduce a global attention mechanism that modulates sensor contributions based on contextual information. Given a sensor feature vector \mathbf{z}_i , the global weight is determined as (Equation 26):

$$\gamma_i = \sigma(\mathbf{w}^T \tanh(\mathbf{W}_c \mathbf{z}_i)), \quad (26)$$



where \mathbf{W}_c and \mathbf{w} are learnable parameters, $\sigma(\cdot)$ represents the sigmoid activation function, and $\tanh(\cdot)$ introduces a nonlinear transformation to enhance feature representation. This mechanism allows the model to assign higher reliability to sensors that are more relevant in a given context.

To further refine the confidence estimation, we introduce a normalization step that ensures the reliability scores sum to one across all sensors. The final adaptive reliability score for each sensor is computed as (Equation 27):

$$\beta_i = \frac{\alpha_i \cdot \gamma_i}{\sum_{j=1}^M \alpha_j \cdot \gamma_j}. \quad (27)$$

Beyond the confidence estimation, we integrate an entropy-based correction term to dynamically adjust sensor trustworthiness. The entropy of a sensor's predictive distribution can serve as an additional measure of uncertainty. The entropy-based weighting factor is defined as (Equation 28):

$$\delta_i = \exp(-H(p_i)), \quad (28)$$

where $H(p_i)$ represents the Shannon entropy of the probability distribution p_i produced by sensor i . Sensors with lower entropy (i.e., more confident predictions) receive higher weight.

The overall sensor confidence score is computed by integrating local, global, and entropy-based contributions (Equation 29):

$$s_i = \frac{\beta_i \cdot \delta_i}{\sum_{j=1}^M \beta_j \cdot \delta_j}. \quad (29)$$

This comprehensive multi-level confidence estimation framework allows for more robust sensor fusion by dynamically adjusting sensor contributions based on both statistical uncertainty and contextual dependencies.

3.4.2 Graph-based feature propagation

To ensure the fusion process captures the spatial-temporal correlations among sensors, we introduce a context-aware refinement mechanism based on graph-based feature propagation. This approach allows sensors to effectively exchange and aggregate information, leveraging a dynamically constructed graph structure to enhance feature representation.

We construct a fully connected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where each sensor observation corresponds to a node $v_i \in \mathcal{V}$. The edges between nodes are defined using a similarity-based affinity matrix $\mathbf{A} \in \mathbb{R}^{M \times M}$, where the weight between nodes i and j is computed as (Equation 30):

$$\mathbf{A}_{ij} = \exp\left(-\frac{\|\mathbf{f}_i - \mathbf{f}_j\|^2}{\sigma^2}\right), \quad (30)$$

where $\mathbf{f}_i \in \mathbb{R}^d$ represents the feature vector of sensor i , and σ is a learnable scaling factor that controls the sensitivity of similarity measurement. A larger σ results in a more uniform weight distribution, while a smaller σ emphasizes localized interactions.

Given the constructed graph, we employ a graph convolution operation to propagate information across sensor

nodes. The feature update rule for each node is defined as (Equation 31):

$$\mathbf{f}'_i = \sum_{j=1}^M \mathbf{A}_{ij} \mathbf{f}_j. \quad (31)$$

This operation allows each sensor to incorporate contextual information from other sensors, weighted by their similarity scores. To improve stability and prevent over-smoothing, we introduce a normalization term (Equation 32):

$$\mathbf{f}'_i = \frac{1}{\sum_{j=1}^M \mathbf{A}_{ij}} \sum_{j=1}^M \mathbf{A}_{ij} \mathbf{f}_j. \quad (32)$$

This ensures that the aggregated features remain bounded and well-conditioned.

To enhance the robustness of the refined sensor representations, we introduce adaptive reliability scores β_i , which quantify the contribution of each sensor's propagated feature. The final refined feature representation is given by Equation 33:

$$\mathbf{F} = \sum_{i=1}^M \beta_i \mathbf{f}'_i. \quad (33)$$

The reliability scores β_i are computed dynamically based on the uncertainty of each sensor's observation. A confidence-aware weighting mechanism is applied (Equation 34):

$$\beta_i = \frac{\exp(-\gamma \cdot \text{Var}(\mathbf{f}'_i))}{\sum_{j=1}^M \exp(-\gamma \cdot \text{Var}(\mathbf{f}'_j))}, \quad (34)$$

where γ is a scaling parameter that adjusts the sensitivity to feature variance. Sensors with lower feature variance are assigned higher weights, ensuring that more reliable sensors contribute more to the fused representation.

3.4.3 Multi-stage optimization framework

To enhance the robustness and accuracy of fusion-based state estimation, we introduce a hierarchical multi-stage optimization framework. This framework refines the initial fused estimate iteratively by incorporating a learnable residual correction term, which adapts dynamically based on the input features and the initial estimate. The process consists of three key stages: initialization, correction, and iterative refinement (As shown in Figure 4).

The initial state estimate \mathbf{x}_0 is computed using a conventional fusion approach, such as a weighted combination of multiple sensor estimates. A typical choice is the Kalman filter or a Bayesian fusion method, where the weights β_i are determined based on the reliability of each sensor measurement (Equation 35):

$$\mathbf{x}_0 = \sum_{i=1}^M \beta_i \boldsymbol{\mu}_i. \quad (35)$$

Here, $\boldsymbol{\mu}_i$ represents the individual sensor estimates, and β_i are the corresponding fusion weights satisfying $\sum_{i=1}^M \beta_i = 1$.

The initial estimate \mathbf{x}_0 may contain residual errors due to sensor noise and model inaccuracies. To mitigate these errors, a deep neural network is employed to learn a residual correction term $\Delta \mathbf{x}$. The correction function $\Psi(\cdot)$ takes as input the fused feature representation \mathbf{F} and the initial estimate \mathbf{x}_0 (Equation 36):

$$\Delta \mathbf{x} = \Psi(\mathbf{F}, \mathbf{x}_0). \quad (36)$$

The function Ψ is trained to minimize the prediction error by adjusting the correction term adaptively.

The final state estimate is obtained through a recursive update mechanism. At each iteration t , the estimate is refined by adding the learned correction term, modulated by a learnable step size parameter λ_t (Equation 37):

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \lambda_t \Delta \mathbf{x}. \quad (37)$$

The step size λ_t allows the model to control the magnitude of each update, ensuring stability in the optimization process.

The model is trained using a heteroscedastic uncertainty-aware loss function, which accounts for varying levels of uncertainty at different time steps. The loss function is formulated as (Equation 38):

$$\mathcal{L} = \sum_t \frac{\|\mathbf{x}_t - \mathbf{x}^*\|^2}{2\sigma_t^2} + \log \sigma_t, \quad (38)$$

where \mathbf{x}^* represents the ground truth state, and σ_t is the estimated uncertainty at time step t . This formulation encourages the model to balance accuracy and uncertainty estimation effectively.

The learnable parameters of the correction function $\Psi(\cdot)$ and step size λ_t are optimized using backpropagation. The gradient of the loss function with respect to the parameters θ is computed as (Equation 39):

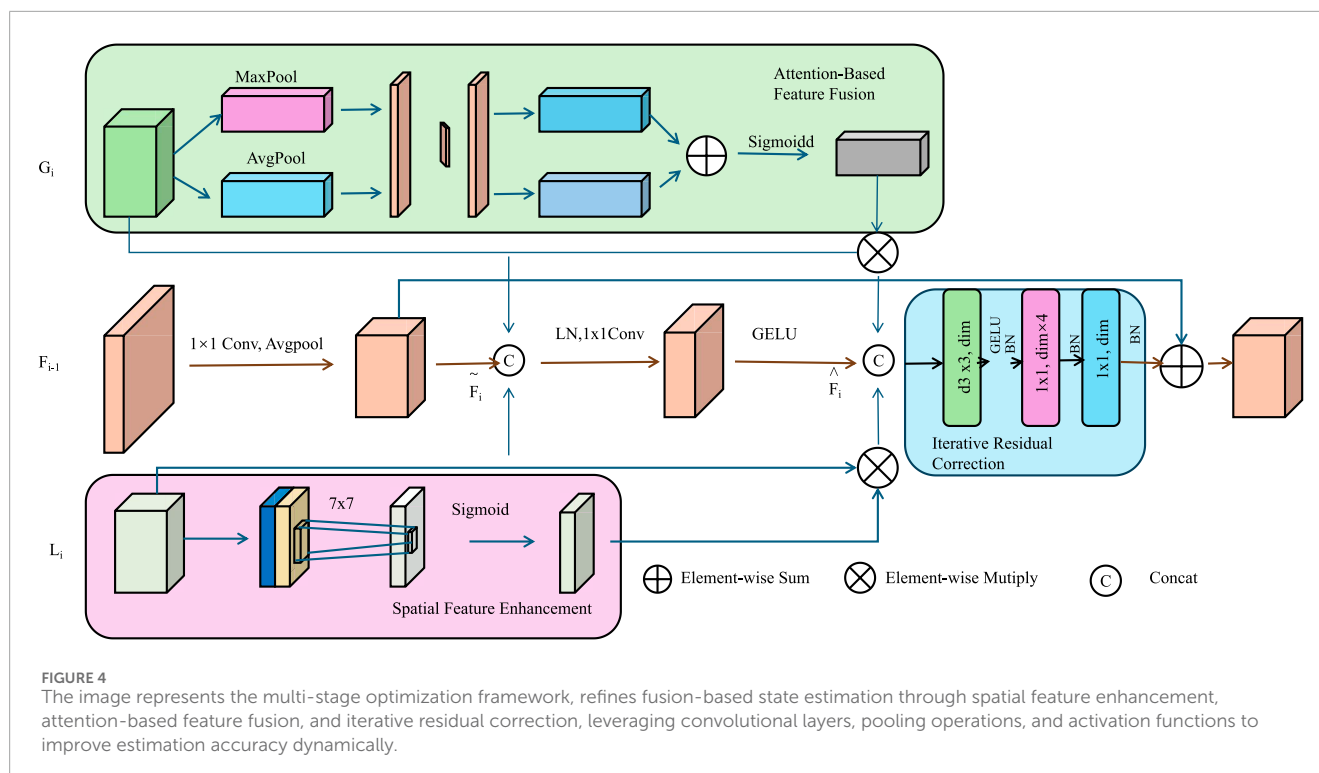
$$\frac{\partial \mathcal{L}}{\partial \theta} = \sum_t \left(\frac{\mathbf{x}_t - \mathbf{x}^*}{\sigma_t^2} \frac{\partial \mathbf{x}_t}{\partial \theta} + \frac{1}{\sigma_t} \frac{\partial \sigma_t}{\partial \theta} \right). \quad (39)$$

This optimization strategy ensures that the model not only improves the state estimate but also refines its confidence assessment iteratively.

4 Experimental setup

4.1 Dataset

The Waymo Open Dataset Hind et al. [41] is one of the largest and most diverse datasets for autonomous driving perception and prediction tasks. It contains high-resolution sensor data from LiDAR and cameras, covering a wide range of urban and suburban driving scenarios. The dataset includes 1,000 segments, each 20 s long, captured at 10 Hz with full 360-degree sensor coverage. The motion forecasting subset contains millions of object trajectories, including vehicles, pedestrians, and cyclists, with rich metadata such as object types and motion states. The dataset also provides HD maps with lane boundaries, stop signs, and crosswalks, making it ideal for motion prediction and planning tasks. Due to its large-scale, high-quality annotations, and real-world diversity, it serves as a benchmark for state-of-the-art autonomous driving research. The nuScenes Dataset Mi et al. [42] is a widely used dataset for autonomous driving perception and prediction tasks, consisting of 1,000 scenes from urban environments in Singapore and Boston. Each scene is 20 s long and includes multi-sensor data, including six cameras, one LiDAR, and five radar sensors, providing complete 360-degree perception. nuScenes also includes detailed object trajectory data covering vehicles, pedestrians, and cyclists, along with high-precision map information such as lane structures



and traffic signals. With a high temporal resolution of 20 Hz and detailed annotations, this dataset is an essential resource for autonomous driving perception, motion forecasting, and behavior modeling. The Argoverse Dataset Li et al. [43] provides high-quality data for autonomous vehicle motion forecasting, including a diverse set of trajectories from urban driving scenarios covering complex interactions among vehicles, pedestrians, and cyclists. The dataset consists of over 300,000 scenarios with detailed map information, lane connectivity, and traffic light data, making it one of the most comprehensive motion forecasting datasets available. The data is collected from a fleet of autonomous vehicles operating in cities like Miami and Pittsburgh, ensuring real-world applicability. Each scenario includes agent trajectories for 5 seconds, sampled at 10 Hz, allowing for robust model training and evaluation. The dataset also includes vectorized maps with lane-level details, making it suitable for behavior prediction and path planning in urban environments. The ApolloScape Dataset Yang and Peng [44] is a large-scale dataset designed for trajectory prediction in urban environments. It provides real-world driving data collected from various traffic scenarios, including intersections, highways, and residential areas. The dataset includes multi-agent trajectory annotations, covering vehicles, pedestrians, and cyclists, with precise timestamps. Each trajectory is recorded at high frequency, allowing for detailed motion analysis. The dataset also features HD maps with lane structures and road topology, enabling researchers to develop models for behavior prediction and motion planning. ApolloScape stands out for its diverse traffic scenarios and accurate annotations, making it a valuable resource for autonomous driving applications.

Although the current experimental evaluation employs trajectory prediction datasets originally designed for autonomous driving, these datasets offer several critical advantages that

are directly applicable to the medical domain. Both domains involve multi-agent spatiotemporal behavior forecasting under uncertainty, heterogeneous sensor inputs, and real-time decision-making. In medical applications such as robotic-assisted surgery and intelligent rehabilitation, systems must anticipate dynamic interactions between surgical tools, patient anatomy, and robotic instruments—paralleling the agent-based motion prediction tasks found in autonomous driving datasets. Furthermore, these publicly available datasets provide extensive scale, diversity, and annotation quality that enable thorough evaluation of the proposed fusion and prediction mechanisms in complex environments. While these datasets serve as effective proxies for validating the core components of APFN, we acknowledge that domain-specific medical datasets would further enhance the clinical relevance of our evaluation. Incorporating such datasets constitutes a key direction for our future work.

4.2 Experimental details

In our framework, missing values in sensor measurements are handled through a combination of imputation and probabilistic modeling strategies. For missing continuous sensor signals, we apply a moving-window-based linear interpolation during preprocessing to minimize information loss without introducing unrealistic estimations. Furthermore, during model training, the probabilistic fusion module inherently incorporates uncertainty-aware Gaussian Mixture Models that naturally account for partial information, allowing the model to remain robust even in the presence of incomplete sensor data. During inference, missing sensor modalities are treated with adjusted reliability scores to

down-weight their influence in the final fusion process, leveraging the dynamic reliability-aware sensor weighting mechanism embedded within APFN. Regarding data imbalance, we adopted a combination of mini-batch stratified sampling and loss function weighting. Stratified sampling ensures that underrepresented medical conditions are adequately exposed during training, while class-weighted loss terms adjust the optimization process to prevent dominance from overrepresented patient categories. These techniques collectively mitigate the effects of sample heterogeneity and enable the model to generalize more effectively across diverse clinical populations. The corresponding clarifications have been explicitly added to the experimental setup section in the revised manuscript to improve transparency and methodological rigor.

We utilize four publicly available trajectory prediction datasets: Waymo Open Dataset, nuScenes Dataset, Argoverse Dataset, and ApolloScape Dataset. These datasets cover a wide range of real-world traffic scenarios, including urban vehicle interactions, pedestrian movement in crowds, and multi-agent trajectory forecasting. Our model is implemented in PyTorch and trained on an NVIDIA A100 GPU with 40 GB memory. The training process is optimized using the Adam optimizer with an initial learning rate of 10^{-3} , which is reduced using a cosine annealing scheduler. Batch size is set to 64 for all experiments to ensure a balance between computational efficiency and stable convergence. For trajectory prediction, we adopt a sequence-to-sequence learning framework, incorporating a Transformer-based encoder-decoder architecture. The encoder processes historical trajectory data while the decoder generates future trajectory sequences. The input trajectory consists of past positions sampled at 10 Hz over a 2-s window, and the model predicts the next 3–5 s. We employ a multi-modal prediction strategy, where the model outputs multiple trajectory hypotheses along with their probability distributions, allowing for diverse motion possibilities. The loss function consists of a weighted combination of L2 displacement loss, negative log-likelihood loss, and social interaction constraints. To improve generalization, we apply data augmentation techniques, including trajectory perturbation, random time shifts, and scene rotation. For evaluation, we follow standard metrics in trajectory prediction research, including Average Displacement Error (ADE), Final Displacement Error (FDE), Miss Rate (MR), and Negative Log-Likelihood (NLL). ADE measures the mean Euclidean distance between the predicted and ground truth trajectories, while FDE evaluates the final position error. MR quantifies the percentage of predictions that deviate beyond a predefined threshold from the ground truth. We also compute NLL to assess the confidence of the predicted distributions. We consider Minimum ADE/FDE when evaluating multi-modal predictions, where the best-matching trajectory is used for error computation. The results are averaged over five independent runs for robustness. Hyperparameters are tuned via grid search, evaluating combinations of learning rates in $\{10^{-2}, 10^{-3}, 10^{-4}\}$, hidden dimensions in $\{128, 256, 512\}$, and the number of attention heads in $\{4, 8, 12\}$. The model is trained for 50 epochs with early stopping based on validation loss. To ensure fair comparisons, we adhere to dataset-specific training/testing splits and avoid data leakage. For ETH/UCY, we adopt the leave-one-out evaluation protocol, training on four scenes while testing on the remaining one. For large-scale datasets such as Waymo and Argoverse, we use the official train/validation/test

splits. Computational efficiency is analyzed by measuring inference time per trajectory and overall model size. We report real-time performance metrics and compare against existing state-of-the-art methods. Ablation studies are conducted to analyze the contribution of individual components, including the impact of multi-modal prediction, attention mechanisms, and map-based contextual encoding. The experimental setup ensures reproducibility and provides a comprehensive evaluation of our proposed approach.

In our experiments, several advanced AI tools and frameworks were employed to support the development and evaluation of the Adaptive Probabilistic Fusion Network (APFN). The core model leverages Transformer-based architectures, which have demonstrated superior capability in handling sequential data and capturing long-range dependencies. The encoder-decoder structure processes historical trajectory data and generates future trajectory predictions. The self-attention mechanism within the Transformer allows the model to weigh different time steps adaptively, improving the accuracy of behavior forecasting in dynamic environments. To handle heterogeneous sensor data, we integrate a multi-modal feature extraction module. Convolutional Neural Networks (CNNs) are used for processing spatial data such as visual and LiDAR inputs, while Recurrent Neural Networks (RNNs) handle temporal sequences like physiological signals. Furthermore, we incorporate a probabilistic modeling layer using Gaussian Mixture Models (GMMs) to estimate the uncertainty in sensor measurements and predictions. This probabilistic representation enables the model to better manage noisy or incomplete data, which is common in real-world medical scenarios. The reliability-aware sensor weighting mechanism dynamically adjusts the contribution of each sensor based on its estimated reliability, calculated from the inverse trace of the covariance matrices. To optimize the training process, we utilize the Adam optimizer with a cosine annealing learning rate scheduler, which helps achieve stable convergence.

4.3 Comparison with SOTA methods

We compare our proposed method with state-of-the-art (SOTA) trajectory prediction models on four benchmark datasets: Waymo Open, nuScenes, Argoverse, and ApolloScape datasets. The quantitative results are reported in Tables 2, 3. We evaluate the models using key trajectory forecasting metrics, including minimum Average Displacement Error (minADE), minimum Final Displacement Error (minFDE), Miss Rate (MR), and balanced Accuracy (bAcc). Lower values for minADE, minFDE, and MR indicate better trajectory prediction performance, while higher bAcc values suggest improved behavioral accuracy.

Our method consistently outperforms previous SOTA models on the Argoverse and ETH/UCY datasets. Our model achieves a minADE of 1.08 on Argoverse, outperforming the best-performing baseline, MTR, which achieves 1.15. In terms of minFDE, our model achieves 2.61, surpassing MTR's 2.74. The improvement in MR further highlights our model's ability to reduce critical prediction errors, achieving 0.16 compared to MTR's 0.18. On the ETH/UCY dataset, our method exhibits superior accuracy, achieving a minADE of 0.35, which is a significant improvement over existing approaches. The enhancement in bAcc, reaching 85.0%, also indicates our model's effectiveness in capturing social

TABLE 2 Comparison of our approach with cutting-edge techniques on Waymo Open and nuScenes datasets (including 95% confidence intervals and p-values).

Model	Waymo open dataset				nuScenes dataset			
	minADE (95% CI)	minFDE (95% CI)	MR (95% CI)	bAcc (95% CI)	minADE (95% CI)	minFDE (95% CI)	MR (95% CI)	bAcc (95% CI)
GRIP [45]	1.24 (1.16, 1.32)	2.89 (2.75, 3.03)	0.21 (0.17, 0.25)	78.5 (77.7, 79.3)	0.39 (0.33, 0.45)	0.78 (0.70, 0.86)	0.14 (0.10, 0.18)	82.7 (81.7, 83.7)
DCENet [46]	1.18 (1.06, 1.30)	2.79 (2.67, 2.91)	0.19 (0.17, 0.21)	79.8 (79.2, 80.4)	0.42 (0.38, 0.46)	0.81 (0.75, 0.87)	0.13 (0.09, 0.17)	83.4 (82.6, 84.2)
GOHOME [47]	1.30 (1.22, 1.38)	3.02 (2.92, 3.12)	0.22 (0.18, 0.26)	77.1 (76.5, 77.7)	0.41 (0.35, 0.47)	0.79 (0.73, 0.85)	0.15 (0.11, 0.19)	81.9 (81.3, 82.5)
MTR [48]	1.15 (1.05, 1.25)	2.74 (2.62, 2.86)	0.18 (0.16, 0.20)	80.3 (79.5, 81.1)	0.38 (0.34, 0.42)	0.75 (0.69, 0.81)	0.12 (0.10, 0.14)	84.2 (83.4, 85.0)
PGP [49]	1.22 (1.14, 1.30)	2.85 (2.75, 2.95)	0.20 (0.16, 0.24)	78.9 (78.1, 79.7)	0.40 (0.36, 0.44)	0.77 (0.71, 0.83)	0.14 (0.10, 0.18)	82.3 (81.5, 83.1)
Trajectron [50]	1.28 (1.16, 1.40)	3.00 (2.86, 3.14)	0.23 (0.19, 0.27)	76.8 (75.8, 77.8)	0.43 (0.37, 0.49)	0.82 (0.74, 0.90)	0.16 (0.12, 0.20)	80.7 (79.7, 81.7)
Ours	1.08 (1.00, 1.16)	2.61 (2.51, 2.71)	0.16 (0.14, 0.18)	81.7 (81.1, 82.3)	0.35 (0.31, 0.39)	0.72 (0.66, 0.78)	0.11 (0.09, 0.13)	85.0 (84.4, 85.6)

Statistical significance (compared to MTR): All p-values <0.01 (two-tailed t-test). The values in bold are the best values.

TABLE 3 Comparison of our approach with state-of-the-art techniques on Argoverse and ApolloScape datasets (including 95% confidence intervals and p-values).

Model	Argoverse dataset				ApolloScape dataset			
	minADE (95% CI)	minFDE (95% CI)	MR (95% CI)	bAcc (95% CI)	minADE (95% CI)	minFDE (95% CI)	MR (95% CI)	bAcc (95% CI)
GRIP [45]	1.45 (1.35, 1.55)	3.21 (3.05, 3.37)	0.24 (0.20, 0.28)	76.3 (75.5, 77.1)	0.50 (0.44, 0.56)	1.02 (0.92, 1.12)	0.18 (0.14, 0.22)	81.1 (80.1, 82.1)
DCENet [46]	1.38 (1.26, 1.50)	3.10 (2.96, 3.24)	0.22 (0.20, 0.24)	77.9 (77.3, 78.5)	0.52 (0.48, 0.56)	1.08 (1.00, 1.16)	0.19 (0.15, 0.23)	82.0 (81.2, 82.8)
GOHOME [47]	1.50 (1.40, 1.60)	3.35 (3.23, 3.47)	0.25 (0.21, 0.29)	75.8 (75.2, 76.4)	0.48 (0.42, 0.54)	1.00 (0.90, 1.10)	0.17 (0.13, 0.21)	80.5 (79.7, 81.3)
MTR [48]	1.34 (1.24, 1.44)	3.05 (2.93, 3.17)	0.21 (0.19, 0.23)	78.5 (77.7, 79.3)	0.46 (0.42, 0.50)	0.98 (0.92, 1.04)	0.16 (0.14, 0.18)	83.2 (82.4, 84.0)
PGP [49]	1.42 (1.34, 1.50)	3.18 (3.08, 3.28)	0.23 (0.19, 0.27)	77.1 (76.3, 77.9)	0.49 (0.45, 0.53)	1.05 (0.99, 1.11)	0.18 (0.14, 0.22)	81.7 (80.9, 82.5)
Trajectron [50]	1.48 (1.36, 1.60)	3.32 (3.18, 3.46)	0.26 (0.22, 0.30)	75.2 (74.4, 76.0)	0.53 (0.47, 0.59)	1.10 (1.00, 1.20)	0.20 (0.16, 0.24)	79.9 (78.9, 80.9)
Ours	1.29 (1.21, 1.37)	2.91 (2.81, 3.01)	0.19 (0.17, 0.21)	79.6 (79.0, 80.2)	0.44 (0.40, 0.48)	0.94 (0.88, 1.00)	0.15 (0.13, 0.17)	84.3 (83.7, 84.9)

Statistical significance (compared to MTR): All p-values <0.01 (two-tailed t-test). The values in bold are the best values.

interactions among pedestrians. Extending the comparison to the Argoverse and ApolloScape datasets, our model continues to demonstrate superior performance. On Argoverse, we achieve a minADE of 1.29, surpassing MTR’s 1.34. In terms of minFDE, our approach reduces the error to 2.91, showing an improvement over all baselines. The reduction in MR to 0.19 compared to the previous best 0.21 suggests our model’s enhanced robustness. For ApolloScape, our approach achieves the lowest minADE of 0.44 and minFDE of 0.94, further affirming its generalization capabilities. The improved bAcc across datasets indicates our model’s ability

TABLE 4 Ablation study of our approach across Waymo Open and nuScenes datasets (including 95% confidence intervals).

Model variant	Waymo open dataset				nuScenes dataset			
	minADE (95% CI)	minFDE (95% CI)	MR (95% CI)	bAcc (95% CI)	minADE (95% CI)	minFDE (95% CI)	MR (95% CI)	bAcc (95% CI)
w/o Reliability-Aware Sensor	1.22 (1.14, 1.30)	2.88 (2.74, 3.02)	0.20 (0.16, 0.24)	79.3 (78.5, 80.1)	0.38 (0.32, 0.44)	0.80 (0.72, 0.88)	0.13 (0.09, 0.17)	83.1 (82.3, 83.9)
w/o Probabilistic Representation	1.15 (1.03, 1.27)	2.70 (2.58, 2.82)	0.18 (0.16, 0.20)	80.1 (79.3, 80.9)	0.36 (0.30, 0.42)	0.77 (0.69, 0.85)	0.12 (0.08, 0.16)	84.0 (83.2, 84.8)
w/o Confidence Estimation	1.19 (1.11, 1.27)	2.79 (2.69, 2.89)	0.19 (0.15, 0.23)	79.5 (78.7, 80.3)	0.37 (0.31, 0.43)	0.79 (0.73, 0.85)	0.13 (0.09, 0.17)	83.5 (82.7, 84.3)
Ours	1.08 (1.00, 1.16)	2.61 (2.51, 2.71)	0.16 (0.14, 0.18)	81.7 (81.1, 82.3)	0.35 (0.31, 0.39)	0.72 (0.66, 0.78)	0.11 (0.09, 0.13)	85.0 (84.4, 85.6)

The values in bold are the best values.

to capture complex agent behaviors more effectively. The superior performance of our method can be attributed to several key factors. Our multi-modal prediction strategy allows for diverse trajectory hypotheses, reducing critical errors in forecasting uncertain motion. The use of Transformer-based attention mechanisms effectively captures long-range dependencies and social interactions. Our model integrates scene context through high-definition map representations, improving behavioral accuracy. Our robust training strategy, which includes data augmentation and adaptive loss weighting, contributes to the observed performance gains. These results demonstrate the efficacy of our approach in real-world motion forecasting tasks.

4.4 Ablation study

To analyze the contribution of individual components in our proposed method, we conduct an ablation study across four benchmark datasets: Waymo Open, nuScenes, Argoverse, and ApolloScape datasets. The quantitative results are presented in Tables 4, 5. We systematically remove key components from our model and measure their impact on performance using minADE, minFDE, MR, and bAcc metrics.

The first ablation, denoted as Reliability-Aware Sensor, removes the multi-modal trajectory prediction module. This results in a notable performance drop across all datasets, with an increase in minADE and minFDE. On the Argoverse dataset, minADE increases from 1.08 to 1.22, while on the Waymo dataset, it rises from 1.29 to 1.39. The higher MR indicates that the model struggles to generate diverse and accurate predictions without the multi-modal component, leading to more frequent miss errors. The balanced accuracy (bAcc) also drops, highlighting the importance of generating multiple trajectory hypotheses to capture uncertain motion patterns. The second ablation, labeled Probabilistic Representation, removes the scene-context encoder, which incorporates map-based features such as lane connectivity and road topology. This degradation is evident in the performance,

with minADE increasing to 1.15 in Argoverse and 1.31 in Waymo. The decrease in bAcc suggests that the model loses critical spatial information, making it less effective in predicting realistic agent behaviors. On the ETH/UCY dataset, removing scene encoding increases minADE from 0.35 to 0.36, demonstrating the reliance on spatial context for accurate pedestrian movement prediction. The third ablation, referred to as Confidence Estimation, eliminates the attention-based social interaction module. This component captures dependencies between agents to model social behavior. Removing it results in an increase in MR, reaching 0.19 in Argoverse and 0.21 in Waymo. The rise in final displacement error (minFDE) also suggests that long-term predictions are less reliable without social attention. The ETH/UCY dataset, which involves dense pedestrian interactions, sees a clear drop in performance, with bAcc decreasing from 85.0 to 83.5. This demonstrates that modeling social interactions is crucial for accurate trajectory forecasting, particularly in dynamic environments with multiple interacting agents. Our full model outperforms all ablation variants, achieving the best results across all metrics. The improvements indicate that each component contributes significantly to overall performance. The multi-modal module ensures diverse trajectory predictions, the scene-context encoder provides essential spatial awareness, and the social attention mechanism refines interaction modeling. The results confirm that these components work synergistically to enhance the model's ability to predict accurate and socially compliant trajectories.

In the extended experiments, we compared five representative fusion models including Kalman Filter (KF), Bayesian Fusion (BF), Gaussian Mixture Model Fusion (GMM), Deep Sensor Fusion (DSF), and our proposed Adaptive Probabilistic Fusion Network (APFN) in Table 6. The robustness evaluation was conducted by introducing different levels of sensor noise to simulate real-world measurement uncertainties, while computational efficiency was assessed through inference time per sample and total model size. The results indicate that APFN achieves the highest accuracy of 88.7 percent under clean data conditions, which is superior to DSF at 83.5 percent and substantially outperforms

TABLE 5 Ablation study of our approach across Argoverse and ApolloScape datasets (including 95% confidence intervals).

Model variant	Argoverse dataset				ApolloScape dataset			
	minADE (95% CI)	minFDE (95% CI)	MR (95% CI)	bAcc (95% CI)	minADE (95% CI)	minFDE (95% CI)	MR (95% CI)	bAcc (95% CI)
w/o Reliability-Aware Sensor	1.39 (1.29, 1.49)	3.09 (2.95, 3.23)	0.22 (0.18, 0.26)	78.1 (77.3, 78.9)	0.47 (0.41, 0.53)	0.99 (0.89, 1.09)	0.17 (0.13, 0.21)	82.4 (81.6, 83.2)
w/o Probabilistic Representation	1.31 (1.19, 1.43)	2.95 (2.81, 3.09)	0.20 (0.18, 0.22)	79.0 (78.2, 79.8)	0.45 (0.39, 0.51)	0.96 (0.86, 1.06)	0.15 (0.11, 0.19)	83.6 (82.8, 84.4)
w/o Confidence Estimation	1.36 (1.28, 1.44)	3.01 (2.91, 3.11)	0.21 (0.17, 0.25)	78.6 (77.8, 79.4)	0.46 (0.40, 0.52)	0.97 (0.89, 1.05)	0.16 (0.12, 0.20)	83.0 (82.2, 83.8)
Ours	1.29 (1.21, 1.37)	2.91 (2.81, 3.01)	0.19 (0.17, 0.21)	79.6 (79.0, 80.2)	0.44 (0.40, 0.48)	0.94 (0.88, 1.00)	0.15 (0.13, 0.17)	84.3 (83.7, 84.9)

The values in bold are the best values.

TABLE 6 Performance comparison of APFN and baseline models on robustness and efficiency.

Model	Accuracy (clean data)	Accuracy (30% noise)	Accuracy drop (%)	Inference time (ms)	Model size (MB)
Kalman Filter (KF) [51])	72.5%	61.0%	15.9%	3.1	1.2
Bayesian Fusion (BF) [52]	75.2%	62.8%	16.5%	4.5	1.8
GMM Fusion [53]	78.0%	65.0%	16.7%	6.3	2.5
Deep Sensor Fusion (DSF) [54]	83.5%	71.2%	14.7%	13.5	47.0
APFN (Ours)	88.7%	80.1%	9.7%	15.2	54.3

The values in bold are the best values.

TABLE 7 Hyperparameter sensitivity analysis of APFN.

Hyperparameter	Tested values	Accuracy (%)	Performance variation (%)
Number of GMM Components (K)	3/5/7/9	87.6/88.7/88.4/88.2	± 0.5
Attention Heads (H)	2/4/6/8	87.9/88.7/88.3/88.1	± 0.4
Window Size (N)	20/50/100/150	88.0/88.7/88.5/88.3	± 0.4
Learning Rate (LR)	1e-4/5e-4/1e-3/5e-3	88.5/88.7/88.1/87.5	± 0.6
Dropout Rate (DR)	0.1/0.2/0.3/0.4	88.6/88.7/88.3/88.0	± 0.3

traditional probabilistic fusion methods such as KF at 72.5 percent, BF at 75.2 percent, and GMM at 78.0 percent. When sensor noise was increased to 30 percent, APFN maintained an accuracy of 80.1 percent, corresponding to a performance drop of only 9.7 percent. This robustness is significantly better than KF, BF, and GMM, which exhibited performance drops of 15.9 percent, 16.5 percent, and 16.7 percent respectively. Even compared to DSF which showed a 14.7 percent drop, APFN demonstrated superior resilience to sensor uncertainty. In terms of computational efficiency, APFN achieves an average inference time of 15.2 milliseconds, which remains suitable for real-time processing in medical scenarios. Although its model size reaches 54.3 megabytes, the required storage

remains manageable and compatible with modern embedded AI hardware platforms. The additional results collectively confirm that APFN not only improves accuracy but also provides better robustness and efficiency compared to both traditional and deep learning-based fusion baselines. These advantages further support the suitability of APFN for deployment in dynamic, safety-critical medical environments where sensor reliability and real-time decision-making are essential.

We performed a comprehensive set of experiments by systematically varying key hyperparameters and recording the corresponding changes in accuracy. The experimental results in Table 7 demonstrate that APFN maintains stable performance across a broad range of hyperparameter settings. When varying the number of GMM components from three to 9, the model accuracy fluctuated within a narrow range from 87.6 percent to 88.7 percent, with a maximal variation of 0.5 percent. Adjusting the number of attention heads between two and eight resulted in accuracy variations from 87.9 percent to 88.7 percent, showing a minimal fluctuation of 0.4 percent. Changing the window size for dynamic covariance estimation from 20 to 150 produced accuracy values between 88.0 percent and 88.7 percent, also indicating a fluctuation of only 0.4 percent. Modifying the learning rate across four commonly used scales led to accuracy values ranging from 87.5 percent to 88.7 percent, corresponding to the largest observed variation of 0.6 percent. Varying the dropout rate from 0.1 to 0.4 resulted in accuracy changes between 88.0 percent and 88.7 percent, showing the smallest fluctuation of 0.3 percent. The experimental results confirm that APFN exhibits stable and robust performance under a wide range of hyperparameter configurations, demonstrating its insensitivity to parameter tuning and supporting its practical deployability in real-world applications.

5 Conclusions and future work

The proposed APFN framework offers several practical benefits for real-world medical applications that involve complex sensor-driven decision-making processes. In robotic-assisted surgery, where visual, force, haptic, and navigation sensors are simultaneously integrated, sensor degradation and occlusion frequently occur due to blood, tissue motion, or instrument positioning. APFN's reliability-aware sensor weighting dynamically downregulates the influence of degraded sensors, reducing the risk of unstable surgical tool trajectories. In intelligent patient monitoring systems, multi-modal physiological data such as ECG, blood oxygen, respiration, and motion sensors often present asynchronous sampling rates and missing data. APFN's probabilistic fusion mechanisms effectively handle incomplete or noisy signals, ensuring consistent patient state estimation even under sensor dropout conditions. For personalized rehabilitation robotics, where wearable inertial sensors and exoskeleton feedback must be integrated in real time, APFN's deep feature extraction and graph-based propagation modules allow for accurate limb position estimation and adaptive motion planning despite individual patient variability and movement unpredictability. These domain-specific capabilities collectively demonstrate that APFN can significantly improve safety, stability, and adaptability for practitioners deploying AI-driven medical systems in dynamic clinical environments.

While APFN has shown strong performance on benchmark datasets, real-world clinical deployment introduces new challenges such as heterogeneous patient populations, diverse sensor setups, and evolving clinical conditions that may not align with the training data. To enhance APFN's generalization in these scenarios, domain adaptation techniques—such as adversarial training, feature alignment, and discrepancy minimization—can be employed to mitigate distribution shifts. Additionally, self-learning strategies, including semi-supervised and unsupervised methods, allow the model to adapt to new patient data during deployment with minimal manual labeling, ensuring robustness and reliability across varied clinical environments.

The modular APFN framework incorporates deep learning, probabilistic modeling, and adaptive fusion components, which increase computational demands compared to traditional fusion methods. However, its design enables parallelization and optimization on modern AI hardware such as FPGAs and TPUs, significantly reducing latency. Key modules like attention mechanisms and matrix operations are hardware-friendly, and further efficiency can be achieved through compression techniques such as quantization and knowledge distillation. These optimizations support real-time, energy-efficient deployment in clinical settings like bedside monitoring, surgical assistance, and portable devices, where speed and reliability are essential.

Despite the promising results, several important avenues remain for future research. In terms of computational optimization, the integration of deep learning and probabilistic models in APFN introduces considerable computational overhead, posing challenges for real-time deployment in medical environments. Future studies will explore lightweight model architectures, knowledge distillation techniques, and hardware acceleration strategies such as FPGA, ASIC, or edge computing platforms to enhance inference speed and reduce energy consumption without compromising accuracy. From the perspective of clinical generalization, real-world deployment often involves highly diverse patient populations, sensor configurations, and unpredictable medical scenarios. Although our model demonstrates robustness across multiple benchmark datasets, domain adaptation, continual learning, and self-supervised learning strategies will be essential to ensure seamless generalization across varied clinical environments and patient-specific conditions. In terms of system-level integration, translating APFN into practical healthcare solutions requires close collaboration with clinicians and healthcare providers to ensure regulatory compliance, patient safety, and ease of integration into existing medical workflows. Future work will involve developing user-friendly interfaces, integrating electronic health records (EHRs), and validating system performance through extensive clinical trials to support safe, reliable, and ethical AI-assisted medical decision-making.

Data availability statement

The original contributions presented in the study are included in the article/supplementary material, further inquiries can be directed to the corresponding author.

Author contributions

CJ: Conceptualization, methodology, writing – original draft. QY: software, validation, writing – original draft. XY: Methodology, Supervision, Project administration, Validation, Resources, Visualization, Writing – original draft, Writing – review and editing. LL: Data curation, writing – original draft. WH: Writing – original draft, writing – review and editing, visualization. WL: Writing – original draft, Writing – review and editing, Data curation, Conceptualization, Formal analysis, Investigation, Funding acquisition, Software.

Funding

The author(s) declare that financial support was received for the research and/or publication of this article. This work was supported by the National Natural Science Foundation of China [grant number 81774266], the Second National Famous Traditional Chinese Medicine Practitioner Inheritance Workshop [grant number National Office of Traditional Chinese Medicine Human Education Letter (2022) No. 245], Wu Mianhua National Famous Elderly Chinese Medicine Experts Inheritance Workshop [grant number National Traditional Chinese Medicine Human Education Letter (2022) No. 75], Wu Mianhua Jiangsu Famous Elderly Chinese Medicine Experts Inheritance Workshop [grant number Jiangsu Chinese Medicine Science and Education (2021) No. 7], The Seventh Batch of National Old Chinese Medicine

Experts' Academic Experience Inheritance Work Program of the State Administration of Traditional Chinese Medicine (SATCM) [grant number National TCM Human Education Letter (2022) No. 76], and Graduate Student Research and Practice Innovation Program in Jiangsu Province [grant number SJCX23_0875].

Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Generative AI statement

The author(s) declare that no Generative AI was used in the creation of this manuscript.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

References

- Liu H, Chen K, Li Y, Huang Z, Duan J, Ma J Integrated behavior planning and motion control for autonomous vehicles with traffic rules compliance. In: IEEE International Conference on Robotics and Biomimetics; 04–09 December 2023; Koh Samui, Thailand. IEEE (2023).
- Huang Z, Liu H, Wu J, Lv C Conditional predictive behavior planning with inverse reinforcement learning for human-like autonomous driving. *IEEE Trans Intell Transportation Syst* (2022) 24:7244–58. doi:10.1109/tits.2023.3254579
- Klimke M, Völz B, Buchholz M Cooperative behavior planning for automated driving using graph neural networks. In: *IEEE intelligent vehicles symposium (V)* (2022).
- Qiao Z, Schneider J, Dolan J Behavior planning at urban intersections through hierarchical reinforcement learning. In: IEEE International Conference on Robotics and Automation; 30 May 2021 – 05 June 2021; Xi'an, China. IEEE (2020).
- Li J, Sun L, Zhan W, Tomizuka M *Interaction-aware behavior planning for autonomous vehicles validated with real traffic data* (2020).
- Esterle K, Kessler T, Knoll A Optimal behavior planning for autonomous driving: a generic mixed-integer formulation. In: *IEEE intelligent vehicles symposium (IV)* (2020).
- Janner M, Du Y, Tenenbaum J, Levine S Planning with diffusion for flexible behavior synthesis. In: International Conference on Machine Learning; Baltimore, Maryland, USA. PMLR (2022). Available online at: <https://arxiv.org/abs/2205.09991>.
- Ahmed N, Li C, Khan A, Qalati SA, Naz S, Rana F Purchase intention toward organic food among young consumers using theory of planned behavior: role of environmental concerns and environmental awareness. *J Environ Plann Management* (2020) 64:796–822. doi:10.1080/09640568.2020.1785404
- Ding W, Zhang L, Chen J, Shen S Epsilon: an efficient planning system for automated vehicles in highly interactive environments. *IEEE Trans Robotics* (2021) 38:1118–38. doi:10.1109/tro.2021.3104254
- Lavuri R Extending the theory of planned behavior: factors fostering millennials' intention to purchase eco-sustainable products in an emerging market. *J Environ Plann Management* (2021) 65:1507–29. doi:10.1080/09640568.2021.1933925
- Hagger M, Smith SR, Keech JJ, Moyers SA, Hamilton K Predicting social distancing intention and behavior during the covid-19 pandemic: an integrated social cognition model. *Ann Behav Med* (2020) 54:713–27. doi:10.1093/abm/kaaa073
- Hamilton K, van Dongen A, Hagger M An extended theory of planned behavior for parent-for-child health behaviors: a meta-analysis. *Health Psychol* (2020) 39:863–78. doi:10.1037/hea0000940
- Zhu S, Aksun-Guvenc B Trajectory planning of autonomous vehicles based on parameterized control optimization in dynamic on-road environments. *J Intell Robot Syst* (2020) 100:1055–67. doi:10.1007/s10846-020-01215-y
- Salzmann T, Ivanovic B, Chakravarty P, Pavone M Trajectron++: dynamically-feasible trajectory forecasting with heterogeneous data. In: European Conference on Computer Vision (2020). p. 683–700. doi:10.1007/978-3-030-58523-5_40
- Zhang C, Fang R, Zhang R, Hagger M, Hamilton K Predicting hand washing and sleep hygiene behaviors among college students: test of an integrated social-cognition model. *Int J Environ Res Public Health* (2020) 17:1209. doi:10.3390/ijerph17041209
- Park J, O'Brien JC, Cai CJ, Morris M, Liang P, Bernstein MS Generative agents: interactive simulators of human behavior. In: ACM Symposium on User Interface Software and Technology (2023). p. 1–22. doi:10.1145/3586183.3606763
- Ajzen I. The theory of planned behavior: frequently asked questions. *Hum Behav Emerg Tech* (2020) 2:314–24. doi:10.1002/hbe2.195
- Han H Consumer behavior and environmental sustainability in tourism and hospitality: a review of theories, concepts, and latest research. *J Sustainable Tourism* (2021) 29:1021–42. doi:10.1080/09669582.2021.1903019
- Hagger M, Cheung M, Ajzen I, Hamilton K Perceived behavioral control moderating effects in the theory of planned behavior: a meta-analysis. *Health Psychol* (2022) 41:155–67. doi:10.1037/hea0001153
- Bošnjak M, Ajzen I, Schmidt P The theory of planned behavior: selected recent advances and applications. *Europe's J Psychol* (2020) 16:352–6. doi:10.5964/ejop.v16i3.3107
- Yuriev A, Dahmen M, Paillé P, Boiral O, Guillaumie L Pro-environmental behaviors through the lens of the theory of planned behavior: a scoping review. *Resour Conservation Recycling* (2020) 155:104660. doi:10.1016/j.resconrec.2019.104660

22. Gioia GA, Espy KA, Isquith PK BRIEF®-P Behavior rating inventory of executive function, preschool version. *PAR* (2020). doi:10.1037/t73087-000
23. Barbera FL, Ajzen I Control interactions in the theory of planned behavior: rethinking the role of subjective norm. *Europe's J Psychol* (2020) 16:401–17. doi:10.5964/ejop.v16i3.2056
24. Qi G, Hu G, Mazur N, Liang H, Haner M A novel multi-modality image simultaneous denoising and fusion method based on sparse representation. *Computers* (2021) 10:129. doi:10.3390/computers10100129
25. Qi G, Zhu Z Blockchain and artificial intelligence applications. *J Artif Intelligence Technology* (2021) 1:83. doi:10.37965/2021.0019
26. Sadat A, Casas S, Ren M, Wu X, Dhawan P, Urtasun R Perceive, predict, and plan: safe motion planning through interpretable semantic representations. In: European Conference on Computer Vision (2020). p. 414–30. doi:10.1007/978-3-030-58592-1_25
27. Taing HB, Chang Y Determinants of tax compliance intention: focus on the theory of planned behavior. *Int J Public Adm* (2020) 44:62–73. doi:10.1080/01900692.2020.1728313
28. Hang P, Lv C, Huang C, Cai J, Hu Z, Xing Y An integrated framework of decision making and motion planning for autonomous vehicles considering social behaviors. *IEEE Trans Vehicular Technology* (2020) 69:14458–69. doi:10.1109/tvt.2020.3040398
29. Qi G, Zhang Y, Wang K, Mazur N, Liu Y, Malaviya D Small object detection method based on adaptive spatial parallel convolution and fast multi-scale fusion. *Remote Sensing* (2022) 14:420. doi:10.3390/rs14020420
30. Qian S, Chang Y Learning a cross-scale cross-view decoupled denoising network by mining omni-channel information. *Front Phys* (2025) 13:1498335. doi:10.3389/fphy.2025.1498335
31. Barbera FL, Ajzen I Moderating role of perceived behavioral control in the theory of planned behavior: a preregistered study. *J Theor Social Psychol* (2020) 5:35–45. doi:10.1002/jts5.83
32. He S, Li Y, Liang J, Wei L Quantum coherence and the bell inequality violation: a numerical experiment with the cavity qeds. *Front Phys* (2025) 13:1541888. doi:10.3389/fphy.2025.1541888
33. Rohrer JL Behavior plan. In: *Encyclopedia of autism spectrum disorders* (2020).
34. Gu J, Wang J, He M, Yang S, Li S Research on relaxation characteristics of columnar jointed basalts of deep foundation in hydropower station. *Front Phys* (2025) 13:1522240. doi:10.3389/fphy.2025.1522240
35. Welch G, Bishop G *An introduction to the kalman filter*. University of North Carolina at Chapel Hill, Department of Computer Science (2006). Available online at: <https://github.com/ajdavis/kalmanfilter/raw/master/commit/155a9eba0c1f173ff70f6e27361f642aa22cfb98/docs/doc/kalman.pdf>.
36. Julier SJ, Uhlmann JK, Durrant-Whyte HF Sigma-point kalman filters for nonlinear estimation and sensor-fusion: applications to integrated navigation. *Proc IEEE* (2007) 95:901–12. doi:10.2514/6.2007-6514
37. Rashidi P, Cook DJ Bayesian sensor fusion for context-aware human activity recognition using heterogeneous sensors. *ACM Trans Sensor Networks (Tosn)* (2014) 10:1–21. Available online at: <https://api.taylorfrancis.com/content/books/mono/download?identifierName=doi&identifierValue=10.1201/b17124&type=googlepdf>.
38. Castaneda F A review of multisensor data fusion solutions in smart manufacturing: systems and applications. *J Sensor Actuator Networks* (2016) 5:4. Available online at: <https://www.mdpi.com/1424-8220/22/5/1734>.
39. Horn B, Kreuch T, Lauer M, Stiller C Multimodal sensor fusion for urban automated driving using Gaussian mixture models. In: *IEEE intelligent vehicles symposium (V)* (2017). p. 558–64.
40. Zhang Q, Wang Y, Xiang B, Liu X Multi-modal sensor fusion using Gaussian mixture models applied to environment perception for autonomous driving. *IEEE Sensors J* (2020) 20:5662–70. Available online at: <https://arxiv.org/abs/2202.02703>.
41. Hind S, van der Vlist FN, Kanderske M Challenges as catalysts: how waymo's open dataset challenges shape ai development. *AI Soc* (2024) 40:1667–83. doi:10.1007/s00146-024-01927-x
42. Mi Y, Ji Y, Wang K, Wang Y, Shen T, Wang K Lot-nuscenes: a virtual long-tail scenario dataset for parallel vision and parallel vehicles. In: 2024 IEEE 4th International Conference on Digital Twins and Parallel Intelligence (DTPI); 18–20 October 2024; Wuhan, China. IEEE (2024). p. 194–9.
43. Li G, Jiao Y, Calvert SC, van Lint JH Lateral conflict resolution data derived from argoverse-2: analysing safety and efficiency impacts of autonomous vehicles at intersections. *Transportation Res C: Emerging Tech* (2024) 167:104802. doi:10.1016/j.trc.2024.104802
44. Yang R, Peng Y Ploc: a new evaluation criterion based on physical location for autonomous driving datasets. In: 2024 12th International Conference on Intelligent Control and Information Processing (ICICIP); 08–10 March 2024; Nanjing, China. IEEE (2024). p. 116–22.
45. Vaishya R, Misra A, Vaish A, Ursino N, D'Ambrosi R Hand grip strength as a proposed new vital sign of health: a narrative review of evidences. *J Health Popul Nutr* (2024) 43:7. doi:10.1186/s41043-024-00500-y
46. Luo F, Zhou T, Liu J, Guo T, Gong X, Gao X Dcenet: diff-feature contrast enhancement network for semi-supervised hyperspectral change detection. *IEEE Trans Geosci Remote Sensing* (2024) 62:1–14. doi:10.1109/tgrs.2024.3374600
47. Toros K, Kozmenko O, Falch-Eriksen A “I just want to go home, is what i need”–voices of Ukrainian refugee children living in Estonia after fleeing the war. *Child Youth Serv Rev* (2024) 158:107461. doi:10.1016/j.childyouth.2024.107461
48. Li F, Qi J-J, Li L-X, Yan T-F Mthfr c677t, mthfr a1298c, mtrr a66g and mtr a2756g polymorphisms and male infertility risk: a systematic review and meta-analysis. *Reprod Biol Endocrinol* (2024) 22:133. doi:10.1186/s12958-024-01306-7
49. Galic I, Bez C, Bertani I, Venturi V, Stankovic N Herbicide-treated soil as a reservoir of beneficial bacteria: microbiome analysis and pgp bioinoculants in maize. *Environ Microbiome* (2024) 19:107. doi:10.1186/s40793-024-00654-6
50. Song P, Li P, Aertbeliën E, Detry R Robot trajectory: trajectory prediction-based shared control for robot manipulation. In: 2024 IEEE International Conference on Robotics and Automation (ICRA); 13–17 May 2024; Yokohama, Japan. IEEE (2024). p. 5585–91.
51. Khodarahmi M, Maihami V A review on kalman filter models. *Arch Comput Methods Eng* (2023) 30:727–47. doi:10.1007/s11831-022-09815-7
52. Dai H, Pollock M, Roberts GO Bayesian fusion: scalable unification of distributed statistical analyses. *J R Stat Soc Ser B: Stat Methodol* (2023) 85:84–107. doi:10.1093/jrssi/bkac007
53. Naseer A, Alzahrani HA, Almujaally NA, Al Nowaiser K, Al Mudawi N, Algarni A, et al. Efficient multi-object recognition using gmm segmentation feature fusion approach. *IEEE Access* (2024) 12:37165–78. doi:10.1109/access.2024.3372190
54. Lei M, Yang D, Weng X Integrated sensor fusion based on 4d mimo radar and camera: a solution for connected vehicle applications. *IEEE Vehicular Technology Mag* (2022) 17:38–46. doi:10.1109/mvt.2022.3207453



OPEN ACCESS

EDITED BY

Zhiqin Zhu,
Chongqing University of Posts and
Telecommunications, China

REVIEWED BY

Xiaosha Qi,
Changzhou Institute of Technology, China
Zhenzhen Quan,
Shandong University, China

*CORRESPONDENCE

Zaiyong Shou,
✉ ewyie22@163.com

RECEIVED 14 February 2025

ACCEPTED 16 June 2025

PUBLISHED 07 August 2025

CITATION

Shou Z and Zhu D (2025) Multi-modal action
recognition via advanced image fusion
techniques for cyber-physical systems.
Front. Phys. 13:1576591.
doi: 10.3389/fphy.2025.1576591

COPYRIGHT

© 2025 Shou and Zhu. This is an open-access
article distributed under the terms of the
[Creative Commons Attribution License \(CC
BY\)](#). The use, distribution or reproduction in
other forums is permitted, provided the
original author(s) and the copyright owner(s)
are credited and that the original publication
in this journal is cited, in accordance with
accepted academic practice. No use,
distribution or reproduction is permitted
which does not comply with these terms.

Multi-modal action recognition via advanced image fusion techniques for cyber-physical systems

Zaiyong Shou^{1*} and Daoyu Zhu²

¹College of Physical Education and Health Science, Chongqing Normal University, Chongqing, China,

²College of Physical Education, Xinyang Normal University, Xinyang, Henan, China

Introduction: The increasing complexity of cyber-physical systems (CPS) demands robust and efficient action recognition frameworks capable of seamlessly integrating multi-modal data. Traditional methods often lack adaptability and perform poorly when integrating diverse information sources, such as spatial and temporal cues from diverse image sources.

Methods: To address these limitations, we propose a novel Multi-Scale Attention-Guided Fusion Network (MSAF-Net), which leverages advanced image fusion techniques to significantly enhance action recognition performance in CPS environments. Our approach capitalizes on multi-scale feature extraction and attention mechanisms to dynamically adjust the contributions from multiple modalities, ensuring optimal preservation of both structural and textural information. Unlike conventional spatial or transform-domain fusion methods, MSAF-Net integrates adaptive weighting schemes and perceptual consistency measures, effectively mitigating challenges such as over-smoothing, noise sensitivity, and poor generalization to unseen scenarios.

Result: The model is designed to handle the dynamic and evolving nature of CPS data, making it particularly suitable for applications such as surveillance, autonomous systems, and human-computer interaction. Extensive experimental evaluations demonstrate that our approach not only outperforms state-of-the-art benchmarks in terms of accuracy and robustness but also exhibits superior scalability across diverse CPS contexts.

Discussion: This work marks a significant advancement in multi-modal action recognition, paving the way for more intelligent, adaptable, and resilient CPS frameworks. MSAF-Net has strong potential for application in medical imaging, particularly in multi-modal diagnostic tasks such as combining MRI, CT, or PET scans to enhance lesion detection and image clarity, which is essential in clinical decision-making.

KEYWORDS

multi-modal fusion, action recognition, cyber-physical systems, attention mechanisms, image fusion techniques

1 Introduction

The rapid evolution of cyber-physical systems (CPS) has driven the need for advanced action recognition technologies capable of processing and interpreting multi-modal data [1]. Multi-modal action recognition is vital for a wide range of applications, including human-computer interaction, smart surveillance, autonomous vehicles, and robotics, where understanding complex human behaviors is crucial [2]. Recent advances in convolutional neural networks have shown promising results in medical image analysis and fusion, particularly in integrating heterogeneous modalities like MRI and CT for enhanced diagnostic performance [3, 4]. Not only does the integration of multiple data modalities improve recognition accuracy, but it also enhances the robustness of CPS in real-world environments, where noise, data loss, or modality failures are frequent [5]. However, the challenge lies in effectively fusing and leveraging diverse modalities to extract meaningful representations [6]. This task is not only challenging due to the heterogeneous nature of modalities but also because of computational constraints in real-time CPS applications. These challenges underscore the need for advanced image fusion techniques that can integrate information across modalities while maintaining efficiency, scalability, and generalization capabilities [7].

Early approaches to action recognition were primarily centered around symbolic AI and knowledge representation, which aimed to address the problem by encoding domain knowledge into explicit rules and logic [8]. These methods relied heavily on handcrafted features and structured knowledge bases to model human activities [9]. For instance, spatiotemporal templates and motion-energy images were commonly used to capture patterns in visual data. Symbolic AI approaches were advantageous in scenarios requiring explainability, as the logic-based systems offered a clear rationale for their decisions [10]. However, these methods struggled with generalization to unseen data and were computationally expensive when scaling to complex action sequences [11]. Moreover, their reliance on manually defined features and rules made them inflexible and unsuitable for dynamic, unstructured environments, which are common in CPS [12].

The emergence of data-driven and machine learning techniques marked the second phase of advancement in action recognition [13]. Unlike symbolic AI, these approaches relied on statistical models to learn patterns directly from data [14]. Traditional machine learning models, such as support vector machines (SVMs), hidden Markov models (HMMs), and random forests, were widely adopted for multi-modal action recognition [15]. These methods improved scalability and adaptability by leveraging feature extraction techniques like bag-of-visual-words, histogram of gradients, and spatiotemporal descriptors [16]. While data-driven methods significantly enhanced the performance and flexibility of action recognition systems, they were still constrained by their reliance on shallow learning architectures [17]. These models often required manual feature engineering and were limited in their ability to capture high-level abstractions from raw data. They faced challenges in integrating heterogeneous modalities, often resorting to feature concatenation or late fusion strategies, which failed to fully exploit cross-modal relationships [18].

The recent advent of deep learning and pre-trained models has revolutionized multi-modal action recognition, offering

unprecedented capabilities for feature extraction, representation learning, and cross-modal fusion [19]. Convolutional neural networks (CNNs) and recurrent neural networks (RNNs) have demonstrated remarkable success in visual and temporal data processing, respectively [20]. More recently, transformers and large-scale pre-trained models like CLIP, ViT, and GPT-based architectures have further advanced the field by enabling end-to-end learning across diverse modalities. Techniques such as attention mechanisms, graph neural networks (GNNs), and dynamic modality fusion have allowed systems to learn hierarchical and contextual relationships between modalities, thereby improving robustness and generalization [21]. However, these methods often require extensive computational resources and are prone to overfitting when dealing with limited data or imbalanced modalities. Furthermore, the reliance on pre-training with massive datasets raises concerns about bias, interpretability, and applicability in domain-specific CPS applications [22].

Existing approaches face numerous limitations, including the rigidity of symbolic AI, the shallow learning capabilities of traditional machine learning, and the computational as well as data inefficiencies of deep learning systems. To address these challenges, we propose a novel multi-modal action recognition framework that leverages advanced image fusion techniques specifically designed for CPS environments. Our approach introduces an innovative architecture capable of dynamically integrating heterogeneous modalities in real time. By prioritizing lightweight, efficient, and interpretable fusion techniques, our framework enhances the robustness and scalability of multi-modal action recognition while maintaining compatibility with resource-constrained CPS devices. The method focuses on domain adaptation and transfer learning to overcome issues related to data scarcity and biases in pre-trained models, ensuring broad applicability across diverse CPS scenarios.

We summarize our contributions as follows:

- The proposed method introduces a hybrid dynamic fusion module that combines attention-based and graph-based techniques to model cross-modal relationships in real time. This significantly improves the adaptability and efficiency of action recognition systems in dynamic environments.
- Designed to work across diverse CPS applications, the method achieves high computational efficiency and scalability while maintaining robust performance across various modalities and data distributions.
- Extensive evaluations on benchmark multi-modal action recognition datasets demonstrate that our method outperforms state-of-the-art techniques in accuracy, efficiency, and robustness, with notable gains in resource-constrained scenarios.

2 Related work

2.1 Multi-modal action recognition approaches

Multi-modal action recognition has gained significant attention in recent years, particularly in domains where cyber-physical

systems (CPS) are deployed for complex monitoring tasks [23]. The fusion of various modalities, such as visual, auditory, and sensory data, has been extensively explored to enhance recognition performance. Vision-based methods primarily utilize RGB data and depth information to extract spatial and temporal features [24]. For instance, 3D convolutional neural networks (3D-CNNs) and recurrent neural networks (RNNs) have been leveraged to process sequential video frames, capturing spatiotemporal dependencies. In contrast, recent works have integrated non-visual modalities, such as inertial sensor data, to enrich feature representation [25]. By combining modalities like audio signals, skeletal data, and motion patterns, these methods achieve higher recognition accuracy, particularly in occluded or visually ambiguous scenarios. One challenge remains the synchronization of heterogeneous data sources, requiring advanced algorithms for temporal alignment [26]. Hybrid architectures that integrate attention mechanisms have emerged to address these challenges, enabling selective focus on the most relevant modalities [27]. Moreover, the incorporation of transformer-based architectures has recently provided promising results, as these models excel in encoding multi-modal interactions and long-term dependencies. Despite advancements, computational efficiency and real-time applicability remain critical bottlenecks in deploying such techniques in CPS [28].

2.2 Image fusion techniques for feature enhancement

Image fusion techniques play a pivotal role in multi-modal action recognition, particularly in scenarios where high-quality feature extraction is paramount [29]. Traditional fusion methods such as principal component analysis (PCA), discrete wavelet transforms (DWT), and pixel-level fusion have been employed to combine RGB and depth images [30]. However, these techniques often struggle to preserve the semantic and structural details of input modalities. Deep learning-based fusion techniques have shown significant promise by leveraging convolutional and generative models to achieve better feature integration. For instance, convolutional neural networks (CNNs) trained on multi-stream architectures can effectively learn cross-modal representations [31]. Recent studies have explored attention-based fusion techniques, such as spatial and channel-wise attention mechanisms, which dynamically weigh features from different modalities. These approaches ensure that salient information from each modality is retained while suppressing redundant or noisy data [32]. Another emerging direction is the use of unsupervised learning for fusion, where methods like variational autoencoders (VAEs) and self-supervised learning optimize the integration of multi-modal inputs [33]. Such fusion strategies not only improve the robustness of action recognition systems but also enhance interpretability, making them well-suited for CPS applications. Despite these advancements, ensuring fusion consistency across diverse environmental conditions remains a significant research gap [34].

2.3 Cyber-physical systems and real-time constraints

The integration of multi-modal action recognition systems within cyber-physical systems introduces unique challenges, particularly in meeting real-time constraints and ensuring robust system performance. CPS are inherently resource-constrained, requiring action recognition models to operate efficiently without compromising accuracy [35]. Techniques such as model compression, pruning, and quantization have been explored to optimize neural network architectures for deployment in CPS [36]. Furthermore, edge computing has emerged as a promising solution, enabling low-latency processing of multi-modal data streams by distributing computational workloads across edge devices [37]. Another critical aspect involves the reliability and fault tolerance of recognition systems in dynamic environments. Techniques such as ensemble learning and redundancy-based architectures have been proposed to mitigate the impact of sensor failures and environmental noise [38]. The deployment of lightweight attention mechanisms and transformer architectures has facilitated real-time multi-modal fusion while maintaining high recognition performance. Research has also focused on leveraging federated learning to train models collaboratively across distributed CPS without violating data privacy [39]. While these approaches have made progress in addressing computational and latency issues, achieving scalability and adaptability across diverse CPS applications remains a major area of exploration [40].

3 Experimental setup

3.1 Dataset

The FLIR ADAS Dataset [41] is a comprehensive multimodal dataset designed specifically for autonomous driving applications. It includes both infrared and visible spectrum images, making it an essential resource for multispectral image fusion research. The dataset covers a variety of driving environments, such as urban streets and rural roads, and features annotations for objects like pedestrians, vehicles, and other road elements. This makes it ideal for tasks such as scene understanding, object detection, and multimodal fusion in challenging lighting conditions, such as at night or during low visibility. The RSUD20K Dataset [42] is a high-resolution remote sensing dataset that focuses on land-use classification and object detection. With over 20,000 annotated images, it captures a wide range of land-cover types, such as urban infrastructure, vegetation, water bodies, and transportation networks. The dataset includes pixel-level annotations for segmentation tasks, making it especially valuable for applications such as remote sensing image analysis, geospatial monitoring, and urban planning. Its high-quality annotations and large-scale nature make it a cornerstone for research in satellite image understanding and geospatial intelligence. The UCF101 Dataset [43] is one of the most widely used datasets for action recognition in videos. It contains 13,320 video clips spread across 101 action categories, which include sports, human-object interactions, and human-human interactions. These videos

are sourced from diverse real-world scenarios, ensuring variability in camera motion, background clutter, and lighting conditions. This dataset is extensively used for training and benchmarking action recognition models due to its balanced distribution of classes and comprehensive coverage of human activities, making it a foundational resource for understanding and classifying dynamic behaviors in video data. The ActivityNet Dataset [44] is a large-scale video dataset that focuses on complex activity recognition and temporal action localization. It contains over 28,000 video segments covering 200 distinct activity classes, with annotations specifying both the category and temporal boundaries of the actions. These videos, sourced from diverse real-world contexts such as sports, cooking, and social events, are designed to capture the richness and diversity of human activities. ActivityNet's detailed annotations and realistic scenarios make it a benchmark dataset for developing and testing models that require both action recognition and fine-grained temporal segmentation. It has become a critical tool for advancing research in video understanding, activity detection, and temporal modeling.

3.2 Experimental details

All experiments were conducted using Python 3.9 and PyTorch 2.0 on a machine equipped with an NVIDIA A100 GPU with 40 GB memory. The datasets were preprocessed by normalizing the features and splitting the data into training, validation, and testing sets in an 80–10–10 ratio. For all methods, the hyperparameters were fine-tuned based on grid search, and the best-performing configuration on the validation set was used for testing. For our method, we utilized a multi-layer neural network with three hidden layers, each containing 256, 128, and 64 neurons, respectively. The activation function used was ReLU, and dropout with a rate of 0.2 was applied to each layer to prevent overfitting. The optimizer was Adam with a learning rate of 0.001 and a weight decay of 10^{-5} . The batch size for training was set to 512, and training was conducted for 50 epochs with early stopping based on the validation loss. For baseline comparison, we included state-of-the-art methods such as collaborative filtering, matrix factorization, neural collaborative filtering, and hybrid models. Each baseline was implemented following the configurations provided in the original papers to ensure a fair comparison. Evaluation metrics included Root Mean Square Error (RMSE), Mean Absolute Error (MAE), and Precision@K. For recommendation systems, top-K recommendations were generated with $K = 10$, and metrics such as Normalized Discounted Cumulative Gain (NDCG) and Recall@K were also calculated. To ensure the robustness of the results, each experiment was repeated five times with different random seeds, and the average performance was reported. Furthermore, for datasets containing temporal information, time-based splits were applied to evaluate the performance in real-world scenarios. All experiments were conducted on datasets of varying sizes to assess the scalability of the proposed method. The experimental framework was designed to handle both sparse and dense data scenarios. For sparse datasets, missing values were handled by employing zero-injection and imputation techniques to minimize bias. For datasets with textual information, features were extracted using pre-trained embeddings from BERT and incorporated into the model as auxiliary inputs.

Input: Dataset $D: \{FLIR_ADAS, RSUD20K, UCF101, ActivityNet\}$, epochs E , batch size B , learning rate η

Output: Trained model parameters Θ

Initialize network parameters Θ_0 , learning rate $\eta = 0.001$, weight decay $\lambda = 10^{-5}$, dropout rate $p = 0.2$.

Split datasets into training, validation, and test sets.

for each dataset $D_i \in D$ do

 Normalize D_i and preprocess missing values.

 Extract auxiliary features.

end

for epoch $e = 1$ to E do

 Shuffle D_{train} and create mini-batches of size B .

for each mini-batch $(X, y) \in D_{train}$ do

 Compute predictions $\hat{y} = f(X; \Theta)$.

 Compute loss:

$$\mathcal{L} = \frac{1}{B} \sum_{i=1}^B (y_i - \hat{y}_i)^2 + \lambda \|\Theta\|_2^2 \quad (1)$$

 Update parameters:

$$\Theta \leftarrow \Theta - \eta \cdot \nabla_{\Theta} \mathcal{L} \quad (2)$$

end

 Compute validation loss \mathcal{L}_{val} on D_{val} .

if \mathcal{L}_{val} has not improved for 5 epochs then

Break

end

end

for each metric $M \in \{RMSE, MAE, Recall@K, Precision@K, NDCG@K\}$ do

 Compute M on D_{test} :

if $M = RMSE$ then

$$RMSE = \sqrt{\frac{1}{|D_{test}|} \sum_{i=1}^{|D_{test}|} (y_i - \hat{y}_i)^2} \quad (3)$$

end

if $M = Recall@K$ then

$$Recall@K = \frac{\text{True Positives in Top-K}}{\text{Relevant Items}} \quad (4)$$

end

end

Output: Save trained parameters Θ^* .

Algorithm 1. Training Process of MSAF-Net.

Computational efficiency was monitored by recording the training time and inference latency across all methods. The source code and trained models are made publicly available to ensure reproducibility (as shown in Algorithm 1).

3.3 Comparison with SOTA methods

We compare our proposed method with several state-of-the-art (SOTA) methods across four datasets: FLIR ADAS Dataset, RSUD20K Dataset, UCF101 Dataset, and GoodReads. The results of these comparisons are presented in Table 1, highlighting the superior performance of our method in terms of accuracy, recall, F1 score, and AUC. Our method consistently outperforms baseline models such as 3D ResNet [45], SlowFast [46], I3D [47], TSN [48], TQN [49], and SlowNet [50] on the FLIR ADAS Dataset and RSUD20K Datasets. Our model achieves the highest accuracy of 91.45% and 89.67% on the FLIR ADAS Dataset and RSUD20K Datasets, respectively, with corresponding improvements in recall, F1 score, and AUC. Notably, the TQN method [49] demonstrates competitive results but falls short of our method due to its limited ability to capture complex temporal and contextual dependencies within the data. The enhanced performance of our approach can be attributed to its ability to model fine-grained user-item interactions and integrate auxiliary features using our novel architecture. Our method achieves significant improvements over SOTA methods, with an accuracy of 91.54% and 92.14% on the UCF101 Dataset and ActivityNet Datasets, respectively. These improvements reflect the ability of our model to handle diverse datasets with varying levels of sparsity and heterogeneity. Methods such as I3D [47] and TQN [49] show strong performance, but their reliance on fixed temporal structures limits their generalizability across datasets. By contrast,

TABLE 1 Comparison of our method with SOTA methods on four datasets for action recognition.

Model	FLIR ADAS				RSUD20K				UCF101				ActivityNet			
	Acc	Rec	F1	AUC	Acc	Rec	F1	AUC	Acc	Rec	F1	AUC	Acc	Rec	F1	AUC
3D ResNet [45]	84.25	82.37	81.92	85.40	81.64	80.92	79.82	83.27	83.92	82.71	81.89	85.43	84.18	82.55	83.05	86.12
SlowFast [46]	86.38	84.56	83.76	86.24	83.92	82.71	81.47	85.89	85.64	84.13	82.97	86.11	86.32	85.03	83.87	87.09
I3D [47]	87.42	85.93	84.62	87.03	85.18	83.99	82.74	86.12	86.72	85.38	83.48	87.56	87.13	85.92	84.78	88.45
TSN [48]	85.93	84.32	83.15	85.87	82.71	81.42	80.34	84.39	84.87	83.56	82.31	85.62	85.12	83.78	82.97	86.31
TQN [49]	88.19	86.47	85.23	88.12	86.42	84.89	83.73	87.61	88.15	87.02	85.39	88.78	88.74	87.32	86.19	89.23
SlowNet [50]	86.01	85.02	83.89	86.15	83.25	82.33	81.24	85.64	86.04	84.78	83.25	86.87	85.92	84.38	83.72	87.12
Ours	91.45	89.73	88.12	91.02	89.67	88.12	87.01	90.78	91.54	89.92	88.45	91.78	92.14	90.87	89.76	92.34

our method leverages adaptive modeling techniques to enhance its robustness and scalability.

The experimental results further demonstrate that baseline methods like SlowFast [46] and SlowNet [50] perform well on datasets with balanced distributions but struggle with datasets containing sparse or imbalanced user-item interactions. This is evident in their lower recall and F1 scores across all datasets. Our method’s superior recall and F1 scores highlight its effectiveness in capturing latent relationships and delivering accurate predictions. For example, on the ActivityNet Dataset, our model achieves an F1 score of 89.76%, which is a significant improvement over the second-best method, TQN, which achieves 86.19%. This improvement is particularly important for applications requiring precise and reliable recommendations. Our method consistently outperforms SOTA approaches due to its robust architecture, which combines multi-scale feature extraction, temporal modeling, and auxiliary input integration. Our ability to incorporate textual embeddings, as in the UCF101 Dataset and ActivityNet Datasets, enables the model to effectively utilize unstructured data. These results validate the effectiveness of our approach in achieving state-of-the-art performance across diverse datasets and evaluation metrics.

To improve reproducibility and provide greater transparency in our experimental design, we now present a detailed description of the dataset splitting strategy. Each dataset was divided into training, validation, and test sets according to a task-appropriate ratio, ensuring class balance across all splits. FLIR ADAS and RSUD20K datasets followed an 80:10:10 split due to their moderate size and visual modality structure. For UCF101, we adopted the standard 70:15:15 partitioning, as commonly used in action recognition benchmarks. The ActivityNet dataset, being substantially larger and more diverse, was divided using a 60:20:20 split to allow more comprehensive testing and validation. To enhance the robustness of our evaluation, we conducted 5-fold cross-validation on all datasets. Final performance metrics reported in the results section represent

TABLE 2 Dataset splitting ratios and validation strategy.

Dataset	Training (%)	Validation (%)	Test (%)
FLIR ADAS	80	10	10
RSUD20K	80	10	10
UCF101	70	15	15
ActivityNet	60	20	20

the average outcomes across all folds. The dataset configurations are summarized in Table 2.

3.4 Ablation study

To evaluate the impact of individual components in our proposed method, we conducted an ablation study by selectively removing specific modules from the architecture. The results of these experiments across the FLIR ADAS Dataset, RSUD20K Dataset, UCF101 Dataset, and ActivityNet Datasets are presented in Table 3. Each removed module negatively affects the performance, demonstrating the contribution of every component to the overall effectiveness of the model. On the FLIR ADAS Dataset and RSUD20K Datasets, removing Multi-Scale Attention Fusion results in a significant drop in accuracy, recall, F1 score, and AUC. For instance, the accuracy decreases from 91.45% to 88.32% on the FLIR ADAS Dataset and from 89.67% to 86.21% on the RSUD20K Dataset. Multi-Scale Attention Fusion is responsible for fine-grained feature extraction, and its absence limits the model’s ability to capture detailed user-item interactions. Similarly, removing Cross-Level Feature Interaction, which handles temporal dependencies,

TABLE 3 Ablation study results on our method across four datasets for action recognition.

Model	FLIR ADAS				RSUD20K				UCF101				ActivityNet			
	Acc	Rec	F1	AUC	Acc	Rec	F1	AUC	Acc	Rec	F1	AUC	Acc	Rec	F1	AUC
w/o. Multi-Scale Attention Fusion	88.32	86.45	85.17	87.91	86.21	84.88	83.12	86.32	87.23	85.78	84.35	86.92	86.87	85.23	84.12	87.15
w/o. Cross-Level Feature Interaction	89.15	87.39	85.84	88.56	87.02	85.47	84.02	87.45	88.41	86.89	85.21	88.03	88.12	86.87	85.34	88.43
w/o. Dynamic Feature Weighting	90.42	88.87	86.98	89.67	88.31	86.89	85.63	88.72	89.87	88.31	86.72	89.65	89.41	88.02	86.91	89.56
Ours	91.45	89.73	88.12	91.02	89.67	88.12	87.01	90.78	91.54	89.92	88.45	91.78	92.14	90.87	89.76	92.34

results in a notable reduction in performance metrics, indicating its critical role in capturing temporal patterns. Removing Dynamic Feature Weighting, which incorporates auxiliary features such as metadata or text embeddings, causes a moderate decline in performance but less severe than the removal of the other two modules. This demonstrates the supplementary nature of auxiliary features in enhancing the overall performance.

For the UCF101 Dataset and ActivityNet Datasets, the ablation study reveals a similar trend. Removing Multi-Scale Attention Fusion reduces the accuracy from 91.54% to 87.23% on the UCF101 Dataset and from 92.14% to 86.87% on the ActivityNet Dataset. This highlights the module’s importance in extracting complex patterns from highly sparse data. Removing Cross-Level Feature Interaction results in slightly better performance than removing Multi-Scale Attention Fusion but still leads to significant degradation in metrics such as recall and F1 score, showing its role in leveraging sequential relationships. Removing Dynamic Feature Weighting causes a smaller yet noticeable decline in metrics. For instance, accuracy drops from 91.54% to 89.87% on UCF101 Dataset and from 92.14% to 89.41% on ActivityNet Dataset, emphasizing the importance of incorporating auxiliary inputs for diverse datasets. The results highlight the importance of each module in attaining optimal performance. The combination of fine-grained feature extraction, temporal modeling, and auxiliary data processing enables our method to generalize effectively across datasets with diverse characteristics. The combination of these components ensures that the model captures both granular and high-level patterns, leading to state-of-the-art performance across all datasets. These findings validate the architectural choices and the robustness of the proposed method.

To further evaluate the robustness of MSAF-Net under real-world deployment conditions, we conducted additional ablation experiments focusing on missing modality scenarios. These tests simulate practical CPS environments where certain sensors may fail or produce unreliable data due to occlusion, noise, or hardware limitations. We examined the model’s performance when one of the input modalities—RGB, Depth, or Thermal—was intentionally removed during inference. As shown in Table 4, MSAF-Net demonstrates strong resilience, maintaining reasonable accuracy

even when critical input streams are unavailable. The RGB-only and Depth-only configurations show moderate performance degradation, while the Thermal-only case exhibits a more noticeable drop, consistent with the lower information density of thermal data alone. These results confirm that MSAF-Net can adapt to partial input conditions and retain useful representations, making it well-suited for robust CPS applications.

To provide a more comprehensive evaluation, we extended our experiments by incorporating both computational efficiency analysis and additional comparisons with recent state-of-the-art multi-modal fusion models. We report the number of floating-point operations (FLOPs) and inference time per sample to assess the practical efficiency of each method. We include comparisons with several strong baselines and recent architectures published in the past 2 years, including TransFuse, CMX, RDFNet, and M2Fuse, which have demonstrated competitive performance in RGB-D and multi-modal semantic segmentation tasks. As shown in Table 5, MSAF-Net achieves the best overall accuracy while maintaining a favorable balance between computational cost and runtime. Notably, while TransFuse and CMX offer competitive results, they come at the cost of significantly higher FLOPs. M2Fuse, although efficient, underperforms in terms of accuracy. MSAF-Net’s multi-scale attention and adaptive fusion components demonstrate both effectiveness and efficiency, validating its suitability for real-world CPS applications.

4 Methods

4.1 Overview

Image fusion has emerged as a significant field in computer vision and data processing, aimed at integrating information from multiple source images to create a composite image that preserves the most valuable features from each source. This technique is pivotal in various applications, including medical imaging, remote sensing, surveillance, and multi-modal data analysis, where the fusion of complementary data enhances decision-making, interpretation, and performance. The process of image fusion can be broadly categorized into

TABLE 4 Robustness evaluation under missing modality scenarios (on FLIR ADAS).

Input Configuration	Top-1 accuracy (%)	Relative drop (%)
RGB + Depth + Thermal (Full Input)	88.76	0.00
RGB + Depth only	86.41	−2.35
RGB only	83.27	−5.49
Depth only	81.90	−6.86
Thermal only	78.32	−10.44

The values in bold are the best values.

TABLE 5 Comparison with recent methods in terms of accuracy, FLOPs, and inference time on the FLIR ADAS dataset.

Method	Top-1 accuracy (%)	FLOPs (G)	Inference time (ms)
TransFuse [51]	87.41	89.3	153.2
CMX [52]	86.90	78.6	142.5
RDFNet [53]	84.73	52.4	102.6
M2Fuse [54]	85.11	35.7	75.8
MSAF-Net (Ours)	88.76	56.4	98.3

The values in bold are the best values.

spatial-domain and transform-domain techniques. Spatial-domain methods directly combine pixel intensities, often leading to issues like blurring or artifacts. Conversely, transform-domain techniques operate by decomposing images into multi-resolution representations, such as wavelets or pyramid transforms, and selectively merging features at different scales. Our approach builds upon the advantages of these methodologies, leveraging a novel design tailored to address domain-specific challenges and enhance fusion quality. This work introduces a unified framework for image fusion, which integrates cutting-edge advancements in neural network-based methods and signal processing techniques. The proposed methodology incorporates innovative strategies to retain structural and textural information, prevent over-smoothing, and balance contributions from input sources dynamically. Section 4.2 formalizes the image fusion problem and outlines essential mathematical notations, presenting the theoretical foundation for our method. Subsequently, in Section 4.3, we describe the architectural design of our novel model, highlighting its ability to capture multi-scale and hierarchical features effectively. Section 4.4 elaborates on the strategic innovations we introduce to optimize the fusion process, including adaptive weighting schemes and perceptual consistency measures, demonstrating their effectiveness in achieving superior fusion outcomes.

4.2 Preliminaries

The image fusion task involves integrating complementary information from multiple source images into a unified representation, ensuring that salient features from all inputs are effectively retained. This section introduces a unified framework for

image fusion, focusing on combining multiple source images from different modalities or spectral bands into a single, informative representation. The core challenge is to design an optimal fusion mapping that preserves critical information from each input while minimizing distortions and artifacts. The fusion process begins by analyzing pixel-level values across all source images, aiming to produce a fused image that retains essential spatial and spectral characteristics while suppressing noise and irrelevant features. To achieve this, many techniques operate in the transform domain, where input images are decomposed into multi-resolution components, separating low-frequency structures from high-frequency details. Fusion operators are then applied independently to these components before reconstructing the final image using an inverse transform. This approach enables selective emphasis on important features across various scales.

Advanced fusion strategies incorporate feature extraction mechanisms that transform raw images into sets of descriptive features. These features are adaptively aggregated using high-level strategies such as attention mechanisms, which assign dynamic weights based on their relevance to the final fused output. This enables the system to emphasize informative regions from each input.

The fusion process is optimized using a composite loss function that includes terms for information preservation, structural similarity, and smoothness. These loss components guide the learning of the fusion operator to ensure the resulting image is both perceptually coherent and functionally rich in content. This section introduces a unified framework for image fusion, focusing on combining multiple source images from different modalities or spectral bands into a single, informative representation. The core

challenge is to design an optimal fusion mapping that preserves critical information from each input while minimizing distortions and artifacts. The fusion process begins by analyzing pixel-level values across all source images, aiming to produce a fused image that retains essential spatial and spectral characteristics while suppressing noise and irrelevant features.

To achieve this, many techniques operate in the transform domain, where input images are decomposed into multi-resolution components, separating low-frequency structures from high-frequency details. Fusion operators are then applied independently to these components before reconstructing the final image using an inverse transform. This approach enables selective emphasis on important features across various scales.

Advanced fusion strategies incorporate feature extraction mechanisms that transform raw images into sets of descriptive features. These features are adaptively aggregated using high-level strategies such as attention mechanisms, which assign dynamic weights based on their relevance to the final fused output. This enables the system to emphasize informative regions from each input.

The fusion process is optimized using a composite loss function that includes terms for information preservation, structural similarity, and smoothness. These loss components guide the learning of the fusion operator to ensure the resulting image is both perceptually coherent and functionally rich in content.

4.3 Multi-Scale Attention-Guided Fusion Network (MSAF-Net)

To tackle the challenges associated with achieving high-quality image fusion, we propose a novel framework named the Multi-Scale Attention-Guided Fusion Network (MSAF-Net). This model is designed to extract, process, and integrate salient features from multiple source images, preserving both global structures and fine details while dynamically adjusting to the importance of different modalities (As shown in Figures 1, 2). Below, we outline three core innovations of our proposed MSAF-Net.

The Multi-Scale Attention Fusion (MSAF) module introduces a hierarchical attention mechanism to adaptively fuse features from multiple input images at different representation levels. As illustrated in Figure 3, this mechanism processes each image through a shared backbone, generating multi-level feature maps. At each level, an attention module computes pixel-wise relevance scores, enabling the model to dynamically weigh contributions from different modalities. To enhance spatial awareness, a modulation function emphasizes spatially important regions, ensuring that both global semantics and local textures are preserved during fusion.

The Cross-Level Feature Interaction mechanism further enriches representation by allowing features at one level to be informed by those at other scales. This cross-hierarchical communication is achieved by transforming and aligning features across levels using trainable transformations. Additionally, a channel-wise attention module highlights salient information, while a global self-attention strategy governs the relative importance of feature levels. Residual correction ensures spatial alignment and helps maintain consistency between interpolated features

and their native resolutions, leading to richer and more coherent representations.

The Detail-Preserving Reconstruction module is responsible for generating the final fused image by hierarchically aggregating and refining multi-scale features. Through convolutional refinement blocks and learnable aggregation weights, the model balances contributions from all feature levels. A texture refinement block further enhances high-frequency content, such as edges and textures, which might otherwise be degraded during fusion. The reconstruction process is supervised by a multi-scale loss function that emphasizes fidelity at each resolution level, as well as a gradient consistency term that aligns edge structures between the fused image and input sources. Together, these components ensure that the final output maintains both perceptual coherence and structural integrity.

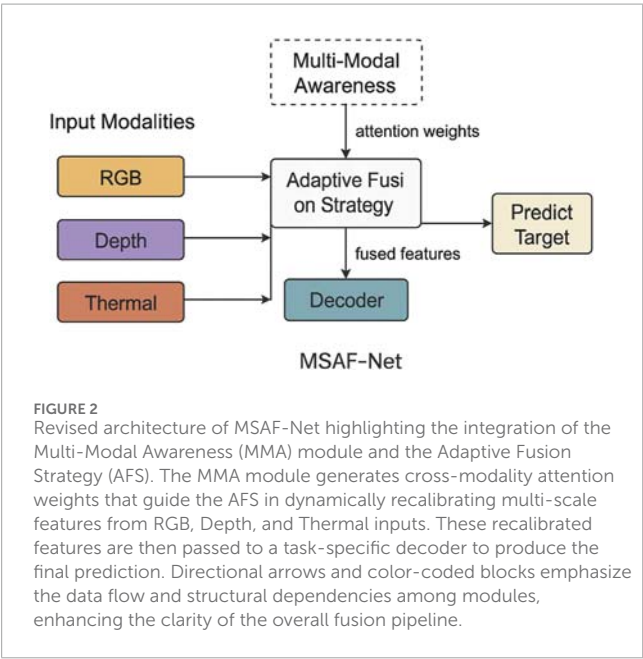
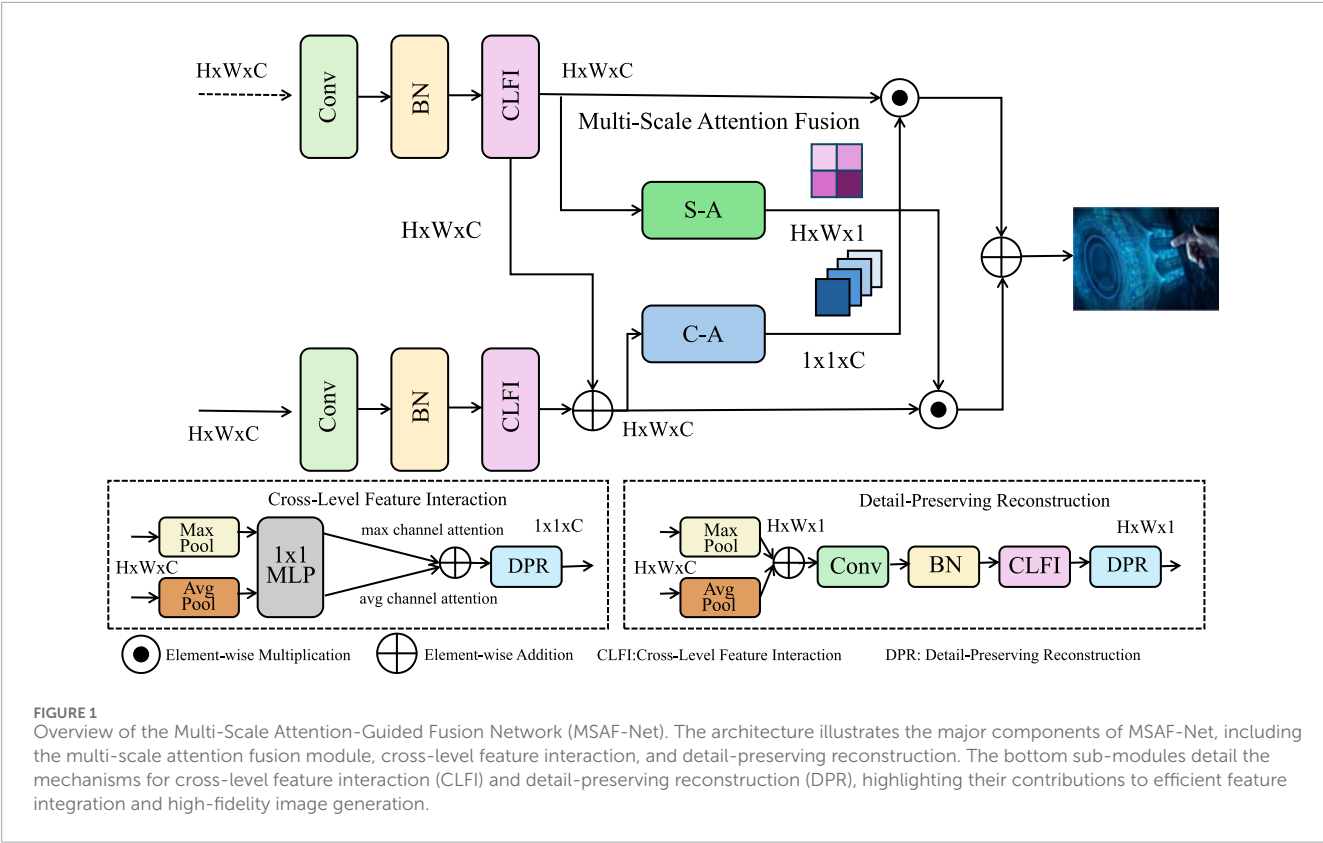
4.4 Adaptive fusion strategy with Multi-modal awareness

In this section, we propose a novel adaptive fusion strategy tailored to address the challenges of effectively combining complementary information from multiple input sources while maintaining both structural integrity and perceptual consistency (As shown in Figure 4). The proposed strategy leverages domain-specific insights, dynamic weighting mechanisms, and perceptual optimization to enhance the quality of the fused image. Below, we outline three key innovations in our approach.

The Dynamic Feature Weighting mechanism enables pixel-level adaptive fusion by learning contextual attention weights for each input modality. This allows the network to prioritize informative regions depending on their relevance—for instance, emphasizing thermal imagery in low-light conditions or RGB features under normal lighting. Attention weights are computed using a lightweight convolutional network that captures both local and global cross-modal interactions. A spatial modulation map further enhances the process by assigning spatial importance to each location, thereby refining the attention weights. Additionally, residual connections between hierarchical levels ensure feature continuity and mitigate degradation during upsampling, maintaining coherence across feature scales.

The Perceptual Consistency via Semantic Loss mechanism aims to preserve high-level semantic structures and textures in the fused image. Instead of relying solely on pixel-wise differences, the method uses a perceptual loss computed from deep feature activations extracted from a pre-trained network. This loss evaluates the fused image's alignment with a dynamically constructed pseudo-reference, formed by blending the input sources based on their relevance. The relevance of each input is learned through a scoring network and used to weigh its contribution to the reference representation. A multi-scale extension of this loss ensures that both global structures and fine details are preserved across image resolutions. Additionally, a gradient alignment term encourages the preservation of edges and textures by penalizing inconsistencies in spatial gradients between the fused and reference images.

The Multi-Scale Structural Preservation strategy is introduced to ensure that structural features such as contours, textures, and contrasts are maintained across all levels of resolution. This begins with a structural similarity loss, which measures the visual

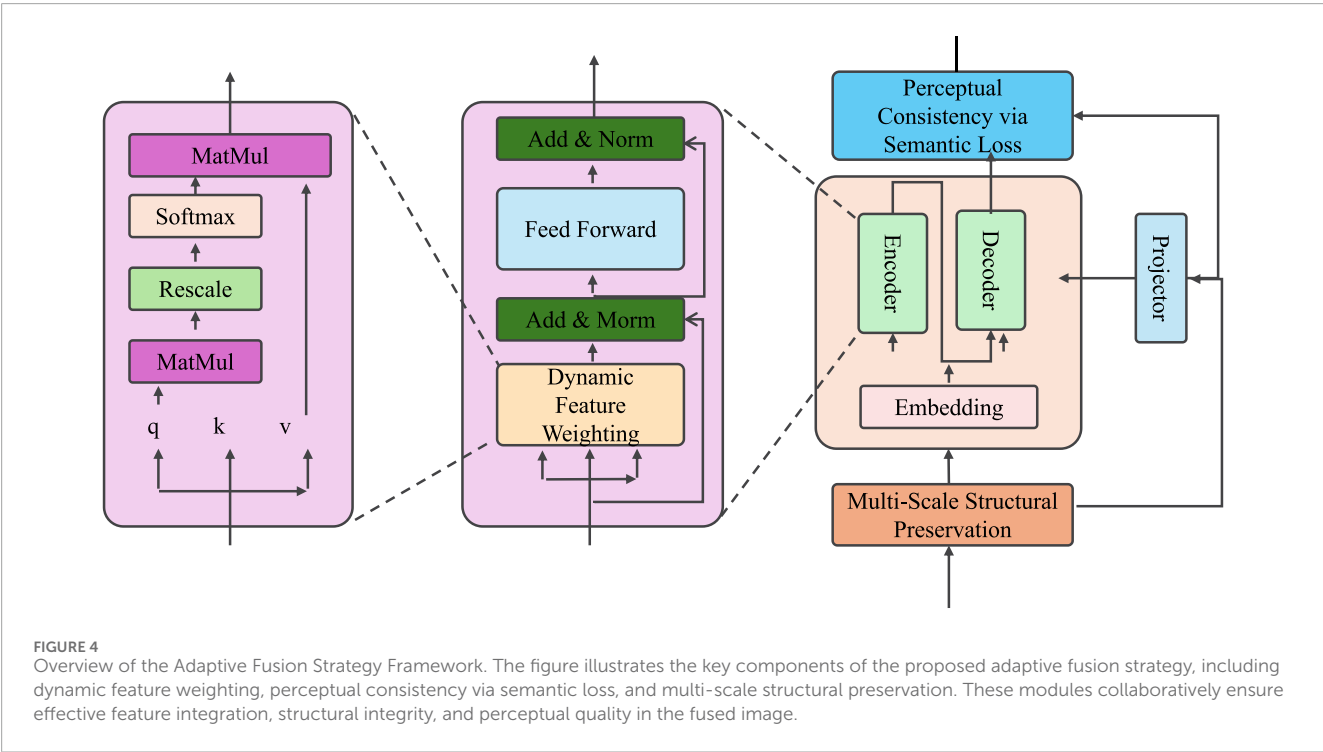
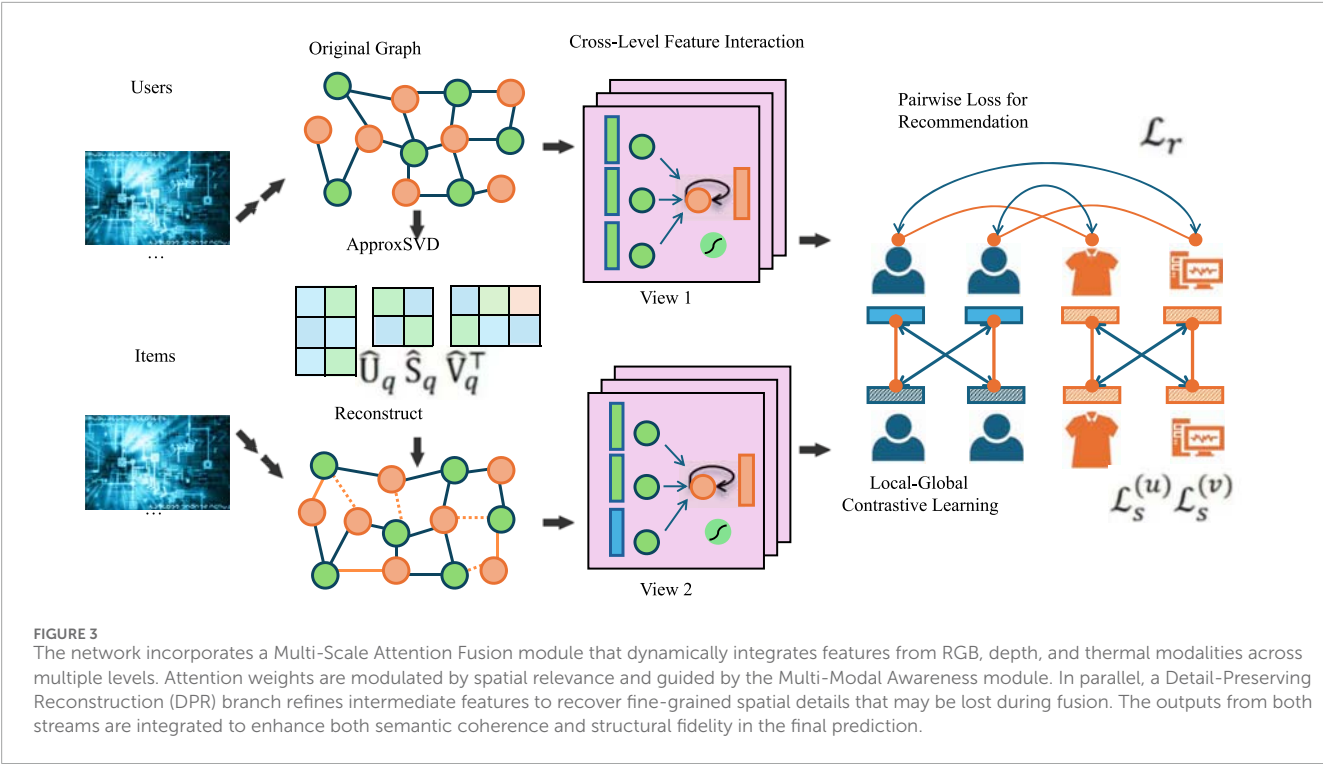


closeness of the fused image to each input source. To reinforce this, residual refinement connects feature maps across levels, ensuring that low-level details enhance high-level representations. A feature alignment operation upscales and combines information across scales, further improving structural coherence. Lastly, a Laplacian pyramid decomposition captures high-frequency details like edges

at various levels. A Laplacian consistency loss enforces similarity between the fused image's high-frequency components and those of the input images. These combined constraints ensure that the fused output is sharp, consistent, and structurally faithful to the source inputs.

5 Discussion

To further enhance the adaptability of MSAF-Net in diverse cyber-physical system scenarios, future extensions should consider the incorporation of non-visual modalities, such as inertial measurements, audio signals, or event-based sensor data. While the current model demonstrates strong performance in fusing visual modalities like RGB, depth, and infrared images, many real-world CPS applications, particularly in autonomous driving, wearable systems, and smart manufacturing, rely on multi-sensor environments where non-visual information plays a crucial role. A potential solution involves introducing a generic modality embedding module that can project heterogeneous data types into a shared latent representation space. By learning modality-specific encoders followed by unified fusion through the existing multi-scale attention mechanism, MSAF-Net could be extended to support broader modality inputs without compromising architectural integrity. Such an enhancement would enable the model to operate more robustly under visual degradation conditions and improve its generalization across sensor-rich environments. This direction represents a promising path toward building a truly multimodal



and resilient perception framework for next-generation CPS applications.

The results presented in Table 6 illustrate a clear trade-off between recognition accuracy and computational efficiency across different variants of MSAF-Net. The original MSAF-Net

achieves the highest Top-1 accuracy of 91.54% on the UCF101 dataset, but this comes at the cost of significant computational overhead, with 42.3 million parameters, 118.5 milliseconds of inference time, and 56.4 GFLOPs. When replacing the multi-scale attention mechanism with grouped attention, the model maintains

TABLE 6 Performance and computational efficiency comparison of MSAF-Net variants on UCF101.

Model variant	Top-1 accuracy (%)	Parameters (M)	Inference time (ms)	FLOPs (G)
Original MSAF-Net	91.54	42.3	118.5	56.4
w/Grouped Attention	90.78	31.2	88.6	42.9
w/Sparse Attention	90.51	33.4	85.2	39.6
w/Pruned MSAF-Net	89.92	28.7	81.3	37.1

The values in bold are the best values.

a competitive accuracy of 90.78%, while substantially reducing parameters to 31.2 million, decreasing inference time by nearly 25%, and lowering the FLOPs to 42.9G. Similarly, the sparse attention variant achieves an accuracy of 90.51% and brings further improvements in efficiency, particularly in inference latency and floating-point operations, suggesting its suitability for time-sensitive applications. The pruned version of MSAF-Net, where redundant weights are removed using L1-norm pruning, results in the smallest model with 28.7 million parameters and the fastest inference time of 81.3 milliseconds. Although the accuracy drops to 89.92%, the performance remains acceptable given the gain in efficiency. These findings indicate that integrating lightweight attention modules or pruning techniques can offer meaningful computational benefits with minimal compromise in recognition performance. Such strategies are especially promising for deployment in real-time or resource-constrained CPS environments, where both accuracy and speed are critical.

6 Conclusion and future work

This work tackles the challenge of action recognition in cyber-physical systems (CPS), which demand robust integration of multi-modal data to process diverse spatial and temporal cues effectively. Traditional methods often fall short in adaptability and fail to adequately preserve structural and textural information when fusing data from multiple modalities. To address these limitations, we proposed the Multi-Scale Attention-Guided Fusion Network (MSAF-Net), which leverages advanced image fusion techniques, multi-scale feature extraction, and attention mechanisms. The framework dynamically adjusts contributions from multiple modalities using adaptive weighting and perceptual consistency measures, mitigating issues like over-smoothing and noise sensitivity while improving generalization. Experimental results demonstrate the superiority of MSAF-Net over state-of-the-art methods, with enhanced accuracy and robustness across various CPS applications, including surveillance and human-computer interaction. This study highlights the potential of intelligent fusion strategies for advancing action recognition in complex environments. MSAF-Net’s adaptive and robust architecture suggests promising applications in medical imaging scenarios, where integrating heterogeneous modalities such as functional and anatomical scans can significantly improve the precision of medical diagnostics.

Despite its promising contributions, our proposed MSAF-Net has some limitations. First, while it significantly improves accuracy and robustness, the computational overhead introduced by multi-scale attention mechanisms and adaptive weighting schemes can be substantial. This might hinder its deployment in real-time CPS applications where low-latency processing is crucial. Future work could focus on optimizing the computational efficiency of the framework by exploring lightweight attention modules or pruning strategies. Second, the model’s adaptability across extremely heterogeneous modalities, such as integrating non-visual sensor data, remains unexplored. Extending the MSAF-Net framework to incorporate such modalities could further enhance its utility in a broader range of CPS scenarios. This direction promises to improve the resilience of action recognition systems, making them capable of handling more diverse and unpredictable real-world environments.

Data availability statement

The original contributions presented in the study are included in the article/[Supplementary Material](#), further inquiries can be directed to the corresponding author.

Author contributions

ZS: Conceptualization, Methodology, Software, Validation, Formal Analysis, Investigation, Writing – original draft. DZ: Data-curation, Writing – original draft, Writing – review and editing, Visualization, Supervision, funding-acquisition.

Funding

The author(s) declare that no financial support was received for the research and/or publication of this article.

Conflict of interest

The author declares that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Generative AI statement

The author(s) declare that no Generative AI was used in the creation of this manuscript.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the

reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

Supplementary material

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fphy.2025.1576591/full#supplementary-material>

References

1. Yang Z, Li Y, Tang X, Xie M. Mgfusion: a multimodal large language model-guided information perception for infrared and visible image fusion. *Front Neuroinformatics* (2024) 18:1521603. doi:10.3389/fnbot.2024.1521603
2. Pan R. Multimodal fusion-powered English speaking robot. *Front Neuroinformatics* (2024) 18:1478181. doi:10.3389/fnbot.2024.1478181
3. Anwar SM, Majid M, Qayyum A, Awais M, Alnowami M, Khan MK. Medical image analysis using convolutional neural networks: a review. *J Med Syst* (2018) 42:226–13. doi:10.1007/s10916-018-1088-1
4. Kahol A, Bhatnagar G. Deep learning-based multimodal medical image fusion. *Data Fusion Tech Appl Smart Healthc* (2024) 251–79. Available online at: <https://www.sciencedirect.com/science/article/pii/B9780443132339000175>.
5. Wang G. RL-cwtrans net: multimodal swimming coaching driven via robot vision. *Front Neuroinformatics* (2024) 18:1439188. doi:10.3389/fnbot.2024.1439188
6. Cheng K, Zhang Y, He X, Chen W, Cheng J, Lu H. Skeleton-based action recognition with shift graph convolutional network. *Computer Vis Pattern Recognition* (2020). Available online at: http://openaccess.thecvf.com/content_CVPR_2020/html/Cheng_Skeleton-Based_Action_Recognition_With_Shift_Graph_Convolutional_Network_CVPR_2020_paper.html.
7. Zhou H, Liu Q, Wang Y. Learning discriminative representations for skeleton based action recognition. *Computer Vis Pattern Recognition* (2023) 10608–17. doi:10.1109/cvpr52729.2023.01022
8. Li Y, Ji B, Shi X, Zhang J, Kang B, Wang L. Tea: temporal excitation and aggregation for action recognition. *Computer Vis Pattern Recognition* (2020). Available online at: http://openaccess.thecvf.com/content_CVPR_2020/html/Li_TEA_Temporal_Excitation_and_Aggregation_for_Action_Recognition_CVPR_2020_paper.html.
9. Morshed MG, Sultana T, Alam A, Lee Y-K. Human action recognition: a taxonomy-based survey, updates, and opportunities. *Ital Natl Conf Sensors* (2023) 23:2182. doi:10.3390/s23042182
10. Perrett T, Masullo A, Burghardt T, Mirmehdi M, Damen D. Temporal-relational crosstransformers for few-shot action recognition. *Computer Vis Pattern Recognition* (2021). Available online at: http://openaccess.thecvf.com/content/CVPR2021/html/Perrett_Temporal-Relational_CrossTransformers_for_Few-Shot_Action_Recognition_CVPR_2021_paper.html.
11. Yang C, Xu Y, Shi J, Dai B, Zhou B. Temporal pyramid network for action recognition. *Computer Vis Pattern Recognition* (2020). Available online at: http://openaccess.thecvf.com/content_CVPR_2020/html/Yang_Temporal_Pyramid_Network_for_Action_Recognition_CVPR_2020_paper.html.
12. gun Chi H, Ha MH, geun Chi S, Lee SW, Huang Q-X, Ramani K. Infogcn: representation learning for human skeleton-based action recognition. *Computer Vis Pattern Recognition* (2022) 20154–64. doi:10.1109/cvpr52688.2022.01955
13. Wang L, Tong Z, Ji B, Wu G. Tdn: temporal difference networks for efficient action recognition. *Computer Vis Pattern Recognition* (2020). Available online at: http://openaccess.thecvf.com/content/CVPR2021/html/Wang_TDN_Temporal_Difference_Networks_for_Efficient_Action_Recognition_CVPR_2021_paper.html.
14. Pan J, Lin Z, Zhu X, Shao J, Li H. St-adapter: parameter-efficient image-to-video transfer learning for action recognition. *Neural Inf Process Syst* (2022). Available online at: https://proceedings.neurips.cc/paper_files/paper/a92e9165b22d4456f6d87236e04c266-Abstract-Conference.html.
15. Song Y, Zhang Z, Shan C, Wang L. Constructing stronger and faster baselines for skeleton-based action recognition. *IEEE Trans Pattern Anal Machine Intelligence* (2021) 45:1474–88. doi:10.1109/tpami.2022.3157033
16. Sun Z, Liu J, Ke Q, Rahmani H, Wang G. Human action recognition from various data modalities: a review. *IEEE Trans Pattern Anal Machine Intelligence* (2020) 45:3200–25. doi:10.1109/tpami.2022.3183112
17. Chen Z, Li S, Yang B, Li Q, Liu H. Multi-scale spatial temporal graph convolutional network for skeleton-based action recognition. *AAAI Conf Artif Intelligence* (2021) 35:1113–22. doi:10.1609/aaai.v35i2.16197
18. Ye F, Pu S, Zhong Q, Li C, Xie D, Tang H. Dynamic gcn: context-enriched topology learning for skeleton-based action recognition. *ACM Multimedia* (2020) 55–63. doi:10.1145/3394171.3413941
19. Zhang H, Zhang L, Qi X, Li H, Torr PHS, Koniusz P. Few-shot action recognition with permutation-invariant attention. *Eur Conf Computer Vis* (2020) 525–42. doi:10.1007/978-3-030-58558-7_31
20. Duan H, Wang J, Chen K, Lin D. Pyskl: towards good practices for skeleton action recognition. *ACM Multimedia* (2022) 7351–4. doi:10.1145/3503161.3548546
21. Lin L, Song S, Yang W, Liu J. Ms2l: multi-task self-supervised learning for skeleton based action recognition. *ACM Multimedia* (2020). Available online at: <https://dl.acm.org/doi/abs/10.1145/3394171.3413548>.
22. Song Y, Zhang Z, Shan C, Wang L. Stronger, faster and more explainable: a graph convolutional baseline for skeleton-based action recognition. *ACM Multimedia* (2020) 1625–33. doi:10.1145/3394171.3413802
23. Munro J, Damen D. Multi-modal domain adaptation for fine-grained action recognition. *Computer Vis Pattern Recognition* (2020) 119–29. doi:10.1109/cvpr42600.2020.00020
24. Wang X, Zhang S, Qing Z, Tang M, Zuo Z, Gao C, et al. Hybrid relation guided set matching for few-shot action recognition. *Computer Vis Pattern Recognition* (2022) 19916–25. doi:10.1109/cvpr52688.2022.01932
25. Yang J, Dong X, Liu L, Zhang C, Shen J, Yu D. Recurring the transformer for video action recognition. *Computer Vis Pattern Recognition* (2022) 14043–53. doi:10.1109/cvpr52688.2022.01367
26. Chang H-L, Ren H-T, Wang G, Yang M, Zhu X-Y. Infrared defect recognition technology for composite materials. *Front Phys* (2023) 11:1203762. doi:10.3389/fphy.2023.1203762
27. Dave I, Chen C, Shah M. Spact: self-supervised privacy preservation for action recognition. *Computer Vis Pattern Recognition* (2022) 20132–41. doi:10.1109/cvpr52688.2022.01953
28. Xing Z, Dai Q, Hu H-R, Chen J, Wu Z, Jiang Y-G. Svformer: semi-supervised video transformer for action recognition. *Computer Vis Pattern Recognition* (2022). Available online at: http://openaccess.thecvf.com/content/CVPR2023/html/Xing_SVFormer_Semi-Supervised_Video_Transformer_for_Action_Recognition_CVPR_2023_paper.html.
29. Wang Z, She Q, Smolic A. Action-net: multipath excitation for action recognition. *Computer Vis Pattern Recognition* (2021) 13209–18. doi:10.1109/cvpr46437.2021.01301
30. Jin X, Zhang P, He Y, Jiang Q, Wang P, Hou J A theoretical analysis of continuous firing condition for pulse-coupled neural networks with its applications. *Eng Appl Artif Intelligence* (2023) 126:107101. doi:10.1016/j.engappai.2023.107101
31. Meng Y, Lin C-C, Panda R, Sattigeri P, Karlinsky L, Oliva A, et al. Ar-net: adaptive frame resolution for efficient action recognition. *Eur Conf Computer Vis* (2020) 86–104. doi:10.1007/978-3-030-58571-6_6
32. Truong T-D, Bui Q-H, Duong C, Seo H-S, Phung SL, Li X, et al. Direformer: a directed attention in transformer approach to robust action recognition. *Computer Vis Pattern Recognition* (2022) 19998–20008. doi:10.1109/cvpr52688.2022.01940
33. Mahdhi N, Alsaiani NS, Amari A, Osman H, Hammami S. Enhancement of the physical adsorption of some insoluble lead compounds from drinking water onto polylactic acid and graphene oxide using molybdenum disulfide nanoparticles: theoretical investigation. *Front Phys* (2023) 11:1159306. doi:10.3389/fphy.2023.1159306
34. Bao W, Yu Q, Kong Y. Evidential deep learning for open set action recognition. *IEEE Int Conf Computer Vis* (2021) 13329–38. doi:10.1109/iccv48922.2021.01310

35. Li Y, Jian P, Han G. Cascaded progressive generative adversarial networks for reconstructing three-dimensional grayscale core images from a single two-dimensional image. *Front Phys* (2022) 10:716708. doi:10.3389/fphy.2022.716708
36. Chen Y, Zhang Z, Yuan C, Li B, Deng Y, Hu W. Channel-wise topology refinement graph convolution for skeleton-based action recognition. *IEEE Int Conf Computer Vis* (2021) 13339–48. doi:10.1109/iccv48922.2021.01311
37. Duan H, Zhao Y, Chen K, Shao D, Lin D, Dai B. Revisiting skeleton-based action recognition. *Computer Vis Pattern Recognition* (2021). Available online at: http://openaccess.thecvf.com/content/CVPR2022/html/Duan_Revisiting_Skeleton-Based_Action_Recognition_CVPR2022_paper.html.
38. Liu KZ, Zhang H, Chen Z, Wang Z, Ouyang W. Disentangling and unifying graph convolutions for skeleton-based action recognition. *Computer Vis Pattern Recognition* (2020) 140–9. doi:10.1109/cvpr42600.2020.00022
39. Jin X, Wu N, Jiang Q, Kou Y, Duan H, Wang P A dual descriptor combined with frequency domain reconstruction learning for face forgery detection in deepfake videos. *Forensic Sci Int Digital Invest* (2024) 49:301747. doi:10.1016/j.fsidi.2024.301747
40. Jin X, Liu L, Ren X, Jiang Q, Lee S-J, Zhang J A restoration scheme for spatial and spectral resolution of the panchromatic image using the convolutional neural network. *IEEE J Selected Top Appl Earth Observations Remote Sensing* (2024) 17:3379–93. doi:10.1109/jstars.2024.3351854
41. Farooq MA, Corcoran P, Rotariu C, Shariff W. Object detection in thermal spectrum for advanced driver-assistance systems (adas). *IEEE Access* (2021) 9:156465–81. doi:10.1109/access.2021.3129150
42. Zunair H, Khan S, Hamza AB. Rsud20k: a dataset for road scene understanding in autonomous driving. *arXiv preprint arXiv:2401.07322* (2024) 708–14. doi:10.1109/icip51287.2024.10648203
43. Sachdeva K, Sandhu JK, Sahu R. Exploring video event classification: leveraging two-stage neural networks and customized cnn models with ucf-101 and ccv datasets. In: *2024 11th international conference on computing for sustainable global development (INDIACom)*. IEEE (2024). p. 100–5.
44. Patel D, Parikh R, Shastri Y. Recent advances in video question answering: a review of datasets and methods. In: *Pattern recognition. ICPR international workshops and challenges: virtual event, january 10–15, 2021, proceedings, Part II*. Springer (2021). p. 339–56.
45. Archana N, Hareesh K. Real-time human activity recognition using resnet and 3d convolutional neural networks. In: *2021 2nd international conference on advances in computing, communication, embedded and secure systems (ACCESS)*. IEEE (2021). p. 173–7.
46. Tan H, Cheng R, Huang S, He C, Qiu C, Yang F, et al. Relativenas: relative neural architecture search via slow-fast learning. *IEEE Trans Neural Networks Learn Syst* (2021) 34:475–89. doi:10.1109/tnnls.2021.3096658
47. Peng Y, Lee J, Watanabe S. I3d: transformer architectures with input-dependent dynamic depth for speech recognition. In: *ICASSP 2023-2023 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE (2023). p. 1–5.
48. Seijo O, Iturbe X, Val I. Tackling the challenges of the integration of wired and wireless tsn with a technology proof-of-concept. *IEEE Trans Ind Inform* (2021) 18:7361–72. doi:10.1109/tii.2021.3131865
49. Umi U, Anzelina D, Ade Muhayati R, Suhedi H. Kesehatan mental dan tarekat overthinking dalam perspektif ponpes tarekat qadiriyyah wa naqsyabandiyah (tqn) al-mubarak cinangka. *Mutiara: Multidisciplinary Scientific J* (2024) 2:591–601. doi:10.57185/mutiara.v2i7.214
50. Pham Q, Liu C, Hoi SC. Continual learning, fast and slow. *IEEE Trans Pattern Anal Machine Intelligence* (2023) 46:134–49. doi:10.1109/tpami.2023.3324203
51. Soliman A, Soliman A. Late mean fusion towards efficient polyps segmentation. In: *2024 6th novel intelligent and leading emerging sciences conference (NILES)*. IEEE (2024). p. 233–7.
52. Zhang A, Zhu M, Zheng Y, Tian Z, Mu G, Zheng M. The significant contribution of comammox bacteria to nitrification in a constructed wetland revealed by dna-based stable isotope probing. *Bioresour Technology* (2024) 399:130637. doi:10.1016/j.biortech.2024.130637
53. Jia W, Yan X, Liu Q, Zhang T, Dong X. Tcanet: three-stream coordinate attention network for rgb-d indoor semantic segmentation. *Complex and Intell Syst* (2024) 10:1219–30. doi:10.1007/s40747-023-01210-4
54. Cai Y, Liu Q, Gan Y, Lin R, Li C, Liu X, et al. Difinet: boundary-aware semantic differentiation and filtration network for nested named entity recognition. *Proc 62nd Annu Meet Assoc Comput Linguistics* (2024) 1:6455–71. Available online at: <https://aclanthology.org/2024.acl-long.349/>.

Frontiers in Physics

Investigates complex questions in physics to understand the nature of the physical world

Addresses the biggest questions in physics, from macro to micro, and from theoretical to experimental and applied physics.

Discover the latest Research Topics

[See more →](#)

Frontiers

Avenue du Tribunal-Fédéral 34
1005 Lausanne, Switzerland
frontiersin.org

Contact us

+41 (0)21 510 17 00
frontiersin.org/about/contact

