

DAE-Fuse: An Adaptive Discriminative Autoencoder for Multi-Modality Image Fusion

Yuchen Guo^{*‡}
Department of Computer Science
BNU-HKBU UIC
 Zhuhai, China
 r130026037@mail.uic.edu.cn

Ruoxiang Xu^{*}
Department of Computer Science
BNU-HKBU UIC
 Zhuhai, China
 r130026173@mail.uic.edu.cn

Rongcheng Li
Department of Computer Science
BNU-HKBU UIC
 Zhuhai, China
 r130026074@mail.uic.edu.cn

Zhenghao Wu
School of Computer Science
University College Dublin
 Dublin, Ireland
 hi@ecwu.xyz

Weifeng Su[†]
Guangdong Provincial Key Laboratory of
Interdisciplinary Research and Application for Data Science
BNU-HKBU UIC
 Zhuhai, China
 wfsu@uic.edu.cn

Abstract—Multi-modality image fusion aims to integrate complementary data information from different imaging modalities into a single image. Existing methods often generate either blurry fused images that lose fine-grained semantic information or unnatural fused images that appear perceptually cropped from the inputs. In this work, we propose a novel two-phase discriminative autoencoder framework, termed *DAE-Fuse*, that generates sharp and natural fused images. In the adversarial feature extraction phase, we introduce two discriminative blocks into the encoder-decoder architecture, providing an additional adversarial loss to better guide feature extraction by reconstructing the source images. While the two discriminative blocks are adapted in the attention-guided cross-modality fusion phase to distinguish the structural differences between the fused output and the source inputs, injecting more naturalness into the results. Extensive experiments on public infrared-visible, medical image fusion, and downstream object detection datasets demonstrate our method’s superiority and generalizability in both quantitative and qualitative evaluations.

Index Terms—Image fusion, Generative model, Multi-modality.

I. INTRODUCTION

Multi-Modality Image Fusion (MMIF), a hot image processing topic in the multimedia and low-level computer vision community, aims to render fused images that maintain the essential information of different modalities. This trait allows the fused images to describe a better visual understanding and also can be applied to subsequent high-level vision tasks, *e.g.*, detection [1], [2], [3], and segmentation [4], [5]. In particular, the Infrared-Visible Image Fusion (IVIF) is a representative fusion task that has been widely applied [6], [7], [8], [9]. Infrared images effectively capture thermal targets in dark environments but lack texture details, which can hinder recognition in applications. On the contrary, visible images maintain most of the textual details but are sensitive to light conditions. The IVIF task aspires to combine the advantages of both images by fusing the thermal radiation information and texture details

into a new image that can thoroughly describe the actual scene, improving the performance of various downstream tasks like Multi-Modality Object Detection (MMOD). The Medical Image Fusion (MIF) aims at combining information from various medical imaging modalities to generate a more comprehensive and detailed representation of anatomical structures that can help diagnosis and treatment [10].

GAN-based models use adversarial learning with zero-sum games in a fused image and source images to fuse two inputs. The usual strategy in MMIF task is that totally two discriminators are employed to discriminate with the fusion results and the two source images [11], [12], [13], [1]. Most of them either fuse the two-dimensional image pairs before input to the model [14], [1], or did not tailor a feature extractor and corresponding loss function to distinctively extract features with different characteristics, weakening to feature extraction ability. Those methods just generate the fused image which is perceptually satisfactory by looking distributionally similar to the original data, but fail to preserve the feature details, resulting in the blurriness within and between functional objects.

More efficient pipelines take comprehensive feature extraction and reconstruction modules into an AE-based manner [15], [16], [17], [18], [19]. They separately encode the two input modalities and fuse the feature embeddings via channel concatenation, then decode the fused embeddings to an output image. By the manually elaborated encoder block and loss functions, the AE-based methods tend to effectively extract both global and local features from different modalities. Usually, the encoder and its loss are shared for both inputs [15], [17], [18], and they concatenate features directly rather than organically combining features from different modalities in the fusion phase. Therefore, bias between modalities may be introduced to the fused images, making the fused images present more obvious traces from the image of a specific modality but are left with inconspicuous information from another modality, which can be verified from the experiments.

In order to solve the aforementioned problems, we devel-

^{*} Equal Contribution, [†] Corresponding Author, [‡] Part of this research was performed while Yuchen Guo was at MMLab@SIAT.

oped a novel end-to-end discriminative autoencoder model for multi-modality image fusion (**DAE-Fuse**), which adopts a two-phase adversarial learning, and a cross-attention fusion module that endows the model with a more balanced fusion capability together with strong generalizability. Qualitative and quantitative experiments show that our model has achieved state-of-the-art on multiple IVIF public datasets, and the superiority can also be generalized to different MIF tasks. More importantly, our approach can boost the performance on downstream MMOD tasks without any fine-tuning.

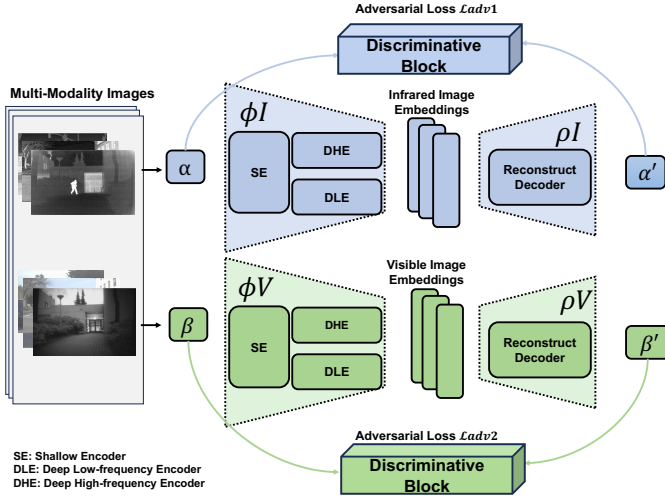


Fig. 1: The workflow of the adversarial feature extraction phase. The cross-attention for fusion purpose is dismissed.

II. METHOD

A. Adversarial Feature Extraction Phase

The multi-level features are extracted by shallow and deep encoders. Specifically, to differentiate the various frequencies features, we deploy a Deep High-frequency Encoder (DHE, termed $\phi_{DH}(\cdot)$) and a Deep Low-frequency Encoder (DLE, termed $\phi_{DL}(\cdot)$) parallelly following the Shallow Encoder (SE, termed $\phi_S(\cdot, \cdot)$). Suppose the embedding from the encoding process is marked as: $\Phi(\cdot)$, and the input of first and second modalities as: α , and β . The encoding process of paired $\{\alpha, \beta\}$ can be formulated as:

$$\begin{aligned}\Phi(\alpha) &= C[\phi_{DH}(\phi_S(\alpha)), \phi_{DL}(\phi_S(\alpha))] \\ \Phi(\beta) &= C[\phi_{DH}(\phi_S(\beta)), \phi_{DL}(\phi_S(\beta))]\end{aligned}\quad (1)$$

where $C(\cdot, \cdot)$ donates channel concatenate operation.

Since the Transformer-based models are good at extracting low-frequency information while CNN-based models are sensitive to high-frequency information [20], [21]. We construct a Vision Transformer [22] for the DLE, and the DHE is implemented by a ResNet18 [23]. Restormer is a Channel-Transformer [24] architecture, which has achieved excellent performance in shallow region reconstruction task without increasing too much computation, so we use a channel-Transformer block for SE to extract shallow features. While

the Reconstruction Decoder (RD, termed $\rho_R(\cdot)$) is responsible for reconstructing the embeddings to image. RD shares the same architecture with SE. The decoding process of paired $\{\alpha, \beta\}$ can be formulated as:

$$\tilde{\alpha} = \rho_R(\alpha), \tilde{\beta} = \rho_R(\beta) \quad (2)$$

where $\tilde{\alpha}$ and $\tilde{\beta}$ represent the images α and β after reconstructing, respectively.

The adversarial process is implemented by two discriminative blocks from different modalities (DM1 and DM2, termed $D_{M1}(\cdot, \cdot)$ and $D_{M2}(\cdot, \cdot)$ respectively). Discriminative blocks is implemented by a stack of Con2D-LeakyReLU-BatchNorm layers and a fully connected layer. Accordingly, the adversarial learning process can be formulated as minimizing the following adversarial objective:

$$\begin{aligned}\min_{AE} \max_{D_{M1}, D_{M2}} & \left(\mathbb{E}[\log(D_{M1}(\alpha))] + \mathbb{E}[\log(D_{M2}(\beta))] \right. \\ & \left. + \mathbb{E}[\log(1 - (D_{M1}(\tilde{\alpha})))] + \mathbb{E}[\log(1 - (D_{M2}(\tilde{\beta})))] \right)\end{aligned}\quad (3)$$

B. Attention-guided Cross-modality Fusion Phase

In this phase, we developed a feature aggregation strategy by calculating the cross-attention weights, which is analogous to the standard attention of Transformer [22]. We use the same structure of discriminative blocks as before. During the adversarial fusion step, the inputs of a discriminative block are a fused image two source images.

1) *Early Fusion*: Owing to the data gap between different modalities, current approaches in MMIF are limited to only incorporating element-wise additions for extracted feature embeddings, which does not capture the important interactions. We deploy a cross-modality attention module, making the different embedding can naturally interact another modality before fusion. After extracting features from encoders ($\Phi(\alpha)$, $\Phi(\beta)$), embeddings of images from two modalities are obtained. Here we use the embedding of α as the Query Q , while the embeddings of β as the Key K and the Value V . Assuming the attention guided embeddings are denoted as: ($\Phi(\hat{\alpha})$) and ($\Phi(\hat{\beta})$).

2) *Adversarial Fusion*: First, the decoder generates fused image from attention-guided embeddings:

$$\mathcal{F}(\alpha, \beta) = \rho_R[C(\Phi(\hat{\alpha}), \Phi(\hat{\beta}))] \quad (4)$$

Then, the adversarial process is adapted to following formulation:

$$\begin{aligned}\min_{AE} \max_{D_{M1}, D_{M2}} & \left(\mathbb{E}[\log(D_{M1}(\alpha))] + \mathbb{E}[\log(D_{M2}(\beta))] + \right. \\ & \left. \mathbb{E}[\log(1 - (D_{M1}(\mathcal{F}(\alpha, \beta))))] + \mathbb{E}[\log(1 - (D_{M2}(\mathcal{F}(\alpha, \beta))))] \right)\end{aligned}\quad (5)$$

C. Loss Function

1) *Phase one*: We construct the loss function for the autoencoder and discriminative blocks separately. The loss

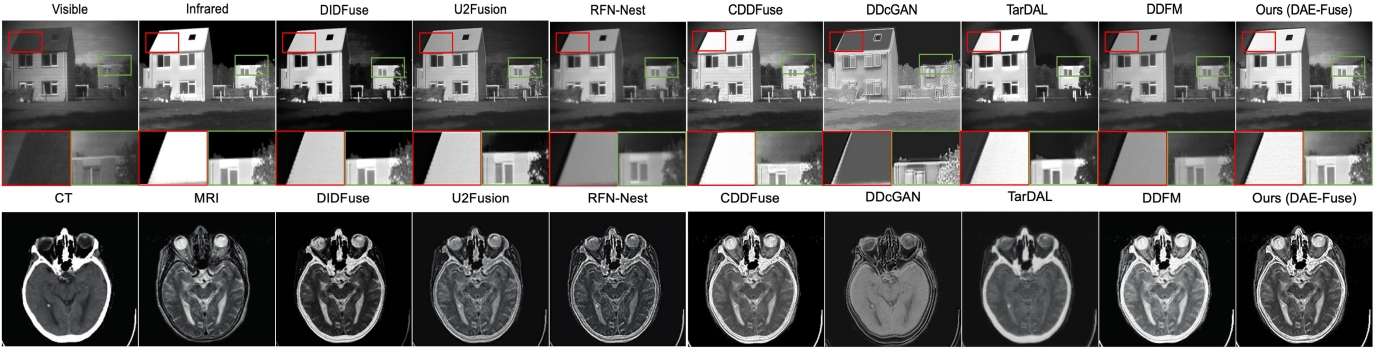


Fig. 2: Qualitative comparison with state-of-the-art methods on TNO and MRI-CT dataset.

function of AE is divided into two parts: adversarial loss and content loss:

$$\mathcal{L}_{AE}^I = \lambda \mathcal{L}_{AE}^{advI} + \sigma \mathcal{L}_{Enc}^{correlation} + (1 - \sigma) \mathcal{L}_{Dec}^{content} \quad (6)$$

where the σ is the hyper-parameter. And the adversarial loss for encoder-decoder is:

$$\mathcal{L}_{AE}^{advI} = \mathbb{E}[\log(1 - (D_{M1}(\tilde{\alpha})))] + \mathbb{E}[\log(1 - (D_{M2}(\tilde{\beta})))] \quad (7)$$

Additionally, we use correlation decomposition loss [16] for differentiate high-frequency feature and low-frequency feature:

$$\mathcal{L}_{Enc}^{correlation} = \frac{(CC(\phi_{DH}(\alpha), \phi_{DH}(\beta))^2}{CC(\phi_{DL}(\alpha), \phi_{DL}(\beta)) + \epsilon} \quad (8)$$

where ϵ set 1.01 to ensure the result always be positive.

The decoder reconstruction loss function consists of the square of the L2 norm and structural similarity index:

$$\mathcal{L}_{Dec}^{content} = \|\alpha - D(\alpha)\|_2^2 + (1 - SSIM(\alpha, D(\alpha))) \quad (9)$$

The adversarial loss of discriminative block DM1 and DM2 are of same structure. Take DM1 as an example:

$$\mathcal{L}_{DM1}^{advI} = \mathbb{E}[-\log(D_{M1}(\alpha))] + \mathbb{E}[-\log(1 - (D_{M1}(\tilde{\alpha})))] \quad (10)$$

As a sum, the total discriminative block losses can be formulated in:

$$\mathcal{L}_{DM}^{advI} = \mathcal{L}_{DM1}^{advI} + \mathcal{L}_{DM2}^{advI} \quad (11)$$

The overall loss function of phase one is defined as:

$$\mathcal{L}^I = \mathcal{L}_{DM}^{advI} + \mathcal{L}_{AE}^I \quad (12)$$

2) *Phase two*: Since the inputs of discriminative blocks have been changed, we represent the adversarial loss of phase two \mathcal{L}_{DM}^{advII} follow the Eq. 5.

The content loss function for decoder in phase two can be formulated as:

$$\mathcal{L}_{AE}^{II} = \mathcal{L}_{text}^{II} + \mathcal{L}_{int}^{II} + \mathcal{L}_{AE}^{advII} \quad (13)$$

And the adversarial loss for the encoder-decoder is defined as:

$$\begin{aligned} \mathcal{L}_{AE}^{advII} = & \mathbb{E}[-\log(1 - (D_{M1}(\mathcal{F}(\alpha, \beta)))] \\ & + \mathbb{E}[-\log(1 - (D_{M2}(\mathcal{F}(\alpha, \beta)))] \end{aligned} \quad (14)$$

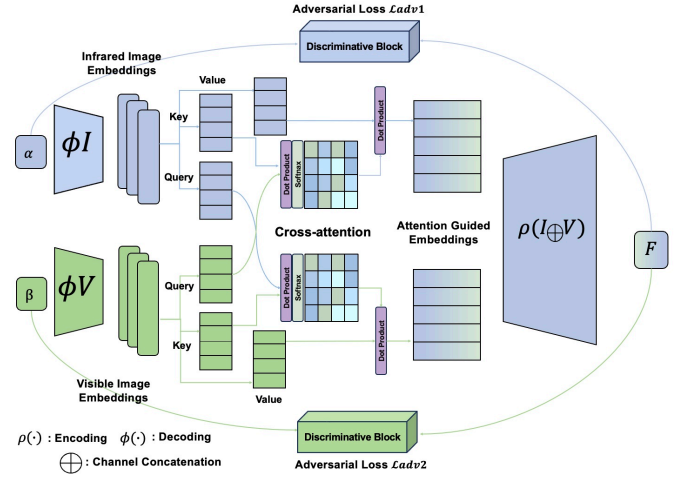


Fig. 3: The workflow of the attention-guided cross-modality fusion phase.

Also, we use the structural content loss [25] as:

$$\mathcal{L}_{text}^{II} = \frac{1}{HW} \|\nabla I_f - \max(|\nabla \alpha|, |\nabla \beta|)\|_1 \quad (15)$$

$$\mathcal{L}_{int}^{II} = \frac{1}{HW} \|I_f - \max(\alpha, \beta)\|_1 \quad (16)$$

Thus, the whole losses in phase two can formulated as:

$$\mathcal{L}^{II} = \mathcal{L}_{DM}^{advII} + \mathcal{L}_{AE}^{II} \quad (17)$$

III. EXPERIMENTS

A. Setup

1) *Datasets and metrics*: The IVIF experiments use three benchmarks: MSRS [26], RoadScene [27], and TNO [28], with only part of MSRS images used for training. MIF experiments utilize data from the Harvard 670 Medical Website [29] for testing.

For the MMOD [1] downstream task, the M3FD dataset, consisting of 4200 infrared-visible image pairs across six categories, is employed. IVIF and MIF tasks are evaluated using eight metrics: entropy (EN) [30], standard deviation (SD) [31], spatial frequency (SF) [32], visual information fidelity

TABLE I: Quantitative comparisons on TNO (IVIF), MRI-CT (MIF), and RoadScene (MMOD) datasets. Bold red indicates the best, and bold blue indicates the second best.

Method	Dataset: TNO							Dataset: MRI-CT							Dataset: RoadScene						
	EN	SD	SF	MI	SCD	VIF	Qabf	EN	SD	SF	MI	SCD	VIF	Qabf	Peo	Car	Lam	Bus	Mot	Tru	mAP@50%
DIDFuse	6.97	45.12	12.59	1.63	1.71	0.58	0.42	4.37	58.34	34.64	1.71	0.69	0.41	0.38	0.791	0.924	0.857	0.833	0.787	0.788	0.830
U2Fusion	6.83	35.66	11.52	1.35	1.71	0.61	0.44	4.21	61.98	32.54	2.08	0.75	0.37	0.46	0.802	0.922	0.870	0.839	0.783	0.786	0.833
RFN-Nest	6.84	34.50	13.23	1.76	1.67	0.55	0.39	4.97	70.36	33.42	1.98	0.68	0.43	0.52	0.813	0.915	0.851	0.829	0.813	0.875	0.849
DDcGAN	6.78	46.33	11.68	1.78	1.72	0.48	0.35	4.26	62.56	30.61	1.72	0.65	0.38	0.42	0.797	0.908	0.832	0.895	0.805	0.872	0.851
TarDAL	7.02	45.63	10.68	2.17	1.62	0.57	0.32	4.35	61.14	28.38	1.94	0.92	0.32	0.56	0.835	0.947	0.854	0.928	0.811	0.874	0.874
CDDFuse	7.12	45.89	13.15	2.11	1.76	0.76	0.54	4.49	71.36	34.02	2.16	1.18	0.44	0.56	0.846	0.928	0.864	0.931	0.813	0.891	0.878
DDFM	7.06	48.42	13.03	2.06	1.66	0.83	0.49	4.77	69.35	32.77	1.98	1.03	0.41	0.54	0.837	0.926	0.869	0.927	0.809	0.882	0.875
Ours	7.17	46.63	13.31	2.21	1.89	0.75	0.55	4.83	76.19	35.56	2.20	1.21	0.49	0.57	0.855	0.931	0.874	0.949	0.822	0.890	0.887

(VIF) [33], sum of correlation of differences (SCD) [32], mutual information (MI) [32], $Q^{AB/F}$ [32], and structural similarity index measure (SSIM) [34]. For MMOD, detection performance is measured by mAP@50% with higher values indicating better results.

2) *Implementation Details*: Our experiments were implemented based on the PyTorch framework and performed on a server with an NVIDIA A100 GPU. In the MMOD downstream testing, the generated 4200 fused images is partitioned into training, validation, and test sets, with an 8:1:1 ratio for a YOLOv8n [35].

B. Infrared-Visible Image Fusion

We tested our model on the three IVIF datasets and compared them with seven state-of-the-art methods including DIDFuse [15], U2Fusion [36], RFN-Nest [19], DDcGAN [14], TarDAL [1], CDDFuse [16] and DDFM [37].

1) *Qualitative comparisons*: As shown in Figure 2, the selected scenario highlights feature extraction and model bias. DDcGAN and TarDAL, both GAN-based models, exhibit noticeable blurriness and detail loss. DIDFuse, an AE-based method, shows a clear bias towards infrared images, darkening the sky and grass. Our DAE-Fuse, however, delivers the best results, preserving rich textural details and balancing both input modalities seamlessly. Notably, it maintains intact roof textures, unlike other methods. CDDFuse, which also uses a parallel encoder, partially retains roof details but overexposes some areas and introduces extra noise on the house wall. In contrast, DAE-Fuse fuses wall textures naturally from both inputs.

2) *Quantitative comparisons*: Afterward, we used the seven metrics to quantitatively compare the results with other models, which are displayed in Table I. DAE-Fuse shows an outstanding performance across all the measurement indices, demonstrating the effectiveness of our method.

C. Generalization on MIF tasks

To validate the generalizability of our DAE-Fuse, we used the same models in our IVIF testing to fuse medical images.

1) *Qualitative comparisons*: Figure 2 compares the MRI-CT fusion results of different models. Similar to IVIF experiments, GAN-based models appear fused but are blurry and



Fig. 4: Example of the distinct detection ability of the source images and our fused image.

detail-poor. AE-based methods generally perform better in feature extraction. However, models like U2fusion and RFN-Nest underweight CT images, leading to darker outlines, whereas DIDFuse emphasizes CT features. Our method integrates both images effectively, preserving all texture details without bias, demonstrating superior generalization.

2) *Quantitative comparisons*: Similarly, seven metrics are adopted to quantitatively compare the result, which are displayed in Table I. DAE-Fuse has the best score on almost all metrics, indicating that our method can be generalized to MIF tasks without any adjustment, and suitable for various kinds of MIF tasks.

D. Downstream MMOD task

Multi-Modality Object Detection is an important downstream task of image fusion. A single modality image *i.e.*, an individual infrared image usually lacks certain features of objects during the detection process. As shown in Figure 4, infrared images exhibit a robust capability for detecting humans but may overlook objects that do not produce thermal radiation. On the other hand, visible images struggle to recognize humans due to reflective lights from vehicles and lamps. After fusing the images from two modalities, by combining the advantages of two types of features, both humans and vehicles are well detected in the fusion image.

IV. CONCLUSION

In conclusion, DAE-Fuse overcomes limitations in image fusion, producing sharp and natural images through adversarial feature extraction and attention-guided fusion. Discriminative blocks in both phases enhance feature extraction and structural fidelity. Public dataset experiments demonstrate DAE-Fuse's superiority over existing methods.

V. ACKNOWLEDGEMENT

Our work was supported in part by the Guangdong Provincial Key Laboratory of Interdisciplinary Research and Application for Data Science, BNU-HKBU United International College, project code 2022B1212010006, and in part by Guangdong Higher Education Upgrading Plan (2021-2025) of "Rushing to the Top, Making Up Shortcomings and Strengthening Special Features" with UIC research grant UICR0400006-24.

REFERENCES

- [1] J. Liu, X. Fan, Z. Huang, G. Wu, R. Liu, W. Zhong, and Z. Luo, "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*, 2022, pp. 5802–5811.
- [2] Y. Sun, B. Cao, P. Zhu, and Q. Hu, "Detfusion: A detection-driven infrared and visible image fusion network," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 4003–4011.
- [3] W. Zhao, S. Xie, F. Zhao, Y. He, and H. Lu, "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)*, June 2023, pp. 13 955–13 965.
- [4] J. Liu, Z. Liu, G. Wu, L. Ma, R. Liu, W. Zhong, Z. Luo, and X. Fan, "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*, 2023, pp. 8115–8124.
- [5] Z. Liu, J. Liu, G. Wu, Z. Chen, X. Fan, and R. Liu, "Searching a compact architecture for robust multi-exposure image fusion," *IEEE TCSVT*, 2024.
- [6] T. Liu, K.-M. Lam, R. Zhao, and G. Qiu, "Deep cross-modal representation learning and distillation for illumination-invariant pedestrian detection," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 1, pp. 315–329, 2021.
- [7] Y. Sun, B. Cao, P. Zhu, and Q. Hu, "Drone-based rgb-infrared cross-modality vehicle detection via uncertainty-aware learning," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 10, pp. 6700–6713, 2022.
- [8] Z. Wang, D. Ziou, C. Armenakis, D. Li, and Q. Li, "A comparative analysis of image fusion methods," *IEEE transactions on geoscience and remote sensing*, vol. 43, no. 6, pp. 1391–1402, 2005.
- [9] Y. Zhang, "Understanding image fusion," *Photogramm. Eng. Remote Sens*, vol. 70, no. 6, pp. 657–661, 2004.
- [10] A. P. James and B. V. Dasarathy, "Medical image fusion: A survey of the state of the art," *Information fusion*, vol. 19, pp. 4–19, 2014.
- [11] J. Ma, W. Yu, P. Liang, C. Li, and J. Jiang, "Fusiongan: A generative adversarial network for infrared and visible image fusion," *Information fusion*, vol. 48, pp. 11–26, 2019.
- [12] H. Zhou, W. Wu, Y. Zhang, J. Ma, and H. Ling, "Semantic-supervised infrared and visible image fusion via a dual-discriminator generative adversarial network," *IEEE Transactions on Multimedia*, vol. 25, pp. 635–648, 2021.
- [13] Y. Rao, D. Wu, M. Han, T. Wang, Y. Yang, T. Lei, C. Zhou, H. Bai, and L. Xing, "At-gan: A generative adversarial network with attention and transition for infrared and visible image fusion," *Information Fusion*, vol. 92, pp. 336–349, 2023.
- [14] J. Ma, H. Xu, J. Jiang, X. Mei, and X.-P. Zhang, "Ddagan: A dual-discriminator conditional generative adversarial network for multi-resolution image fusion," vol. 29. IEEE, 2020, pp. 4980–4995.
- [15] P. Li, "Didfuse: deep image decomposition for infrared and visible image fusion," in *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, 2021, pp. 976–976.
- [16] Z. Zhao, H. Bai, J. Zhang, Y. Zhang, S. Xu, Z. Lin, R. Timofte, and L. V. Gool, "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)*, Jun. 2023, pp. 5906–5916.
- [17] H. Li and X.-J. Wu, "Densefuse: A fusion approach to infrared and visible images," *IEEE Transactions on Image Processing*, vol. 28, no. 5, pp. 2614–2623, 2018.
- [18] H. Li, X.-J. Wu, and T. Durrani, "Nestfuse: An infrared and visible image fusion architecture based on nest connection and spatial/channel attention models," *IEEE Transactions on Instrumentation and Measurement*, vol. 69, no. 12, pp. 9645–9656, 2020.
- [19] H. Li, X.-J. Wu, and J. Kittler, "Rfn-nest: An end-to-end residual fusion network for infrared and visible images," vol. 73. Elsevier, March 2021, pp. 72–86.
- [20] N. Park and S. Kim, "How do vision transformers work?" *arXiv preprint arXiv:2202.06709*, 2022.
- [21] Z. Pan, J. Cai, and B. Zhuang, "Fast vision transformers with hilo attention," *Advances in Neural Information Processing Systems*, vol. 35, pp. 14 541–14 554, 2022.
- [22] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [23] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [24] S. W. Zamir, A. Arora, S. Khan, M. Hayat, F. S. Khan, and M.-H. Yang, "Restormer: Efficient transformer for high-resolution image restoration," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 5728–5739.
- [25] L. Tang, J. Yuan, and J. Ma, "Image fusion in the loop of high-level vision tasks: A semantic-aware real-time infrared and visible image fusion network," *Information Fusion*, vol. 82, pp. 28–42, 2022.
- [26] L. Tang, J. Yuan, H. Zhang, X. Jiang, and J. Ma, "Piafusion: A progressive infrared and visible image fusion network based on illumination aware," *Information Fusion*, vol. 83-84, pp. 79–92, 2022.
- [27] A. Toet and M. A. Hogervorst, "Progress in color night vision," *Optical Engineering*, vol. 51, no. 1, pp. 010 901–010 901, 2012.
- [28] H. Xu, J. Ma, Z. Le, J. Jiang, and X. Guo, "Fusiondn: A unified densely connected network for image fusion," in *Proceedings of the Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI)*, 2020, pp. 12 484–12 491.
- [29] <http://www.med.harvard.edu/AANLIB/home.html>, Harvard Medical Website.
- [30] J. W. Roberts, J. A. Van Aardt, and F. B. Ahmed, "Assessment of image fusion procedures using entropy, image quality, and multispectral classification," *Journal of Applied Remote Sensing*, vol. 2, no. 1, p. 023522, 2008.
- [31] A. M. Eskicioglu and P. S. Fisher, "Image quality measures and their performance," *IEEE Transactions on communications*, vol. 43, no. 12, pp. 2959–2965, 1995.
- [32] J. Ma, Y. Ma, and C. Li, "Infrared and visible image fusion methods and applications: A survey," *Information fusion*, vol. 45, pp. 153–178, 2019.
- [33] Y. Han, Y. Cai, Y. Cao, and X. Xu, "A new image fusion performance metric based on visual information fidelity," *Information fusion*, vol. 14, no. 2, pp. 127–135, 2013.
- [34] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [35] G. Jocher, A. Chaurasia, and J. Qiu, "Ultralytics YOLOv8, version 8.0.0," <https://github.com/ultralytics/yolov8>, 2023.
- [36] H. Xu, J. Ma, J. Jiang, X. Guo, and H. Ling, "U2fusion: A unified unsupervised image fusion network." IEEE, 2020.
- [37] Z. Zhao, H. Bai, Y. Zhu, J. Zhang, S. Xu, Y. Zhang, K. Zhang, D. Meng, R. Timofte, and L. Van Gool, "Ddfm: Denoising diffusion model for multi-modality image fusion," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023, pp. 8082–8093.