

基于语义损失的红外与可见光图像融合算法

丁华彬¹, 丁麒文²

(1. 吉林化工学院 信息与控制工程学院, 吉林 吉林 132000; 2. 华北光电技术研究所, 北京 100015)

摘要: 提出了一种基于语义损失的红外与可见光图像融合算法, 通过语义损失引导生成图像包含更多语义信息, 满足高级视觉任务需求。首先使用预训练的分割网络对融合图像进行分割, 分割结果与标签图构成语义损失, 在语义损失和内容损失的共同引导下, 迫使融合网络在保证融合图像质量的前提下同时兼顾图像语义信息量, 融合图像满足高级视觉任务需求。同时本文还设计了一种新的特征提取模块, 通过残差密集连接实现特征重用, 提高细节描述能力, 进一步减轻融合框架, 从而提高图像融合的时间效率。实验结果表明, 本文算法在主观视觉效果和定量指标方面优于现有融合算法, 且融合图像包含更丰富的语义信息。

关键词: 图像融合; 语义损失; 神经网络; 红外图像; 可见光图像

中图分类号: TP391.41 文献标志码: A 文章编号: 1001-8891(2023)09-0941-07

Fusion Algorithm of Infrared and Visible Images Based on Semantic Loss

DING Huabin¹, DING Qiwen²

(1. School of Information and Control Engineering, Jilin Institute of Chemical Technology, Jilin 132000, China;

2. North China Research Institute of Electric-optics, Beijing 100015, China)

Abstract: In this study, we propose an infrared and visible image fusion algorithm based on semantic loss, to ensure that the generated images contain more semantic information through semantic loss, thereby satisfying the requirements of advanced vision tasks. First, a pre-trained segmentation network is used to segment the fused image, with the segmentation result and label map determining the semantic loss. Under the joint guidance of semantic and content losses, we force the fusion network to guarantee the quality of the fused image by considering the amount of semantic information in the image, to ensure that the fused image meets the requirements of advanced vision tasks. In addition, a new feature extraction module is designed in this study to achieve feature reuse through a residual dense connection to improve detail description capability while further reducing the fusion framework, which improves the time efficiency of image fusion. The experimental results show that the proposed algorithm outperforms existing fusion algorithms in terms of subjective visual effects and quantitative metrics and that the fused images contain richer semantic information.

Key words: image fusion, semantic loss, neural network, infrared image, visible image

0 引言

由于硬件设备的理论和技术限制, 在单个传感器或单个拍摄设置下拍摄的图像不能有效、全面地描述成像场景。图像融合能够将不同源图像中有意义的信息组合起来生成一幅包含更丰富的信息, 更有利于后续应用的图像^[1]。由于融合图像具有优良的特性, 图像融合作为一种图像增强方法在摄影可视化、目标跟踪^[2]、医学诊断^[3]、遥感监测等领域得到广泛应用。

在深度学习的普及之前, 图像融合就已经得到了深入的研究。早期实现图像融合的方法主要采用相关的数学变换来对图像进行分解, 并根据特定的融合规则进行融合。典型的传统融合方法包括基于多尺度变换的方法、基于稀疏表示的方法、基于子空间的方法、基于显著性的方法等。然而, 这些方法的局限性也变得越来越明显。一方面, 为了保证后续特征融合的可行性, 传统的方法被迫对不同的源图像采用相同的变换来提取特征。但是, 该操作没有考虑到源图像的特

收稿日期: 2022-08-18; 修订日期: 2022-11-24.

作者简介: 丁华彬 (1998-), 男, 山东青岛人, 硕士研究生, 研究方向为图像处理, 深度学习. E-mail: dhh41416@163.com

征差异，这可能会导致所提取的特征的表达性较差。另一方面，传统的特征融合策略过于粗糙，使得融合性能非常有限。将深度学习引入图像融合就是为了克服传统方法的这些局限性。在基于深度学习的方法中，主要有基于自动编码器的框架、基于卷积神经网络的框架和基于生成对抗网络的框架3种主体框架。

虽然最近的基于深度学习的图像融合算法可以在相应的融合任务中能够取得更好的效果，但仍存在一些影响图像融合性能的缺陷。一方面，现有的融合算法倾向于追求更好的视觉质量和更高的评价指标，但很少系统地考虑融合后图像是否可以满足高级的视觉任务需求。研究表明，仅考虑视觉质量和定量指标并不能帮助用于高级视觉任务。虽然有些研究通过分割图像掩模指导图像融合过程，但该掩模只分割了部分显著目标，在增强语义信息方面效果有限。其次，在现有的融合方法中，大多只使用一个特征提取层的输出作为图像融合组件。这种架构直接造成了源图像中的信息丢失，无法更好地提取特征细节，直接影响了最终的融合性能。第三，现有的融合方法由于计算复杂度和参数量大，在运行时间和存储空间方面通常缺乏竞争力。

为了克服上述挑战，本文提出了一种基于语义损失引导的融合网络，来执行红外和可见光图像融合任务。该方法在保证图像融合质量的前提下同时兼顾了高级视觉任务，融合后图像满足高级视觉任务需求。具体来说，引入图像分割网络^[4]（Fast Semantic Segmentation Network, Fast-SCNN）来对融合图像进行分割，分割结果与label图构成语义损失，融合网络在语义损失和内容损失的共同指导下，生成融合图像。此外，本文还设计了一种新的特征提取模块，通过残差密集连接实现特征重用，提高细节描述能力，同时以进一步减轻融合框架，从而提高图像融合的时间效率。

1 本文融合算法

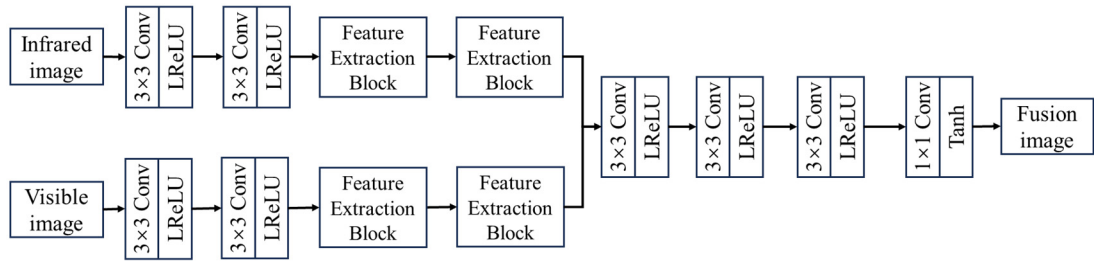


图2 融合网络架构

Fig.2 Fusion network architecture

1.1 整体框架

给定一对预配准完成的红外图像和可见光图像，在损失函数的引导下，通过特征提取、集成和重建实现图像融合。融合后图像的质量在很大程度上取决于损失函数。为了提升融合性能，本文设计了一个由内容损失和语义损失组成的联合损失函数来约束融合网络，该算法的整体框架如图1所示。

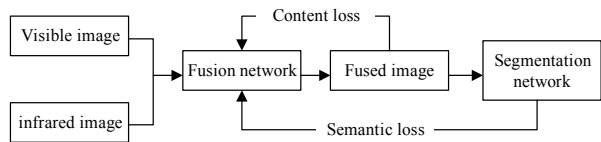


图1 本文融合算法整体框架

Fig.1 The overall framework of the proposed image fusion algorithm

首先，为了更好地从红外和可见光图像中提取具有丰富细节信息的深度特征，本文设计了一种基于Sobel算子残差密集块的图像融合网络，通过特征重用将源图像中的互补信息充分集成。然后，通过特征集成和图像重建模块对融合后的图像进行重建，融合图像集成了红外和可见光特征，包含了丰富的细节信息。

此外，为了满足高级视觉任务需求，本文引入了语义损失来约束融合网络，使融合后图像包含更丰富的语义信息。该算法使用预训练完成的Fast-SCNN分割模型来对融合后图像进行分割，分割结果与源图像语义标签之间的差异就可以反映融合后图像所包含语义信息的丰富度。

1.2 融合网络

为了轻量化融合网络，更好地提取深度信息，本文提出了一种基于Sobel算子残差密集块的图像融合网络，如图2所示。该融合网络由特征提取器和图像重建器组成，通过特征提取、集成和图像重建实现图像融合。

如图2所示,特征提取器包含两个并行的红外和可见光的特征提取流,每个特征提取流包含两个共同的卷积层和两个残差密集块。采用公共卷积层,卷积核大小为 3×3 ,激活函数为漏式校正直线性单元(LReLU),提取浅层特征。接下来是两个残差密集块,用于提取细节信息。残差密集块的具体设计如图3所示。每个残差密集块包含3个分支,以提高所提取的特征的多样性,使其能够充分利用模块内每个卷积层所提取的深度特征。每个分支的 1×1 卷积层产生低维特征映射,用于降低输入特征的维数,从而降低模块中内部特征的卷积计算成本。

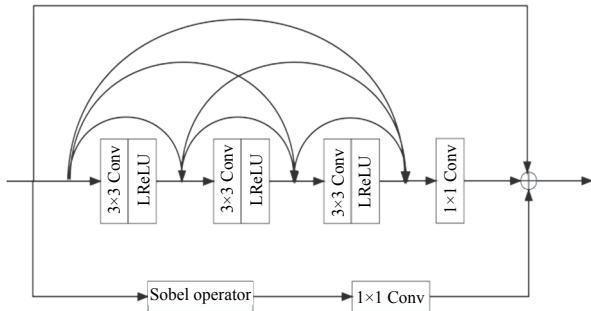


图3 特征提取器具体设计

Fig.3 Specific design of the feature extractor

然后,通过concat策略将红外图像和可见光图像的细粒度特征进行集成,并将结果输入到图像重构器中,以实现特征聚合和图像重建。图像重构器由3个串联的 3×3 卷积层和一个 1×1 卷积层组成。所有 3×3 卷积层均采用LReLU作为激活函数, 1×1 卷积层的激活函数为Tanh。为了不造成融合过程中的信息丢失,该网络没有引入降采样,图像边缘填充设置为0,步幅设置为1,融合后图像与源图像大小一致。

1.3 分割网络

分割网络采用预训练完成的快速分割卷积神经网络(Fast-SCNN),该网络在现有的双分支快速分割方法的基础上,引入了“学习降采样”模块,该模块同时计算多个分辨率分支的低水平特征。该网络结合了高分辨率的空间细节和低分辨率提取的深度特征,在Cityscapes数据集^[5]上分割速度可以达到123.5帧,准确率为68.0%。

1.4 损失函数

本文融合网络旨在加强融合图像中的语义信息,同时提高视觉质量和评价指标。为了实现这些目标,本文从两个角度设计了损失函数,对融合网络进行约束。一方面,融合网络需要充分整合源图像中的互补信息,如红外图像中的突出目标和可见光图像中的纹理细节。为此,本文设计内容损失来提升融合图像的视觉质量。另一方面,融合后的图像包含更丰富的语

义信息,可以有效地促进高级视觉任务。为此,构建了一个语义损失来反映融合图像语义信息丰富程度。

1.4.1 内容损失

内容损失决定了提取的信息类型以及重建中各种类型信息之间的主要和次要关系。为最大限度地减少信息损失,以有效地保存红外图像的热辐射信息和可见图像的纹理细节信息。内容损失由两部分组成,定义如下:

$$L_{\text{content}} = aL_{\text{structure}} + L_{\text{intensity}} \quad (1)$$

式中: $L_{\text{structure}}$ 项表示两幅图像的结构相似性,保证融合图像包含丰富的纹理细节; $L_{\text{intensity}}$ 约束融合的图像,以保持与源图像相似的强度分布; a 为权衡结构损失和强度损失的正参数。

强度损失在像素水平上测量融合图像和源图像之间的差异。因此,将红外图像和可见光图像的强度损失定义为:

$$L_{\text{intensity}} = \frac{\delta \cdot \|I_f - I_r\|_F^2 + (1 - \delta) \|I_f - I_v\|_F^2}{H \cdot W} \quad (2)$$

式中: I_f 为融合图像; I_r , I_v 分别为红外图像和可见光图像; H 和 W 分别表示输入图像的高度和宽度; $\|\cdot\|_F$ 代表矩阵弗罗贝尼乌斯范数; δ 是一个控制两项之间的权衡的正参数。

在像素级上,结构相似度^[6](Structural Similarity Index Measure, SSIM)指数度量是最流行和最有效的度量方法,它根据亮度、对比度和结构信息的相似性来模拟失真。因此,我们选择它来约束输入图和输出图像之间的结构相似性。 $L_{\text{structure}}$ 定义如下:

$$L_{\text{structure}} = \frac{\lambda \cdot (1 - \text{SSIM}(I_f, I_r)) + (1 - \lambda) (1 - \text{SSIM}(I_f, I_v))}{H \cdot W} \quad (3)$$

式中: $\text{SSIM}(I_x, I_y)$ 表示 I_x 和 I_y 之间的结构相似性; λ 为控制两项之间的权衡的正参数。

1.4.2 语义损失

为了更好地提高融合图像中的语义信息,本文设计了一个语义损失来约束融合网络。该算法使用预训练好的Fast-SCNN分割模型来对融合后图像进行分割,分割网络输出分割结果 I_s 。通过计算分割结果与Label图之间交叉熵可得到语义损失定义如下:

$$L_{\text{semantic}} = - \sum_{h=1}^H \sum_{w=1}^W y_c \log(I_s) \quad (4)$$

式中: y_c 为分割标签进行one-hot编码转换而来。

最后,构造一个联合损失来引导融合网络的训练,该联合损失定义为:

$$L=L_{\text{content}}+\beta L_{\text{semantic}} \quad (5)$$

式中： β 是一个表征语义损失 L_{semantic} 重要性的常数。 β 越大，损失函数中语义损失所占比重越大，融合后图像语义信息更加接近源图像。

1.5 训练细节

本文使用 MFNet 数据集^[7]来训练融合网络，MFNet 数据集包含 1569 对图像对（820 张在白天拍摄，749 张在夜间拍摄），空间分辨率为 480×640 。MFNet 数据集为 8 个对象提供了语义标签，即汽车、人、自行车、曲线、汽车停车站、护栏、警戒锥和凸起。为增强数据可信度，所有的数据集图像在被输入网络之前都被归一化到 $[0,1]$ 范围内。

训练中的所有参数设置如下：内容损失超参数 $\alpha=10$ ，语义损失权重参数 $\beta=3$ 。该算法使用 Adam 优化器，bathsize 大小为 8，epoch 大小为 10，初始学习率 0.001，平滑参数 $\beta_1=0.9$ ， $\beta_2=0.99$ ， $\varepsilon=1e^{-8}$ ，权重衰减为 0.0002，在损失函数的指导下优化我们的融合模型。该算法在 PyTorch 平台^[8]上实现。

由于 MFNet 数据集包含三通道彩色可见光图像，为更好地处理彩色图像，本文首先将可见图像转换为 YCbCr 颜色空间。然后，利用不同的融合算法来融合可见光图像和红外图像的 Y 通道。最后，利用可见光

图像的 Cb 和 Cr 通道将融合后的图像转换为 RGB 颜色空间^[9]。

2 实验与分析

为了全面评估所提出的算法，我们在 MFNet 数据集上进行了广泛的定性和定量实验。我们选择 3 种较为先进的融合算法，即 DDcGAN^[10]（Dual-Discriminator Conditional Generative Adversarial Network），Nestfuse^[11]，GANMcM^[12]（Generative Adversarial Network with Multiclassification Constraints）与本文算法在主观和客观两个方面对融合结果进行比较，然后再对其融合图像分割性能进行对比。所有这 3 种方法的实现都是公开的，参数均按照原论文中默认参数设置。

2.1 主观评价

主观评价方法以人体视觉系统为基础，对融合图像的质量进行评价，根据图像细节、物体完整性、图像失真等标准一致地比较不同的融合方法，在融合质量评价过程中能够更直接，更可靠。本文从 MFNet 数据集中随机挑选 3 张分别为白天、黑夜和黄昏的图片进行实验，实验对比结果如图 4，5，6。

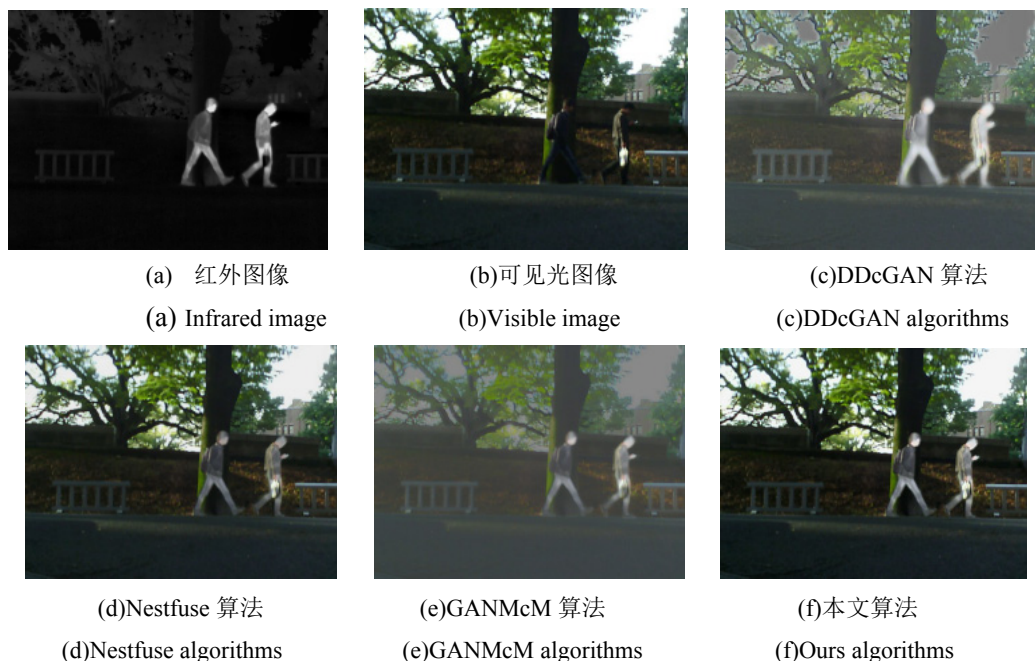


图4 “00057D”图像不同算法的融合结果

Fig.4 Fusion results of different algorithms for 00057D image



图5 “00510D”图像不同算法的融合结果

Fig.5 Fusion results of different algorithms for 00510D image

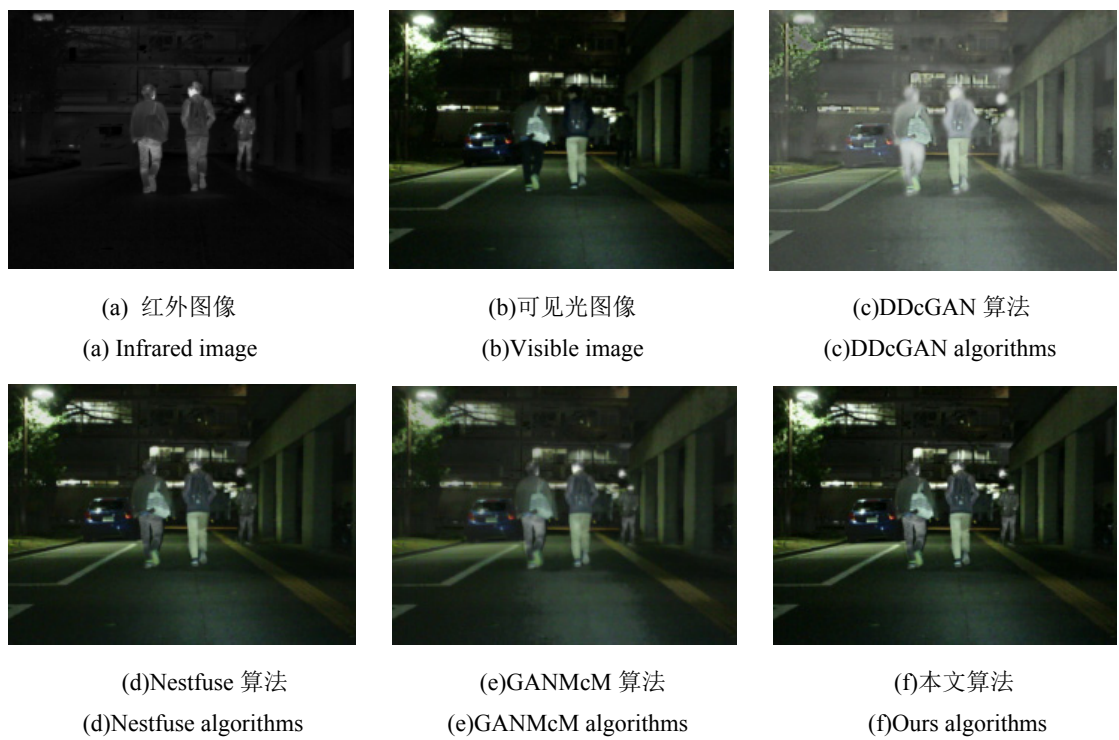


图6 “01347N”图像不同算法的融合结果

Fig.6 Fusion results of different algorithms for 01347N image

在日间场景中,可以利用红外图像的热辐射信息作为可见光图像的补充信息。因此,具有良好视觉质量的融合图像应包含丰富的可见光图像的纹理细节,并增强红外图像中的突出目标。从实验结果可以看出DDcGAN 和 GANMcM 算法并没有有效保留可见光的亮度信息,导致融合后图像偏暗,Nestfuse 算法虽

然保留了可见光图像中的亮度信息,但没有更好地保留红外图像中的热辐射信息,融合后图像目标显著性不强。

在夜间场景中,由于环境限制,红外和可见光图像都只能提供有限的场景信息。因此,能否自适应地整合红外和可见光图像中可用信息将成为评判融合

算法的标准之一。实验结果可以看出,所有算法都在一定程度上合并了红外和可见光图像中的互补信息,但不同算法的融合结果仍有一些细微的差异。DDcGAN 算法虽然保留了目标的热辐射信息,但目标轮廓不明显。GANMcM 算法在背景中丢失了部分纹理信息,在融合后图像中引入了一些无用信息。实验表明,本文算法能够更好地融合红外和可见光图像中的有用信息,具有良好的视觉效果,从融合后图像的背景中可以看出,本文算法较其它算法能够更好地提取纹理细节,并避免受到无用信息的干扰。

2.2 客观评价

由于主观评价方法存在人工干预、时间消耗较长、成本高、不可重复性等缺点,为保证评价的全面性和可信度,需要同时使用能够定量、自动度量融合图像质量的客观评价方法。本文采用了4种具有代表性的图像质量评价指标,分别为熵(Entropy, EN)、互信息(Mutual Information, MI)、标准差(Standard Deviation, SD)、空间频率(Spatial Frequency, SF)和边缘相似度($Q^{AB/F}$)。

1) 熵

熵表示在融合图像中所包含的信息量,定义如下:

$$EN = -\sum_{l=0}^{L-1} p_l \log_2 p_l \quad (6)$$

2) 互信息

互信息度量通过测量两幅图像间依赖性,度量从源图像传输到融合图像的信息量,定义如下:

$$MI = MI_{A,F} + MI_{B,F} \quad (7)$$

式中: $MI_{A,F}$ 和 $MI_{B,F}$ 分别表示从红外图像和可见光图像传输到融合图像的信息量,定义如下:

$$MI_{X,F} = \sum_{x,f} p_{X,F}(x,f) \log \frac{p_{X,F}(x,f)}{p_X(x)p_F(f)} \quad (8)$$

3) 空间频率

空间频率表示了图像的纹理信息,数学表达式如下:

$$SF = \sqrt{RF^2 + CF^2} \quad (9)$$

4) 标准差

标准差反映了图像的对比度特性,数学定义如下:

$$SD = \sqrt{\sum_{i=1}^M \sum_{j=1}^N [F(i,j) - \mu]^2} \quad (10)$$

5) 边缘相似度

边缘相似度($Q^{AB/F}$)测量从源图像传输到融合图像的边缘信息的量,定义如下:

$$Q^{AB/F} = \frac{\sum_{i=1}^N \sum_{j=1}^M Q^{AF}(i,j)w^A(i,j) + Q^{BF}(i,j)w^B(i,j)}{\sum_{i=1}^N \sum_{j=1}^M [w^A(i,j) + w^B(i,j)]} \quad (11)$$

为保证实验可靠性,本文在 MFNet 数据集中随机挑选 20 张图片进行实验,获取 20 张图片性能指标的平均值,实验数据如表 1,其中加粗的值表示所有方法中的性能最佳的方法。从表中可以看出,本文算法在 4 个评价指标 EN、SF、 $Q^{AB/F}$ 和 MI 值上的表现最佳。虽然该算法的 SD 不是最大的,但结果与其他算法相差不大。实验结果表明,本文算法在图像融合评价指标上优于其它算法。

表 1 图像融合客观指标结果

Table 1 Objective index results of fusion image					
	MI	SD	SF	EN	$Q^{AB/F}$
DDcGAN	1.9848	6.2586	0.0231	5.2684	0.1859
Nestfuse	2.9764	7.5843	0.0301	6.2386	0.4862
GANMcM	2.5130	8.2627	0.0236	6.1895	0.3238
Ours	3.8672	7.9858	0.0458	6.5841	0.5842

2.3 分割性能评价

为了验证融合图像的语义分割性能,本文选取 MFNet 数据集中 20 张图片经过不同融合算法生成融合图片,然后将融合后图片进行语义分割后获取 mIoU (Mean Intersection Over Union) 值来衡量分割性能,实验结果如表 2。实验结果显示,本文算法相较于其它算法有更高的 mIoU,这表明本文算法生成的融合图像中包含了丰富的语义信息,使分割网络能够更准确地描述成像场景。

表 2 融合图像分割性能指标 (mIoU)

Table 2 mIoU of fusion image				
	DDcGAN	Nestfuse	GANMcM	Ours
mIoU	75.33	76.32	75.68	77.98

此外,我们还提供了一些可视化的例子来显示红外、可见光 and 不同融合图像的分割结果,如图 7。从结果中可以看出,红外图像提供了更多关于行人等突出目标的信息,而可见图像可以更好地描述背景。本文融合算法可以集成来自源图像的互补信息,并实现一个更全面的对成像场景的描述。

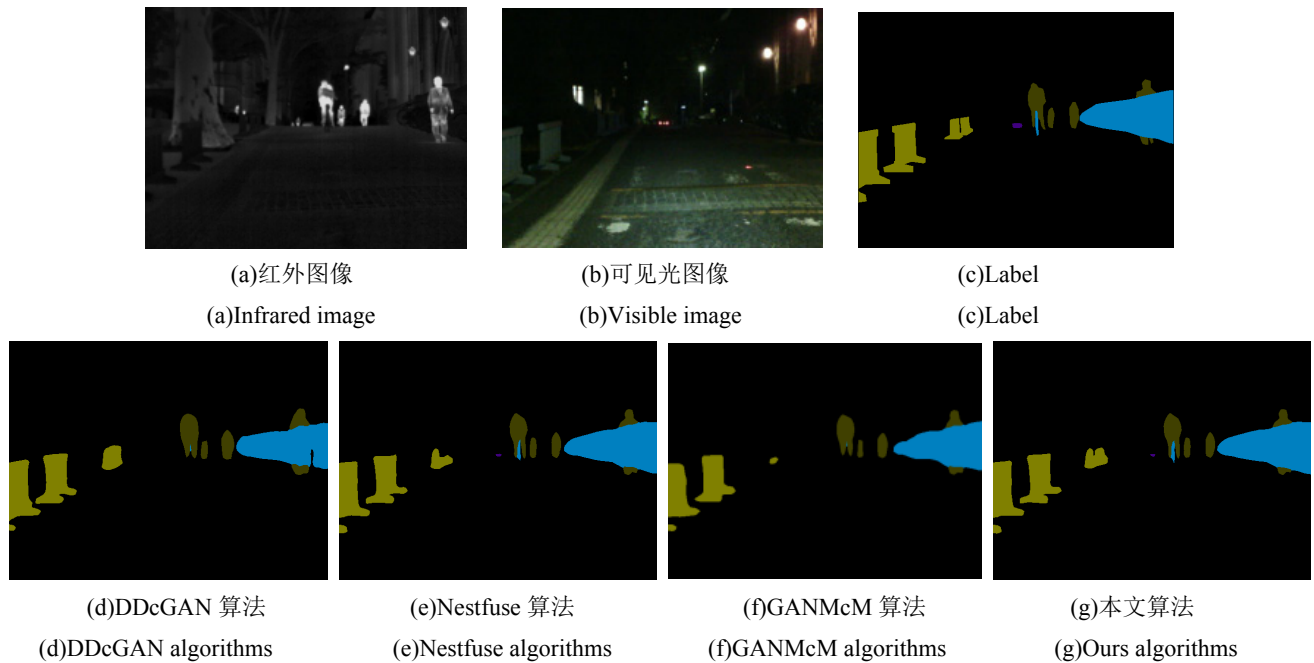


图7 不同算法融合图分割可视化图

Fig.7 Different algorithms merge the graph segmentation visualization diagram

3 结论

本文提出了一种基于语义损失的红外与可见光图像融合算法,引入语义损失,使融合网络在语义损失和内容损失的共同引导下,生成富含语义信息的高质量图像,融合图像满足高级视觉任务需求。同时,本文设计了一种新的特征提取模块,通过残差密集连接实现特征重用,提高细节描述能力,实现了显著的目标强度和纹理细节的有效保留。通过比较和泛化实验表明,本文融合算法在主观和定量指标方面都优于其它先进算法,且相较于其它算法融合后图像语义信息含量最高。

参考文献:

- [1] MA J, MA Y, LI C. Infrared and visible image fusion methods and applications: a survey[J]. *Information Fusion*, 2019, **45**: 153-178.
- [2] ZHU Y, LI C, LUO B, et al. Dense feature aggregation and pruning for RGBT tracking[C]// *The 27th ACM International Conference. ACM*, 2019: 465-472.
- [3] Bhatnagar G, WU Q, ZHENG L. Directive contrast based multimodal medical image fusion in nset domain[J]. *IEEE Transactions on Multimedia*, 2014, **9**(5): 1014-1024.
- [4] Poudel R, Liwicki S, Cipolla R. Fast-SCNN: fast semantic segmentation network[J/OL]. arXiv Preprint arXiv:1902.04502, 2019. <https://doi.org/10.48550/arXiv.1902.04502>
- [5] Cordts M, Omran M, Ramos S, et al. The cityscapes dataset for semantic urban scene understanding[C]// *Conference on Computer Vision and Pattern Recognition (CVPR). IEEE*, 2016: 3213-3223.
- [6] ZHOU W, Bovik A C, Sheikh H R, et al. Image quality assessment: from error visibility to structural similarity[J]. *IEEE Trans Image Process*, 2004, **13**(4): 600-612.
- [7] HA Q, Watanabe K, Karasawa T, et al. MFNet: Towards real-time semantic segmentation for autonomous vehicles with multi-spectral scenes[C]// *IEEE/RSSJ International Conference on Intelligent Robots and Systems (IROS). IEEE*, 2017: 4714-4722.
- [8] Paszke A, Gross S, Massa F, et al. PyTorch: an imperative style, high-performance deep learning library[J/OL]. *Advances in Neural Information Processing Systems*, 2019, <https://doi.org/10.48550/arXiv.1912.01703>.
- [9] Prabhakar K R, Srikar V S, Babu R V. DeepFuse: A deep unsupervised approach for exposure fusion with extreme exposure image pairs[C]// *IEEE International Conference on Computer Vision (ICCV). IEEE Computer Society*, 2017: 4714-4722.
- [10] MA J, XU H, JIANG J, et al. DDcGAN: a dual-discriminator conditional generative adversarial network for multi-resolution image fusion[J]. *IEEE Transactions on Image Processing*, 2020, **29**: 4980-4995.
- [11] LI H, WU X J, Durrani T. NestFuse: an infrared and visible image fusion architecture based on nest connection and spatial/channel attention models[J]. *IEEE Transactions on Instrumentation and Measurement*, 2020(99): 1-1.
- [12] MA J, ZHANG H, SHAO Z, et al. GANMcC: a generative adversarial network with multiclassification constraints for infrared and visible image fusion[J]. *IEEE Transactions on Instrumentation and Measurement*, 2021, **70**: 1-14, Doi: 10.1109/TIM.2020.3038013.