

文章编号: 1671-7449(2024)06-0575-11

去伪影空洞残差与双向注意力金字塔增强的 图像级联语义分割算法

王国刚, 罗香莲

(山西大学 物理电子工程学院, 山西 太原 030006)

摘要: 针对ICNet空洞卷积引入网格伪影, 金字塔池化模块难以充分融合不同尺度特征层间的上下文信息, 同一尺度特征和跨尺度特征均交互不足的问题, 提出去伪影空洞残差与双向注意力金字塔增强的图像级联语义分割算法。该算法提出去伪影空洞残差网络, 消除空洞卷积引入的网格伪影, 加强网络的表征能力; 构建双向金字塔池化模块, 沿自顶向下的路径将语义信息传递到低层, 并沿自底向上的路径将位置信息传递到高层; 构造注意力金字塔增强模块, 创建自注意力、向上注意力、向下注意力机制, 建立跨空间与跨尺度的特征交互。Cityscapes数据集上的实验结果表明, 所提算法的总体精确率、平均像素精确度、平均交并比分别比ICNet提高了0.67%、5.58%、4.40%, 而且该算法的以上3个评价指标均优于其它8种主流比对算法。

关键词: 深度学习; 语义分割; ICNet; 双向金字塔; 注意力金字塔增强

中图分类号: TP391.4 **文献标识码:** A **doi:** 10.3969/j.issn.1671-7449.2024088

引用格式: 王国刚, 罗香莲. 去伪影空洞残差与双向注意力金字塔增强的图像级联语义分割算法[J]. 测试技术学报, 2024, 38(6): 575-585.

WANG Guogang, LUO Xianglian. Image cascade semantic segmentation algorithm based on deartifacting dilated residual and bidirectional attention pyramid enhancement[J]. Journal of Test and Measurement Technology, 2024, 38(6): 575-585.

Image Cascade Semantic Segmentation Algorithm Based on Deartifacting Dilated Residual and Bidirectional Attention Pyramid Enhancement

WANG Guogang, LUO Xianglian

(College of Physics and Electronic Engineering, Shanxi University, Taiyuan 030006, China)

Abstract: A novel algorithm named as image cascade semantic segmentation algorithm based on deartifacting dilated residual and bidirectional attention pyramid enhancement is proposed to avoid the disadvantages of an introduction of gridding artifacts for dilated convolutions, the difficulty in fully integrating context information among different scale feature layers for the pyramid pooling module, and lacking sufficient interaction for the same scale features and the cross-scale features in ICNet algorithm. The characterization ability of the network is strengthened since the deartifacting dilated residual network is constructed to remove gridding artifacts

收稿日期: 2024-03-05

基金项目: 国家自然科学基金资助项目(11804209); 山西省自然科学基金资助项目(201901D111031, 201901D211173); 山西省高校科技创新计划(2019L0064, 2020L0051)

作者简介: 王国刚(1977—), 男, 副教授, 博士, 主要从事图像处理与视频分析研究。E-mail: kingguogang@sxu.edu.cn.

introduced by dilated convolutions. After building the bidirectional pyramid pooling module, the semantic information is transmitted to the lower level along a top-down path while the location information is conveyed to the higher level along a bottom-up path. Moreover, the feature interactions for both cross-space and cross-scale are established by the attention pyramid enhancement module, which introduces self-attention, upward attention, and downward attention mechanism. On the Cityscapes dataset, the experimental results show that the overall accuracy, the mean pixel accuracy and the mean intersection over union of the proposed algorithm have been improved by 0.67%, 5.58% and 4.40% respectively compared to ICNet. Compared with the state-of-the-art methods, the presented method achieves outperformance on the above evaluation indexes.

Key words: deep learning; semantic segmentation; ICNet; bidirectional pyramid; attention pyramid enhancement

0 引言

作为计算机视觉领域的一项关键技术,语义分割的目的是预测图像中每个像素点对应的语义类别标签。目前,语义分割在遥感测绘^[1-3]、无人驾驶^[4-7]、医疗影像分析^[8-10]等领域得到了广泛应用。

基于机器学习的图像分割方法主要包括基于边缘、阈值、区域、聚类和图论的分割方法。该方法通过人工经验提取图像特征,复杂场景下分割精度低,难以满足实际场景中多语义类别的需求。近年来,卷积神经网络在计算机视觉(Computer Vision, CV)领域表现出了强大的特征学习能力。2015年,Long等^[11]提出全卷积神经网络(Fully Convolutional Networks, FCN),该网络将深度学习引入图像语义分割领域,通过端到端的方式构建语义分割模型。

当前主流语义分割算法大多采用类似FCN的框架,大致可分为两类:一类算法使用编解码结构,通过编码器提取语义信息,再经解码器恢复空间分辨率,如SegNet^[12]和GCN(Global Convolutional Network)^[13]算法,该类算法一定程度上弥补了空间位置信息的损失,但缺乏足够的上下文信息,导致分割精度较低;另一类算法通过池化或卷积形式的金字塔获取多尺度上下文信息以提高分割准确率,如PSPNet^[14]和OCNet^[15]算法,该类算法往往模型复杂,推理速度慢。

ICNet^[16]通过低分辨率分支获得语义信息,并增加中、高分辨率分支以获取细节信息,这在保证分割性能的同时,提高了推理速度。然而,ICNet空洞卷积引入网格伪影,金字塔池化模块难以充分融合不同尺度特征层间的上下文信息,同

一尺度特征和跨尺度特征均交互不足。

针对以上问题,提出去伪影空洞残差与双向注意力金字塔增强的图像级联语义分割算法(Image Cascade Semantic Segmentation Algorithm Based on Deartifacting Dilated Residual and Bidirectional Attention Pyramid Enhancement, IDB)。该算法构建去伪影空洞残差网络,去除空洞卷积引入的网格伪影,提高网络的表征能力;构造双向金字塔池化模块,引入自顶而下和自底向上的路径以加强不同尺度特征层之间上下文信息的融合;提出注意力金字塔增强模块,创建自注意力、向上注意力、向下注意力机制,实现跨空间与跨尺度的特征交互。实验结果表明, IDB算法在Cityscapes数据集上的总体精确率、平均像素精确度、平均交并比均优于其它8种主流对比算法。

1 相关模型

作为一种典型的实时语义分割模型, ICNet在PSPNet的基础上采用低、中、高三分辨率分支策略来降低时间成本,其结构如图1所示。ICNet的低分辨率分支利用PSPNet获得粗略特征,中分辨率分支与高分辨率分支使用少量卷积层获取细节特征。金字塔池化模块(Pyramid Pooling Module, PPM)和级联特征融合模块(Cascade Feature Fusion, CFF)融合各支路特征,再经预测网络得到语义分割图。

2 IDB方法

2.1 去伪影空洞残差网络

ICNet在ResNet50的后两个模块res-4和res-5中分别连续采用扩张率为2和4的空洞卷积。设 F 为特征图, K 为卷积核, d 为扩张率, 则空洞卷

积可定义为

$$(F*K)(p) = \sum_{a+db=p} F(a)K(b), \quad (1)$$

式中： b 、 a 和 p 分别为卷积层、输入特征图和输出特征图中的点。

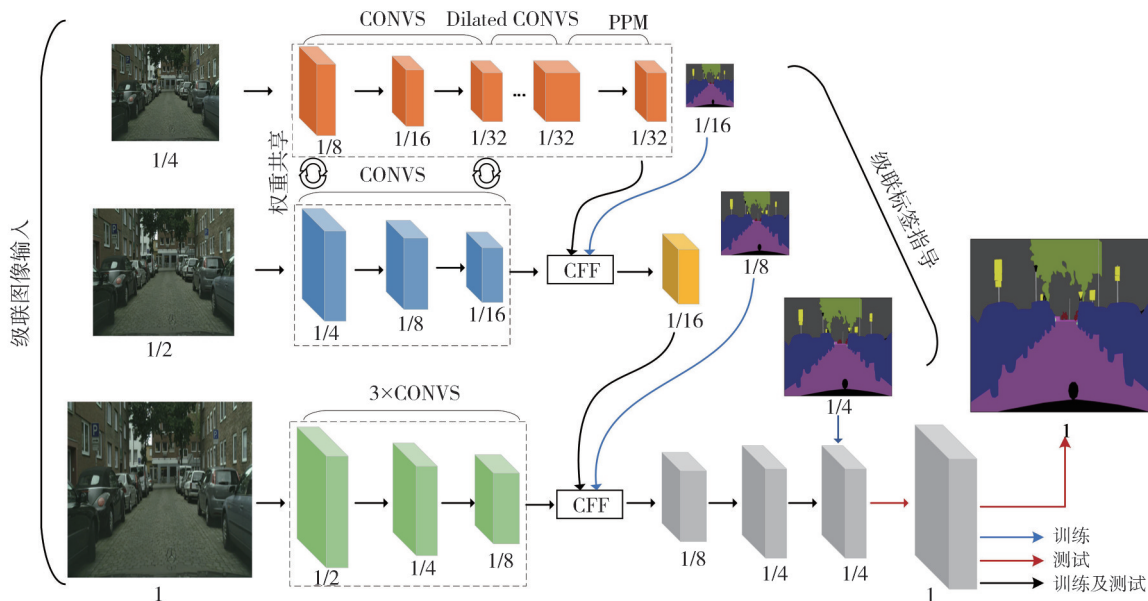


图1 ICNet网络架构

Fig. 1 ICNet network structure

空洞卷积在卷积核的相邻权值之间填充了零值,这扩充了特征图上卷积核的计算范围,从而提取了更大区域的特征信息。

虽然空洞卷积可以扩大感受野,但有可能会引入网格伪影,如图2所示。图2(a)绿色部分表示输入的单个像素,图2(b)表示扩张率为2的空洞卷积,图2(c)展示了网格伪影。ResNet50在7×7卷积层后采用最大池化层。最大池化层导致的高幅高频效应向后续特征层传播,这使网格伪影问题更严重。

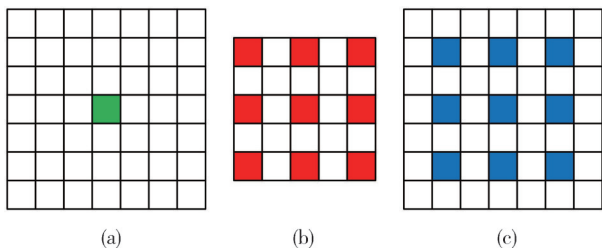


图2 网格伪影

Fig. 2 A gridding artifact

为消除网格效应,提高分割精度,IDB采用去伪影空洞残差网络(Dearfacting Dilated Residual Network, DDRN)提取特征。若 s 、 d 、 n 、 c 分别表示步长、扩张率、重复次数和输出通道数,则DDRN的网络结构可由表1表示。DDRN将最大池化层替换为两个3×3的卷积层以消除高幅高频

效应。此外,DDRN在网络末端添加扩张率为2和1的两个卷积层以滤除混叠伪影。

表1 去伪影空洞残差网络结构

Tab. 1 Dearfacting dilated residual network structure

层号	操作	c	n	d	s
1	Conv7×7	16	1	1	1
2	Conv3×3	16	1	1	1
3	Conv3×3	32	1	1	2
4	Bottleneck	256	3	1	2
5	Bottleneck	512	4	1	2
6	Bottleneck	1 024	6	2	1
7	Bottleneck	2 048	3	4	1
8	Conv3×3	512	1	2	1
9	Conv3×3	512	1	1	1

通过替换最大池化层和添加不同扩张率的卷积层,DDRN消除了网格伪影,加强了网络的表征能力,提升了IDB模型的分割性能。

2.2 双向金字塔池化

PPM拼接各池化特征图可获取不同区域的特征信息,但缺乏足够的轮廓细节,而且不同尺度特征图之间的上下文信息交互不足。针对此问题,IDB构建双向金字塔池化模块(Bidirectional Pyramid Pooling Module, BPPM),如图3所示。该模块的基本单元BPPM层,沿自顶向下的路径将语义信息传递到低层,并沿自底向上的路径将位置信息传递到高层。这使不同尺度特征间的语

义信息和位置信息得到了充分融合,提高了IDB 算法的分割准确率。

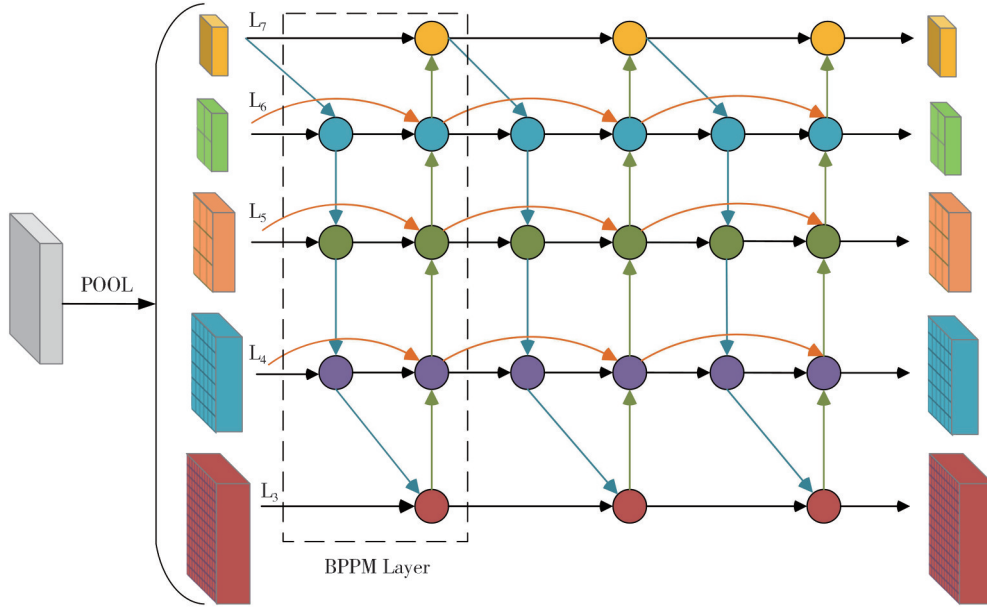


图3 双向金字塔池化模块

Fig.3 Bidirectional pyramid pool module

基于不同尺度子区域, BPPM对输入特征图进行全局平均池化得到 12×12 , 6×6 , 3×3 , 2×2 , 1×1 的池化特征图 L_3 、 L_4 、 L_5 、 L_6 、 L_7 。将 L_3 、 L_4 、 L_5 、 L_6 、 L_7 作为输入, 第一个BPPM层沿自上而下的路径得到中间特征层 L_i^{td} , 再自下而上逐层下采样并融合同层特征以得到输出特征层 L_j^{out} , 如式(2)~(6)所示。

$$L_6^{td} = Conv_{3 \times 3} \left(\frac{\omega_{6,in} \cdot L_6 + \omega_{7,in} \cdot Re\ size_{up}(L_7)}{\omega_{6,in} + \omega_{7,in} + \epsilon} \right), \quad (2)$$

$$L_j^{out} = Conv_{3 \times 3} \left(\frac{\omega'_{j,in} \cdot L_j + \omega'_{j,td} \cdot L_j^{td} + \omega'_{j-1,out} \cdot Re\ size_{down}(L_{j-1}^{out})}{\omega'_{j,in} + \omega'_{j,td} + \omega'_{j-1,out} + \epsilon} \right), \quad (5)$$

$$L_7^{out} = Conv_{3 \times 3} \left(\frac{\omega'_{7,in} \cdot L_7 + \omega'_{6,out} \cdot Re\ size_{down}(L_6^{out})}{\omega'_{7,in} + \omega'_{6,out} + \epsilon} \right), \quad (6)$$

式中: $\omega'_{j,in}$, $\omega'_{j,td}$, $\omega'_{j,out}$ 分别为第 j 输入层、中间层、输出层权重; $Re\ size_{down}(\cdot)$ 表示下采样; $j \in \{4, 5, 6\}$ 。

BPPM包含3个BPPM层。与首个BPPM层操作类似, BPPM依次经第2、3层后得到输出特征图。

IDB构建双向金字塔池化模块, 该模块引入池化特征图, 并双向加权以融合多尺度特征。这获取了更多轮廓细节, 增强了不同尺度特征图之间上下文信息的交互。

2.3 注意力金字塔增强模块

PPM通过自适应全局平均池化得到不同尺度的特征图, 但同一尺度特征和跨尺度特征均交互不足。针对此问题, 构建注意力金字塔增强模块

$$L_i^{td} = Conv_{3 \times 3} \left(\frac{\omega_{i,in} \cdot L_i + \omega_{i+1,td} \cdot Re\ size_{up}(L_{i+1}^{td})}{\omega_{i,in} + \omega_{i+1,td} + \epsilon} \right), \quad (3)$$

式中: $\omega_{i,in}$, $\omega_{i,td}$ 分别为第 i 输入层和中间层权重; $Re\ size_{up}(\cdot)$ 和 $Conv_{3 \times 3}(\cdot)$ 分别表示上采样和 3×3 卷积; ϵ 为0.0001; $i \in \{3, 4, 5\}$ 。

$$L_3^{out} = L_3^{td}, \quad (4)$$

(Attention Pyramid Enhancement, APE)。该模块创建自注意力(Self-Attention, SA)、向下注意力(Downward Attention, DA)和向上注意力(Upward Attention, UA)3种机制。SA捕捉同一尺度特征图上共现的目标特征; DA将高层特征图中的语义信息传递到低层; UA把低层特征图中的位置信息融入到高层。

2.3.1 自注意力机制

X' 、 X 分别为输入特征图和调整尺寸后的输入特征图, 设 $X' \in \mathbf{R}^{C \times H \times W}$ 、 $X \in \mathbf{R}^{C \times HW}$ 。若 $C' \times C$ 的可学习矩阵 W_q 、 W_k 、 W_v 分别表示查询、键、值的权重矩阵, 则查询、键、值矩阵可分别表示为 $Q = W_q X$ 、 $K = W_k X$ 、 $V = W_v X$ 。

令

$$\begin{aligned}
 \mathbf{Q} = (q_i) &= \begin{pmatrix} q_1 \\ \vdots \\ q_{\frac{C'}{N}} \\ \vdots \\ q_{\frac{C'}{N} \times (n-1) + 1} \\ \vdots \\ q_{\frac{C'}{N} \times n} \\ \vdots \\ q_{\frac{C'}{N} \times (N-1) + 1} \\ \vdots \\ q_{C'} \end{pmatrix} = \begin{pmatrix} \mathbf{Q}_1 \\ \vdots \\ \mathbf{Q}_n \\ \vdots \\ \mathbf{Q}_N \end{pmatrix}, \\
 \mathbf{K} = (k_j) &= \begin{pmatrix} k_1 \\ \vdots \\ k_{\frac{C'}{N}} \\ \vdots \\ k_{\frac{C'}{N} \times (n-1) + 1} \\ \vdots \\ k_{\frac{C'}{N} \times n} \\ \vdots \\ k_{\frac{C'}{N} \times (N-1) + 1} \\ \vdots \\ k_{C'} \end{pmatrix} = \begin{pmatrix} \mathbf{K}_1 \\ \vdots \\ \mathbf{K}_n \\ \vdots \\ \mathbf{K}_N \end{pmatrix}. \quad (7)
 \end{aligned}$$

则归一化的注意力图可由式(8)给出。从而可得SA的输出特征图,如式(9)~(10)所示。

$$\begin{aligned}
 f_{SA}(\mathbf{Q}, \mathbf{K}, N) &= (\pi_1, \dots, \pi_n, \dots, \pi_N) \cdot \\
 &\left(\text{softmax} \begin{bmatrix} \mathbf{Q}_1^T \mathbf{K}_1 \\ \vdots \\ \mathbf{Q}_n^T \mathbf{K}_n \\ \vdots \\ \mathbf{Q}_N^T \mathbf{K}_N \end{bmatrix} \right), \quad (8)
 \end{aligned}$$

$$\mathbf{X}_{ss}' = f_{SA}(\mathbf{Q}, \mathbf{K}, N) \cdot \mathbf{V}^T, \quad (9)$$

$$\mathbf{X}_{ss} = \text{reshape}(\mathbf{W}_s \cdot (\mathbf{X}_{ss}')^T), \quad (10)$$

式中: π_n 为第 n 个可学习的聚合权重; \mathbf{W}_s 为可学习的 $C \times C'$ 权重矩阵; $\mathbf{X}_{ss} \in \mathbf{R}^{C \times H \times W}$ 。

2.3.2 向下注意力机制

\mathbf{X}_i'' 表示输入的高层特征图, \mathbf{X}_i'' 上采样后的特征图为 \mathbf{X}_i' , \mathbf{X}_i' 调整尺寸后的特征图为 \mathbf{X}_i 。 \mathbf{X}_d' 表示输入的低层特征图, \mathbf{X}_d' 调整尺寸后的特征图为 \mathbf{X}_d 。

设 $\mathbf{X}_i'' \in \mathbf{R}^{C \times H_i \times W_i}$ 、 $\mathbf{X}_i' \in \mathbf{R}^{C \times H_i \times W_d}$ 、 $\mathbf{X}_i \in \mathbf{R}^{C \times H_d \times W_d}$ 、 $\mathbf{X}_d' \in \mathbf{R}^{C \times H_d \times W_d}$ 、 $\mathbf{X}_d \in \mathbf{R}^{C \times H_d \times W_d}$ 。若查询、键、值的权重矩阵分别由可学习的 $C' \times C$ 的矩阵 \mathbf{W}_q 、 \mathbf{W}_k 、 \mathbf{W}_v 表示,则查询、键、值矩阵可分别由 $\mathbf{Q} = \mathbf{W}_q \mathbf{X}_d$ 、 $\mathbf{K} = \mathbf{W}_k \mathbf{X}_i$ 、 $\mathbf{V} = \mathbf{W}_v \mathbf{X}_i$ 表示。

一般而言,不同尺度的特征图提取不同的上下文信息或语义信息。特征图语义信息不同时,负欧式距离表达相似度比点积更有效^[17]。对 \mathbf{Q} 、 \mathbf{K} 采用式(7)所示的分块,使用负欧式距离进行相似度计算,可得归一化的注意力图,如式(11)所示。于是可得DA的输出特征图,如式(12)~式(13)所示。

$$\begin{aligned}
 f_{DA}(\mathbf{Q}, \mathbf{K}, N) &= \\
 &(\pi_1, \dots, \pi_n, \dots, \pi_N) \cdot \\
 &\left(\text{softmax} \begin{bmatrix} -\|\mathbf{Q}_1^T - \mathbf{K}_1^T\|^2 \\ \vdots \\ -\|\mathbf{Q}_n^T - \mathbf{K}_n^T\|^2 \\ \vdots \\ -\|\mathbf{Q}_N^T - \mathbf{K}_N^T\|^2 \end{bmatrix} \right), \quad (11)
 \end{aligned}$$

$$\mathbf{X}_{sd}' = f_{DA}(\mathbf{Q}, \mathbf{K}, N) \cdot \mathbf{V}^T, \quad (12)$$

$$\mathbf{X}_{sd} = \text{reshape}(\mathbf{W}_d \cdot (\mathbf{X}_{sd}')^T), \quad (13)$$

式中: π_n 为第 n 个可学习的聚合权重; \mathbf{W}_d 为可学习的 $C \times C'$ 权重矩阵; $\mathbf{X}_{sd} \in \mathbf{R}^{C \times H_d \times W_d}$ 。

2.3.3 向上注意力机制

\mathbf{K} 、 \mathbf{V} 均表示输入的低层特征图, \mathbf{Q} 表示输入的高层特征图, \mathbf{Q}_i 为 \mathbf{Q} 的第 i 个通道。设 $\mathbf{Q} \in \mathbf{R}^{C \times H_i \times W_i}$ 、 $\mathbf{K} \in \mathbf{R}^{C \times H_d \times W_d}$ 、 $\mathbf{V} \in \mathbf{R}^{C \times H_d \times W_d}$ 。 \mathbf{K} 经全局平均池化得到权重 \mathbf{W} ,如式(14)所示。 \mathbf{W} 的第 i 个分量 W_i 与 \mathbf{Q}_i 数乘、拼接、卷积后得到 \mathbf{Q}_w 。 \mathbf{V} 下采样后与 \mathbf{Q}_w 求和,再经 3×3 卷积得到UA的输出 \mathbf{X}_{st} ,见式(15)~式(16)。

$$\mathbf{W} = \text{GAP}(\mathbf{K}), \quad (14)$$

$$\mathbf{Q}_w = \text{Conv}_{3 \times 3}(\text{concat}(W_i \cdot \mathbf{Q}_i)), \quad (15)$$

$$\mathbf{X}_{st} = \text{Conv}_{3 \times 3}(\mathbf{Q}_w + \text{SConv}_{3 \times 3}(\mathbf{V})), \quad (16)$$

式中: $\mathbf{W} \in \mathbf{R}^{C \times 1}$ 、 \mathbf{Q}_w 、 $\mathbf{X}_{st} \in \mathbf{R}^{C \times H_i \times W_i}$; $\text{Conv}_{3 \times 3}$ 和 $\text{SConv}_{3 \times 3}$ 分别表示 3×3 卷积和带步长的 3×3 卷积。

2.3.4 注意力金字塔增强

如图4所示, \mathbf{X}_1 、 \mathbf{X}_2 、 \mathbf{X}_3 、 \mathbf{X}_4 、 \mathbf{X}_5 表示APE模块的输入特征图。 \mathbf{X}_1 经SA、UA作用后得 \mathbf{Y}_1 , \mathbf{X}_5 经SA、DA作用后得 \mathbf{Y}_5 。 \mathbf{X}_2 、 \mathbf{X}_3 、 \mathbf{X}_4 分别经SA、DA、UA作用后得 \mathbf{Y}_2 、 \mathbf{Y}_3 、 \mathbf{Y}_4 。从下到上,依次堆叠 \mathbf{Y}_1 、 \mathbf{Y}_2 、 \mathbf{Y}_3 、 \mathbf{Y}_4 、 \mathbf{Y}_5 中的第 i 层得到 \mathbf{Z}_i 。自上而下依次拼接 \mathbf{X}_i 与 \mathbf{Z}_i 得到 \mathbf{C}_i , \mathbf{C}_i 经卷积得输出特征图 \mathbf{O}_i 。

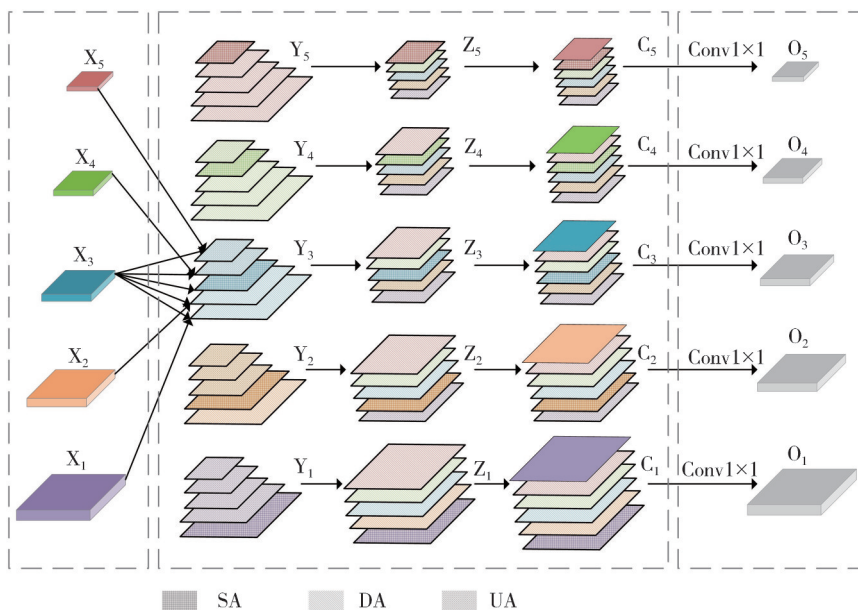


图4 APE模块

Fig. 4 APE module

APE创建SA、DA、UA机制以建立跨空间与跨尺度的特征交互，这丰富了特征图的上下文信息，提高了模型的分割性能。

2.4 IDB网络模型

IDB网络包含低、中、高3个分辨率分支，如图5

所示。低、中分辨率分支构建DDRNet网络，低分辨率分支创建BPPM和APE模块。DDRNet消除空洞卷积引入的网格伪影，提高网络的表征能力；BPPM引入自顶而下和自底向上的路径，加强多尺度特征的融合；APE构造自注意力、向上注意力、向下注意力机制，建立跨空间与跨尺度的特征交互。

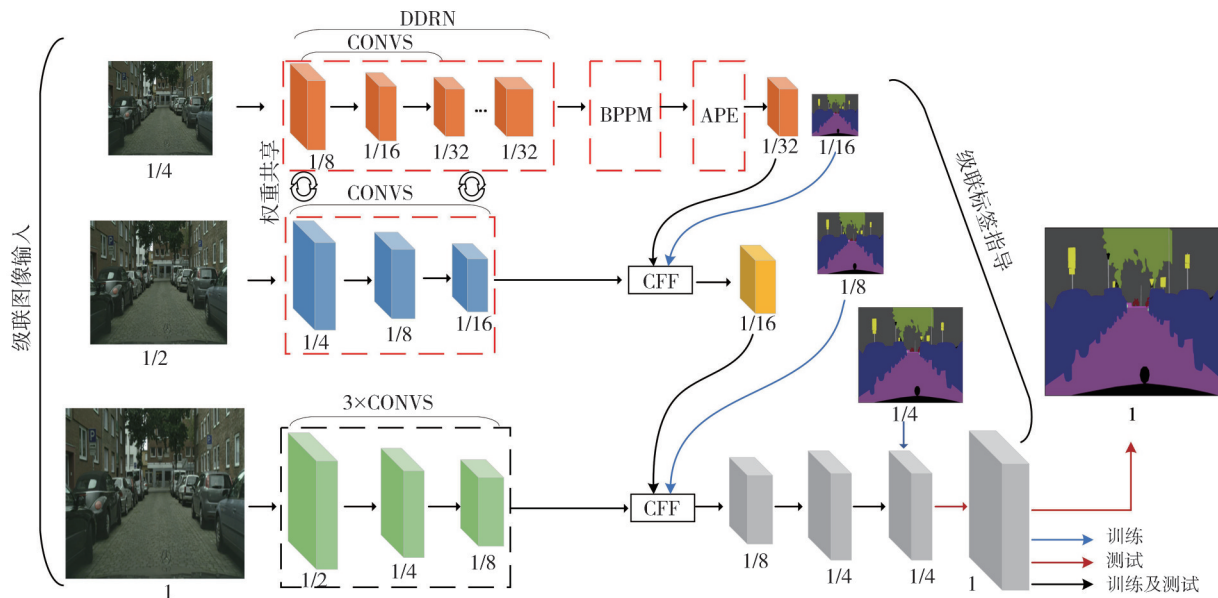


图5 IDB网络模型

Fig. 5 IDB network model

3 实验结果与分析

3.1 数据集与评价指标

实验采用 Cityscapes 数据集评估模型性能。

该数据集包含 2 975 幅训练图像、500 幅验证图像和 1 525 幅测试图像，共有 19 个语义类别。实验采用的测试集为验证集。

评价指标为总体精确率^[18](Overall Accuracy, OA)、平均像素精确度 (Mean Pixel Accuracy,

MPA)、平均交并比(Mean Intersection over Union, MIoU)、频率权重交并比(Frequency Weighted Intersection over Union, FWIoU)和 FPS。各评价指标的数学表达式如式(17)~式(21)所示。

$$OA = \frac{\sum_{i=0}^k p_{ii}}{\sum_{i=0}^k \sum_{j=0}^k p_{ij}}, \quad (17)$$

$$MPA = \frac{1}{k+1} \sum_{i=0}^k \frac{p_{ii}}{\sum_{j=0}^k p_{ij}}, \quad (18)$$

$$MIoU = \frac{1}{k+1} \frac{\sum_{i=0}^k p_{ii}}{\sum_{j=0}^k p_{ij} + \sum_{j=0}^k p_{ji} - p_{ii}}, \quad (19)$$

$$FWIoU = \frac{1}{\sum_{i=0}^k \sum_{j=0}^k p_{ij}} \frac{\sum_{i=0}^k \sum_{j=0}^k p_{ij} p_{ii}}{\sum_{j=0}^k p_{ij} + \sum_{j=0}^k p_{ji} - p_{ii}}, \quad (20)$$

式中: $k+1$ 为类别总数; p_{ij} 表示实际类别为 i 的像素被标记为类别 j 的像素数量。

$$FPS = \frac{N}{\sum_{j=1}^N T_j}, \quad (21)$$

式中: N 为图像数量; T_j 为处理第 j 幅图像所用的

时间。

3.2 实验设置

硬件配置: Intel(R) Xeon(R) Platinum CPU, NVIDIA GeForce RTX 3090 GPU。深度学习框架为 pytorch1.10, GPU 加速库为 CUDA11.3 和 Cudnn8.2.0。

训练阶段, 选取 SGD (Stochastic Gradient Descent) 优化器更新参数。动量和权重衰减分别设置为 0.9 和 0.0001。初始学习率为 0.01, 利用 poly 策略动态调整学习率。批处理大小为 8, epoch 设置为 200。

3.3 定量分析

为评估 IDB 算法性能, 在 Cityscapes 数据集上与 ENet^[19], CGNet^[20], LEDNet^[21], ICNet^[16], DUNet^[22], BiSeNet^[23], EncNet^[24], OCNet^[15] 8 种主流算法进行对比, 各类别 IoU 如表 2 所示。由表 2 可知, IDB 算法 6 个目标类别的 IoU 值为最优, 4 个目标类别的 IoU 值为次优。且 IDB 算法的 MIoU 值为 72.1%, 分别比 ENet、CGNet、LEDNet、ICNet、DUNet、BiSeNet、EncNet、OCNet 算法的 MIoU 值提高了 13.75%、8.19%、7.06%、3.04%、2.84%、2.39%、1.26%、1.11 百分点。

表 2 9 种算法的各类别交并比
Tab. 2 Each class IoU of nine algorithms

类别	IoU/%								
	ENet	CGNet	LEDNet	ICNet	DUNet	BiSeNet	EncNet	OCNet	IDB
road	96.92	97.02	97.13	97.32	96.95	<u>97.52</u>	96.84	96.98	97.70
swalk	77.11	77.56	77.80	79.13	80.02	<u>81.03</u>	79.14	79.69	82.08
build	88.85	89.39	89.24	89.64	89.96	<u>90.37</u>	90.25	90.23	91.11
wall	33.64	40.95	46.60	<u>50.30</u>	42.82	43.09	32.54	34.71	52.98
fence	41.30	46.38	47.78	48.45	48.70	50.92	49.62	50.33	<u>50.57</u>
pole	52.03	52.12	50.94	48.51	62.78	53.28	<u>58.47</u>	57.05	54.94
tlight	42.78	53.55	49.93	53.37	<u>67.89</u>	60.36	68.66	67.24	59.49
sign	61.62	63.77	61.77	64.53	76.01	70.38	<u>75.01</u>	74.60	72.43
veg	90.46	90.40	90.35	90.89	91.75	91.29	<u>91.62</u>	91.48	91.61
terrain	55.64	57.14	55.37	<u>60.53</u>	59.47	60.06	61.72	60.09	59.80
sky	93.74	93.23	92.75	93.42	93.20	<u>93.87</u>	93.38	92.48	94.11
person	68.58	71.35	70.58	72.79	81.01	74.83	<u>79.48</u>	78.77	76.07
rider	36.70	45.94	47.12	47.48	58.50	50.69	<u>57.08</u>	56.21	55.25
car	90.57	90.75	90.72	92.46	92.80	92.83	<u>93.25</u>	93.06	93.70
tuck	32.73	44.89	53.64	72.78	42.41	63.96	57.92	51.11	<u>71.05</u>
bus	45.58	61.59	65.06	76.62	66.08	75.35	75.66	79.85	<u>78.03</u>
train	27.03	31.12	41.26	60.84	35.08	56.61	53.18	68.86	<u>67.06</u>
mbike	11.12	39.66	42.65	44.24	<u>54.21</u>	47.72	56.04	50.54	49.92
bike	62.35	67.52	65.07	68.89	76.27	70.23	<u>76.02</u>	75.52	71.94
MIoU/%	58.35	63.91	65.04	69.06	69.26	69.71	70.84	<u>70.99</u>	72.10

注: 加粗、带下划线字体分别表示各行最优与次优值。

表3给出了9种算法各类别PA值的对比结果。由表3可知, IDB算法5个目标类别的PA值为最优, 3个目标类别的PA值为次优。且IDB算法的MPA值达到了80.05%, 相较于ENet、

CGNet、LEDNet、ICNet、DUNet、BiSeNet、EncNet、OCNet算法, 分别提高了13.06、6.18、4.96、4.23、0.94、1.68、1.73、1.4个百分点。

表3 9种算法的各类别像素准确率
Tab.3 Each class PA of nine algorithms

类别	PA/%								
	ENet	CGNet	LEDNet	ICNet	DUNet	BiSeNet	EncNet	OCNet	IDB
road	98.62	98.59	98.66	<u>98.82</u>	98.35	98.67	98.56	98.42	99.13
swalk	86.76	87.89	88.31	88.83	89.03	90.62	87.61	<u>89.60</u>	89.04
build	95.94	94.90	94.69	96.55	<u>96.14</u>	95.58	96.10	96.04	96.08
wall	36.50	48.14	53.03	<u>57.93</u>	51.43	50.03	34.98	37.22	63.04
fence	54.83	58.76	65.75	56.10	61.33	<u>63.45</u>	58.69	59.25	58.30
pole	62.08	64.17	63.94	56.84	<u>71.43</u>	63.86	71.58	70.10	65.46
tlight	50.05	66.94	59.08	60.74	<u>78.87</u>	72.25	80.09	77.20	68.93
sign	67.39	75.50	73.09	70.45	<u>85.15</u>	78.02	83.92	85.23	80.88
veg	95.54	95.39	95.38	96.05	95.44	96.36	<u>96.44</u>	96.13	96.45
terrain	65.77	69.26	62.75	68.44	75.45	68.17	69.62	<u>69.68</u>	69.03
sky	96.98	97.12	<u>97.52</u>	96.03	96.67	97.77	96.95	96.99	97.32
person	81.90	85.22	83.56	84.46	89.62	87.29	<u>90.51</u>	91.14	86.78
rider	44.16	59.86	61.09	58.40	71.23	64.60	68.30	67.45	<u>70.10</u>
car	95.86	96.11	95.61	95.56	97.00	96.66	<u>97.08</u>	96.74	97.41
tuck	40.42	55.89	65.99	78.87	55.28	73.79	60.60	53.80	<u>76.08</u>
bus	58.09	83.46	76.61	82.45	72.57	89.54	85.21	84.09	<u>85.49</u>
train	46.87	36.61	56.02	63.56	64.62	60.62	58.83	<u>71.78</u>	74.66
mbike	13.05	48.98	52.62	48.20	66.87	57.48	<u>65.98</u>	65.26	60.38
bike	82.06	80.79	83.01	82.36	86.66	84.27	<u>86.96</u>	88.31	86.49
MPA/%	66.99	73.87	75.09	75.82	<u>79.11</u>	78.37	78.32	78.65	80.05

注: 加粗、带下划线字体分别表示各行最优与次优值。

表4给出了9种算法在4项评价指标上的结果。由表4可知, 与ICNet基准算法相比, IDB算法的OA、MPA、MIoU分别提高了0.63、4.23、3.04个百分点。由表4还可看出, IDB算法的OA、MPA、MIoU均优于其它对比算法。

表4 9种算法客观评价指标对比

Tab4 Comparison of objective evaluation indices of nine algorithms

Models	OA/%	MPA/%	MIoU/%	FPS/(帧·s ⁻¹)
ENet	93.78	66.99	58.35	<u>119.49</u>
CGNet	94.09	73.87	63.91	87.18
LEDNet	94.10	75.09	65.04	90.44
ICNet	94.62	75.82	69.06	62.76
DUNet	94.80	<u>79.11</u>	69.26	105.11
BiSeNet	<u>94.89</u>	78.37	69.71	134.53
EncNet	94.86	78.32	70.84	93.15
OCNet	94.83	78.65	<u>70.99</u>	75.88
IDB	95.25	80.05	72.10	20.48

注: 加粗、带下划线字体分别表示各列最优与次优值。

3.4 定性分析

为定性分析IDB算法的分割性能, 选取5幅图像, 并与ENet、CGNet、LEDNet、ICNet、DUNet、

BiSeNet、EncNet、OCNet算法进行对比。实验结果如图6所示, 黄色框为效果提升明显部分。

由第1列对比图可知, ENet、DUNet、EncNet、OCNet算法对卡车存在错误分割问题; ENet、CGNet、LEDNet、ICNet、BiSeNet算法将行人的一部分错误分割成卡车。相比之下, IDB算法对卡车和行人的分割效果更好。

由第2列对比图可知, IDB算法对汽车尾部的分割效果优于ENet、CGNet、LEDNet、ICNet、DUNet、BiSeNet、EncNet、OCNet算法, 更接近标签图。

由第3列对比图可知, IDB以外的8种算法不同程度地将车道错误分割成人行道。ENet、CGNet、LEDNet、DUNet、BiSeNet、EncNet、OCNet算法将骑行者的一部分错误分割成行人。相比之下, IDB算法对车道和骑行者的分割结果更准确。

由第4列对比图可知, 除IDB算法外, 其余算法均将右边的骑行者错误分割成行人。ENet、CGNet、LEDNet、DUNet、BiSeNet、EncNet、OCNet算法对

左边的骑行者存在错误分割现象。只有 IDB 算法没有错误分割问题, 分割效果优于其它算法。

由第 5 列对比图可知, ENet、CGNet、LEDNet、DUNet、BiSeNet、EncNet、OCNet 算法对汽车存在错误分割现象。ICNet 算法虽然能分割出汽车, 但在车轮轮廓处, IDB 算法的分割结果与标签图更接近。

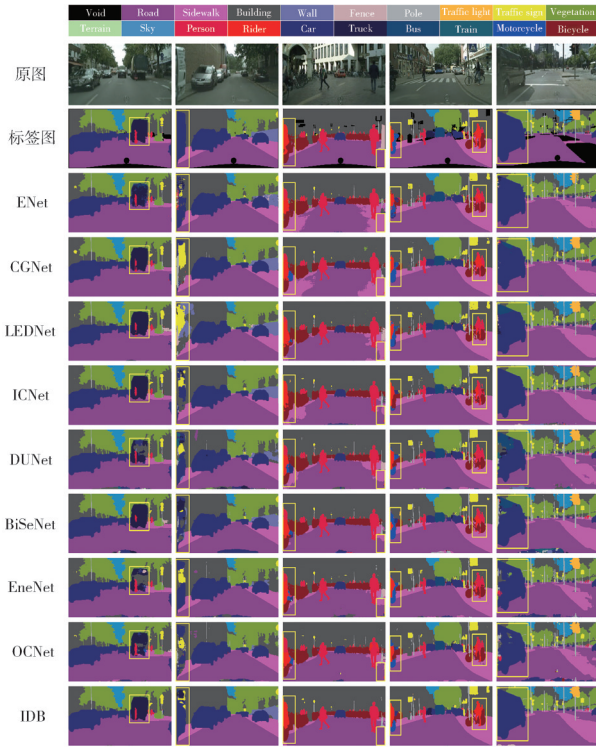


图 6 9 种算法在 Cityscapes 数据集上的分割结果

Fig. 6 Segmentation results of nine algorithms on the Cityscapes dataset

3.5 消融实验

为验证 IDB 算法中 DDRN、BPPM 和 APE 模块的有效性, 进行了 4 组消融实验, 结果如表 5~表 8 所示。

表 5 IDB 算法的消融实验

Tab. 5 Ablation experiment of the IDB algorithm

	APE	BPPM	DDRN	OA/%	MIoU/%	FWIoU/%
1				94.62	69.06	90.15
2		✓		94.74	69.74	90.38
3			✓	94.71	67.61	90.36
4				95.03	70.64	90.87
5	✓	✓		94.82	69.08	90.58
6	✓		✓	95.12	71.72	91.06
7		✓	✓	95.23	71.62	91.30
8	✓	✓	✓	95.25	72.10	91.30

注: 加粗字体表示各列最优值。

表 6 DDRN 的消融实验

Tab. 6 Ablation experiment of the DDRN

Backbone	OA/%	MIoU/%	FWIoU/%
ResNext-50	93.98	67.93	89.08
ResNest-50	94.27	68.21	89.64
Dilated-ResNet-50	94.62	69.06	90.15
Xception	94.66	70.09	90.29
DDRN	95.03	70.64	90.87

注: 加粗字体表示各列最优值。

表 7 BPPM 的消融实验

Tab. 7 Ablation experiment of the BPPM

方法	Backbone	OA/%	MIoU/%	FWIoU/%
PPM	DDRN	95.03	70.64	90.87
PPM+PANet	DDRN	94.94	70.96	90.77
PPM+BiFPN	DDRN	95.17	71.58	91.15
BPPM	DDRN	95.23	71.62	91.30

注: 加粗字体表示各列最优值。

表 5 给出了 DDRN、BPPM 和 APE 模块的消融实验结果。由表中第 1、8 行可看出, IDB 算法的 OA、MIoU、FWIoU 分别比 ICNet 提高了 0.63、3.04、1.15 百分点。由表中第 5、8 行可看出, 缺少 DDRN 时, OA、MIoU、FWIoU 分别下降了 0.43、3.02、0.72 百分点, 这表明 DDRN 消除了空洞卷积引入的网格伪影, 提高了网络的表征能力。由表中第 6、8 行可知, 缺少 BPPM 时, OA、MIoU、FWIoU 分别下降了 0.13、0.38、0.24 百分点, 这说明 BPPM 通过自顶向下和自底向上的路径加强了多尺度特征的融合能力。由表中第 7、8 行可知, 缺少 APE 模块时, OA、MIoU 分别下降了 0.02、0.48 百分点, 这说明 APE 模块通过创建的自注意力、向上注意力、向下注意力机制提高了跨空间与跨尺度的特征交互能力。

表 6 给出了 DDRN 的消融实验结果, 表中第 3 行为基准网络。由表 6 可知, 与采用 ResNext-50、ResNest-50、Dilated-ResNet-50、Xception 的网络相比, 采用 DDRN 网络的 OA 分别提高了 1.05、0.76、0.41、0.37 百分点, MIoU 分别提升了 2.71、2.43、1.58、0.55 百分点, FWIoU 分别提升了 1.79、1.23、0.72、0.58 百分点。可以看出, 5 种骨干网络中, Backbone 为 DDRN 的网络分割性能最好, 验证了 DDRN 模块的有效性。

表 7 给出了 BPPM 的消融实验结果。由表 7 可知, 相较于 PPM、PPM+PANet、PPM+BiFPN, 采用 BPPM 分割方法的 OA 分别提升了 0.20、0.29、0.06 百分点, MIoU 分别提高了 0.98、0.66、0.04 百分点, FWIoU 分别提高了 0.43、0.53、0.15 百分点。可以看出, 4 种方法中, BPPM 的特征融合效果最佳, 验证了 BPPM 的有效性。

表8给出了APE模块的消融实验结果。表中第1~4行、第5~8行的骨干网络分别采用Dilated-ResNet-50和DDRN。由表中第1~4行可知,与ICNet、BPA^[25]和FPN相比,采用APE分割方法的MPA分别提高了1.38、0.2、0.14百分点,MIoU分别提高了0.68、0.52、0.31百分点。

由表中第5~8行可看出,与第5、6、7行相比,采用APE分割方法的OA分别提高了0.02、0.26、0.17百分点,MIoU分别提升了0.48、0.86、0.54百分点。综上可知,各模块中,APE的特征融合效果最好,证明了APE模块的有效性。

表8 APE模块的消融实验

Tab.8 Ablation experiment of the APE module

方法	Backbone	OA/%	MPA/%	MIoU/%	FWIoU/%
ICNet	Dilated-ResNet-50	94.62	75.82	69.06	90.15
ICNet+BPA	Dilated-ResNet-50	94.72	77.0	69.22	90.36
ICNet+FPN	Dilated-ResNet-50	94.78	77.06	69.43	90.46
ICNet+APE	Dilated-ResNet-50	94.74	77.2	69.74	90.38
BPPM	DDRN	95.23	80.61	71.62	91.30
BPPM+BPA	DDRN	94.99	79.56	71.24	90.83
BPPM+FPN	DDRN	95.08	79.69	71.56	91.02
BPPM+APE	DDRN	95.25	80.05	72.10	91.30

注:加粗字体表示各列最优值。

4 结论

本文提出一种去伪影空洞残差与双向注意力金字塔增强的图像级联语义分割算法。该算法构建去伪影空洞残差网络,消除空洞卷积引入的网格伪影,增强网络的特征提取能力;构造双向金字塔池化模块,引入自顶向下和自底向上的路径,充分融合多尺度特征;设计注意力金字塔增强模块,通过自注意力、向上注意力、向下注意力机制建立跨空间与跨尺度的特征交互。实验结果表明,Cityscapes数据集上所提算法的总体精确率、平均像素精确度、平均交并比均高于其它8种主流算法;与基准算法相比,以上3种评价指标分别提高了0.67%、5.58%、4.40%。

参考文献:

[1] JIANG B, AN X, XU S, et al. Intelligent image semantic segmentation: a review through deep learning techniques for remote sensing image analysis[J]. Journal of the Indian Society of Remote Sensing, 2023, 51(9): 1865-1878.

[2] SUN Y, ZHENG W. HRNet-and PSPNet-based multiband semantic segmentation of remote sensing images [J]. Neural Computing and Applications, 2023, 35(12): 8667-8675.

[3] SU Z, LI W, MA Z, et al. An improved U-Net method for the semantic segmentation of remote sensing images [J]. Applied Intelligence, 2022, 52(3): 3276-3288.

[4] HUANG T, SONG S, LIU Q, et al. A novel multi-exposure fusion approach for enhancing visual semantic segmentation of autonomous driving [J]. Proceedings of the Institution of Mechanical Engineers, Part D: Journal of Automobile Engineering, 2023, 237(7): 1652-1667.

[5] FAN J, GAO B, GE Q, et al. SegTransConv: transformer and CNN hybrid method for real-time semantic segmentation of autonomous vehicles [J]. IEEE Transactions on Intelligent Transportation Systems, 2024, 25(2): 1586-1601.

[6] YASMIN S, DURRANI M Y, GILLANI S, et al. Small obstacles detection on roads scenes using semantic segmentation for the safe navigation of autonomous vehicles [J]. Journal of Electronic Imaging, 2022, 31(6): 061806.

[7] JULIUS FUSIC S, HARIHARAN K, SITHARTHAN R, et al. Scene terrain classification for autonomous vehicle navigation based on semantic segmentation method [J]. Transactions of the Institute of Measurement and Control, 2022, 44(13): 2574-2587.

[8] RAHMAN M M, MARCULESCU R. Multi-scale hierarchical vision transformer with cascaded attention decoding for medical image segmentation [C]//Medical Imaging with Deep Learning, PMLR, 2024: 1526-1544.

[9] HATAMIZADEH A, TANG Y, NATH V, et al. UNETR: transformers for 3D medical image segmentation [C]//Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, 2022: 1748-1758.

- [10] TYAGI P, MASSON P, JAIN A, et al. Automated knowledge transfer for medical image segmentation using deep learning [C]//2024 International Conference on Communication, Computer Sciences and Engineering (IC3SE), 2024: 954-960.
- [11] SHELHAMER E, LONG J, DARRELL T. Fully convolutional networks for semantic segmentation [C]//IEEE Transactions on Pattern Analysis and Machine Intelligence, 2017: 640-651.
- [12] BADRINARAYANAN V, KENDALL A, CIPOLLA R. SegNet: a deep convolutional encoder-decoder architecture for image segmentation [J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2017, 39(12): 2481-2495.
- [13] PENG C, ZHANG X, YU G, et al. Large kernel matters--improve semantic segmentation by global convolutional network [C]//IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017: 1743-1751.
- [14] ZHAO H, SHI J, QI X, et al. Pyramid scene parsing network [C]//IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017: 6230-6239.
- [15] YUAN Y, HUANG L, GUO J, et al. OCNNet: object context for semantic segmentation [J]. International Journal of Computer Vision, 2021, 129 (8) : 2375-2398.
- [16] ZHAO H, QI X, SHEN X, et al. ICNet for real-time semantic segmentation on high-resolution images [C]//Lecture Notes in Computer Science. Cham: Springer International Publishing, 2018: 418-434.
- [17] ZHANG Y, HARE J, PRÜGEL-BENNETT A. Learning to count objects in natural images for visual question answering [J/OL]. [2024-03-02]. arXiv: 1802.05766v1. <https://doi.org/10.48550/arXiv.1802.05766>.
- [18] 彭秀平, 仝其胜, 林洪彬, 等. 一种面向散乱点云语义分割的深度残差-特征金字塔网络框架[J]. 自动化学报, 2021, 47(12): 2831-2840.
- [19] PENG Xiuping, TONG Qisheng, LIN Hongbin, et al. A deep residual-feature pyramid network framework for scattered point cloud semantic segmentation [J]. Acta Automatica Sinica, 2021, 47 (12) : 2831-2840. (in Chinese)
- [20] PASZKE A, CHAURASIA A, KIM S, et al. ENet: a deep neural network architecture for real-time semantic segmentation [J/OL]. [2024-03-02]. arXiv: 1606.02147v1. <http://arxiv.org/abs/1606.02147v1>.
- [21] WU T, TANG S, ZHANG R, et al. Cgnet: a light-weight context guided network for semantic segmentation [J]. IEEE Transactions on Image Processing: a Publication of the IEEE Signal Processing Society, 2021, 30: 1169-1179.
- [22] WANG Y, ZHOU Q, LIU J, et al. Lednet: A light-weight encoder-decoder network for real-time semantic segmentation [C]// IEEE International Conference on Image Processing (ICIP), 2019: 1860-1864.
- [23] TIAN Z, HE T, SHEN C, et al. Decoders matter for semantic segmentation: data-dependent decoding enables flexible feature aggregation [C]//IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019: 3121-3130.
- [24] YU C, WANG J, PENG C, et al. BiSeNet: bilateral segmentation network for real-time semantic segmentation [C]//Lecture Notes in Computer Science. Cham: Springer International Publishing, 2018: 334-349.
- [25] ZHANG H, DANA K, SHI J, et al. Context encoding for semantic segmentation [C]//IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018: 7151-7160.
- [26] LIU S, QI L, QIN H, et al. Path aggregation network for instance segmentation [C]//IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018: 8759-8768.