

中图法分类号: 文献标识码: A 文章编号: 1006-8961(XXXX)XX-0001-14

论文引用格式: Swin-EdgeNet: medical image segmentation with multi-scale edge enhancement based on transformer[J/OL]. Journal of Image and Graphics, XXXX:1-14. DOI: 10.11834/jig.250592. (刘善洁, 童孟军. Swin-EdgeNet:基于Transformer的多尺度边缘增强医学图像分割模型[J/OL]. 中国图象图形学报, XXXX:1-14. DOI: 10.11834/jig.250592.) [DOI:10.11834/jig.250592]

Swin-EdgeNet: 基于Transformer的多尺度边缘增强医学图像分割模型

刘善洁, 童孟军

浙江农林大学, 浙江省杭州市 310000

摘要: **目的** 传统医学图像分割方法依赖手工特征提取, 难以应对复杂多样的医学图像, 而现有深度学习方法在边缘细节保留和多尺度特征提取方面仍存在不足。本文旨在提出一种新型深度学习框架, 以提升医学图像分割的精度与鲁棒性。**方法** 提出融合Swin Transformer与U-Net的混合架构, 其中编码器采用Swin Transformer V2-B以捕获分层特征及长范围依赖关系, 解码器引入多尺度边缘密集模块(multi-scale edge dense module, MS-EDM)以增强边缘保留能力, 并结合双重注意力机制(包括细节增强双重注意力(detail-enhanced dual attention, DEVA)与双重注意力(dual attention, DA))以动态聚焦关键区域。为进一步优化模型效率, 将传统Conv2d层替换为KanConv, 在保持性能的同时减少参数数量。**结果** 在ISIC2018皮肤病变、肺部分割、眼底血管分割和息肉分割四个公开数据集上的实验表明, 本文方法在多项评价指标上均优于当前主流方法。具体而言, 在ISIC2018数据集中, Dice和IoU分别达到89.87%与83.97%; 在肺部分割任务中Dice系数达97.62%; 在眼底血管分割中灵敏度达80.19%, 显著优于对比方法; 在息肉分割任务中Dice系数达93.80%, 召回率提升6.58%。可视化结果进一步验证了本方法在边缘保留与小目标识别方面的优势。**结论** 本文提出的融合Swin Transformer与多尺度边缘增强机制的分割模型, 能够有效综合Transformer的全局建模能力与CNN的局部特征提取优势, 在多种医学图像分割任务中均表现出更优的分割精度与泛化能力。

关键词: 医学图像分割; 神经网络架构; Swin Transformer; U-Net; 多尺度特征提取; 注意力机制

Swin-EdgeNet: medical image segmentation with multi-scale edge enhancement based on transformer

Abstract: Objective Medical image segmentation plays a critical role in computer-aided diagnosis, treatment planning, and clinical research. Traditional segmentation methods often rely on handcrafted features and classical image processing algorithms, which struggle with the complexity, heterogeneity, and low contrast frequently present in modern medical imaging data. While deep learning, particularly Convolutional Neural Networks (CNNs) and more recently Transformer-based architectures, has revolutionized the field by enabling automatic feature extraction and achieving state-of-the-art performance, significant challenges remain. Existing deep learning models often exhibit insufficient capability in preserving fine boundary details and effectively leveraging multi-scale contextual information, leading to suboptimal segmentation of irregular or small anatomical structures and pathological regions. To address these limitations, this study aims to design and vali-

收稿日期: 2025-11-24; 修回日期: 2026-03-19

* 通信作者: 童孟军, 通信作者, 男, 教授, 主要研究方向为大数据分析和人工智能。E-mail: tmj@zafu.edu.cn

基金项目: 的规范中文全称(项目编号:……)(不同基金之间用分号隔开)Supported by: 基金项目的英文全称(主要基金项目的中英文名称可在学报网站下载中心查找核对)

date a novel deep learning framework that synergistically combines the strengths of hierarchical feature learning, long-range dependency modeling, and precise edge delineation to enhance the accuracy, robustness, and generalizability of medical image segmentation across diverse tasks and datasets. **Method** We propose Swin-EdgeNet, a hybrid encoder-decoder architecture for medical image segmentation that innovatively integrates a Swin Transformer backbone with a U-Net style decoder enhanced with specialized modules. The encoder is built upon Swin Transformer V2-B, which serves as a powerful feature extractor. Its self-attention mechanism within shifted windows efficiently captures hierarchical feature representations and establishes long-range spatial dependencies across the image, providing a rich global context that is often missing in pure CNN architectures. The corresponding decoder is designed to progressively recover spatial resolution and generate precise segmentation maps. To specifically address the challenge of boundary ambiguity, we introduce a Multi-Scale Edge Dense Module (multi-scale edge dense module, MS-EDM) within the decoder. This module employs parallel convolutional pathways with different receptive fields (dilated convolutions) to simultaneously capture and fuse features at multiple scales. The dense connections within MS-EDM facilitate feature reuse and strengthen gradient flow, ensuring that fine edge details and local textures are preserved and enhanced throughout the upsampling process. Furthermore, we incorporate a dual attention mechanism to dynamically focus the model's capacity on semantically critical and structurally salient regions. This mechanism comprises two components: a Detail-Enhanced Visual Attention (detail-enhanced dual attention, DEVA) module that refines feature maps along the channel dimension by emphasizing interdependent channel maps, and a Dual Attention (dual attention, DA) module that combines spatial and channel attention to selectively aggregate contextual information from relevant positions and feature channels. This dual strategy allows the model to suppress irrelevant background noise while accentuating subtle structures and pathological areas. To optimize the model's computational efficiency without compromising performance, we replace the standard Conv2d layers in specific parts of the decoder with KanConv layers, a parameter-efficient alternative that reduces the overall model parameter count. The entire network is trained end-to-end using a combined loss function, typically involving Dice Loss and Cross-Entropy Loss, to handle class imbalance effectively. **Result** Extensive experiments and quantitative evaluations were conducted on four public benchmark datasets encompassing distinct medical image segmentation tasks to validate the effectiveness and generalization capability of our proposed Swin-EdgeNet. On the ISIC 2018 skin lesion segmentation dataset, our method achieved a Dice Similarity Coefficient (DSC) of 89.87% and an Intersection over Union (IoU) of 83.97%, outperforming existing state-of-the-art methods like Rolling-Unet. The high recall rate of 93.44% indicates a superior ability to capture the complete lesion area, minimizing false negatives. For the lung segmentation task, the model demonstrated exceptional performance, attaining a DSC of 97.62% and an IoU of 95.44%, surpassing strong baselines such as DCSAU-Net (DSC: 97.56%, IoU: 95.35%) and HTC-Net (DSC: 97.48%, IoU: 95.19%). It also maintained high scores in Precision (97.12%), Recall (95.94%), and F1-Score (96.44%), indicating a robust balance between accuracy and coverage. In the challenging domain of retinal fundus image vessel segmentation, where thin capillary structures are difficult to detect, our method achieved a notably high Sensitivity of 80.19%, significantly exceeding that of models like U-Net++ (61.56%) and M2SNet (43.46%). This highlights its enhanced capability for segmenting low-contrast, fine-grained structures. The model also attained an F1-Score of 77.77% and an AUC of 88.63%. For the polyp segmentation task, Swin-EdgeNet achieved a DSC of 93.80% and a Recall of 91.93%, representing an improvement of 2.12% in DSC and 6.58% in Recall over models like Trans-UNet. The F1-Score for polyp segmentation reached 89.81%. Qualitative visual analysis further corroborated these quantitative findings. As illustrated in the result figures, our method produces segmentation masks with noticeably sharper boundaries and more accurate contours for skin lesions, retains more complete small vessel branches in retinal images, and provides more precise coverage for flat and sessile polyps compared to other models, which often suffer from under-segmentation or over-segmentation. **Conclusion** This paper presents Swin-EdgeNet, a robust and effective deep-learning framework for medical image segmentation. By successfully integrating the global contextual understanding of the Swin Transformer encoder with the detailed, localized feature recovery of a progressively enhanced U-Net decoder, the model effectively bridges the gap between high-level semantics and low-level spatial precision. The introduced Multi-Scale Edge Dense Module (MS-EDM) and the dual attention mechanism collectively address key challenges in medical imaging: preserving critical edge information, leveraging multi-scale context, and dynamically focusing on diagnostically relevant regions. The

replacement of standard convolutions with KanConv also contributes to a more parameter-efficient design. Comprehensive experimental results across four diverse and challenging clinical tasks—skin lesion segmentation, lung segmentation, retinal vessel segmentation, and polyp segmentation—consistently demonstrate that our proposed method achieves superior or highly competitive performance compared to current state-of-the-art approaches. Swin-EdgeNet shows strong generalization ability and holds significant promise for improving the accuracy and reliability of computer-aided diagnosis systems.

Key words: medical image segmentation; neural network architecture; swin transformer; U-Net; multi-scale feature extraction; attention

刘善洁¹, 童孟军¹

1 浙江农林大学, 浙江省杭州市 310000

0 引言

医学图像分割是临床诊断和手术规划的关键环节(Litjens等, 2017)。早期的阈值分割(Otsu, 1979)、区域生长(Adams等, 1994)和边缘检测(Canny, 1986)等方法依赖手工特征, 在处理复杂病理结构时往往因抗噪性和泛化力不足而受限(Litjens等, 2017)。

深度学习的兴起为医学图像分割带来了革命性突破。全卷积网络(Lecun等, 2015)首次实现了端到端的像素级预测, 为后续研究奠定了基础。U-Net(Ronneberger等, 2015)凭借其独特的对称编码器-解码器结构和跳跃连接, 在医学图像分割领域取得了显著成功。国内学者也对从卷积神经网络到视觉Transformer的演进历程进行了全面的系统梳理(石军等, 2025)。在此基础之上, 研究者们提出了一系列改进方案: UNet++(Zhou等, 2018)通过嵌套密集连接增强了特征传播效率; Attention U-Net(Oktay等, 2018)引入了注意力机制, 使模型能够聚焦于重要区域; 3D U-Net(Çiçek等, 2016)扩展了三维医学图像的处理能力。

然而, CNN固有的局部感受野限制了其对长程依赖关系的建模能力, 在处理具有复杂空间结构的解剖目标时表现不佳(Wang等, 2021)。为此, 研究者引入了具有全局建模能力的Transformer(Vaswani等, 2017)。

为了平衡局部细节与全局语义, CNN与Transformer的混合架构应运而生: TransUNet(Chen等, 2021)率先将Transformer嵌入U-Net瓶颈层;

UNETR(Hatamizadeh等, 2022)则利用纯Transformer编码器结合CNN解码器处理3D容积数据; nnFormer(Zhou等, 2021)进一步将Transformer集成至nnU-Net框架以提升特征表达。此外, 针对医学图像普遍存在的边界模糊问题, XBound-Former(Wang等, 2022)等研究开始探索基于Transformer的多尺度边界增强机制。虽然Swin Transformer(Liu等, 2021)通过窗口机制降低了计算复杂度, 但如何在保持高效推理的同时, 进一步解决跨域泛化难(Ghafoorian等, 2017; Oktay等, 2018)及人体器官结构复杂、微小病灶边缘模糊导致分割精度不足(蒋婷等, 2024)的问题, 仍是当前亟待解决的挑战。

针对上述问题, 本研究提出了一种新型医学图像分割框架Swin-EdgeNet。该框架以Swin Transformer V2(Liu等, 2022)作为编码器, 充分利用其强大的全局上下文建模能力; 在解码器中引入多尺度边缘密集模块和双重注意力机制, 增强对边缘细节和小目标的感知能力; 同时采用KanConv替代传统卷积, 在保持性能的前提下降低模型复杂度。通过在多个公开数据集上的系统验证, 证明了该方法在保持优异分割性能的同时, 具有更好的泛化能力和边缘保持特性。

1 方法

1.1 整体架构

本文提出的Swin-EdgeNet模型是一种专为医学图像分割设计的深度神经网络架构, 其核心思想是通过层次化的特征提取与精细化的特征融合实现高精度的像素级分割。如图1所示, 模型整体采用编码器-解码器结构, 但与传统U-Net不同, 其编码器采用先进的Swin Transformer V2-B作为主干网络; 而解码器部分则创新性地结合了多尺度特征融合模块(multi-scale edge dense module, MS-EDM)和双重

注意力机制(detail-enhanced dual attention, DEVA), 逐步恢复空间细节并抑制无关噪声, 最终实现病灶区域的精准定位。

在编码阶段, 模型首先通过一个动态调整的前置卷积层将输入图像适配到 Swin Transformer 的输入维度。当处理医学影像常见的单通道数据时, 该层会自动调整卷积核参数, 避免传统 RGB 三通道预训练模型的特征错位问题。在编码器末端, 模型引入了 MS-EDM 模块作为桥梁。通过并行执行不同膨胀率的空洞卷积捕获多尺度上下文信息, 并利用门控融合机制自适应地加权整合各分支特征, 显著提升了模型对病灶大小和形状的鲁棒性。解码器部分

采用对称的 U 型结构, 为了解决深层特征在传递过程中细节丢失的问题, 在解码器的上采样层级嵌入了细节增强双重注意力(DEVA)机制, 利用跳跃连接传递的浅层特征动态校准注意力权重, 从而有效抑制背景噪声并聚焦微小病灶。特别地, 为了在引入上述复杂注意力模块的同时防止模型参数量过度膨胀, 在解码器的上采样阶段策略性地引入了 Kan-Conv 替代部分标准卷积层(Conv2d)。这一设计作为一种高效的参数控制机制, 在保持模型非线性表达能力的同时, 实现了分割精度与计算负载的有效平衡。下面将详细描述各个核心模块的设计细节。

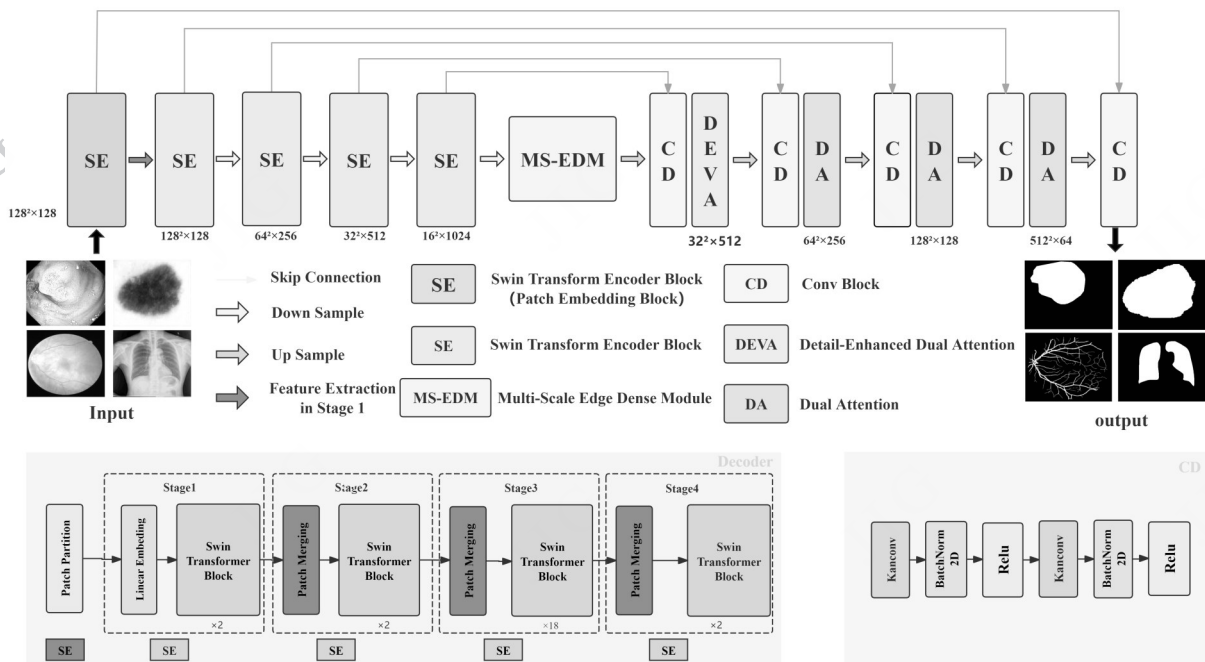


图 1 提出的 Swin-EdgeNet 架构

Fig. 1 The architecture of our proposed Swin-EdgeNet

1.2 编码器结构

医学图像通常具有高分辨率且包含多尺度的解剖结构, 这对骨干网络的显存效率与训练稳定性提出了极高要求。为此, 本文选用 Swin Transformer V2 作为编码器, 旨在利用其强大的全局建模能力捕获长程依赖, 同时克服传统 Transformer 在处理高分辨影像时的容量瓶颈。如图 1 所示, 编码器首先通过 Patch Partition 模块对输入图像进行分块, 随后经四个 Stage 构建多尺度特征。除 Stage 1 外, 其余 Stage 均通过 Patch Merging 层进行下采样。

相较于 V1, Swin Transformer V2 在三个关键技

术上实现突破, 使其更适配医学分割任务:

1. 为解决部分注意力头被少数像素对主导的问题, V2 引入缩放余弦注意力函数, 其像素对 (i, j) 的相似度计算如下:

$$\text{Sim}(q_i, k_j) = \cos(q_i, k_j) / \tau + B_{ij} \quad (1)$$

其中 B_{ij} 是像素 i 和 j 的相对位置偏差, τ 是一个可学习的标量, 在不同的头和层之间不共享, 且值设置为大于 0.01。余弦函数固有的归一化特性确保了注意力权值的平稳分布。

2. V2 摒弃了直接学习参数化偏置的方法, 采用小规模元网络生成连续位置偏置:

$$B(\Delta x, \Delta y) = G(\Delta x, \Delta y) \quad (2)$$

其中G为默认两层MLP(含ReLU激活)。该设计通过元网络为任意相对坐标生成偏置值,使其能够灵活适应不同窗口大小的任务。为增强大范围坐标的外推能力,模型将对数间隔坐标转换引入:

$$\widehat{\Delta x} = \text{sign}(x) \cdot \log(1 + |\Delta x|) \quad (3)$$

$$\widehat{\Delta y} = \text{sign}(y) \cdot \log(1 + |\Delta y|) \quad (4)$$

式中 $\Delta x, \Delta y$ 与 $\widehat{\Delta x}, \widehat{\Delta y}$ 分别为线性间隔与对数间

隔坐标。

3. V2将预归一化调整为后归一化方案,即在残差连接前进行归一化操作,此举有效提升了训练稳定性。

本文采用SwinV2-B作为编码器,其配置参数为:初始通道数 $C=128$,各Stage块数配置 $\text{block}=(2, 2, 18, 2)$ 。该配置在保证强大的特征提取能力的同时,与后续解码器中的KanConv参数控制策略相配合,实现了模型性能与计算开销的良好平衡。

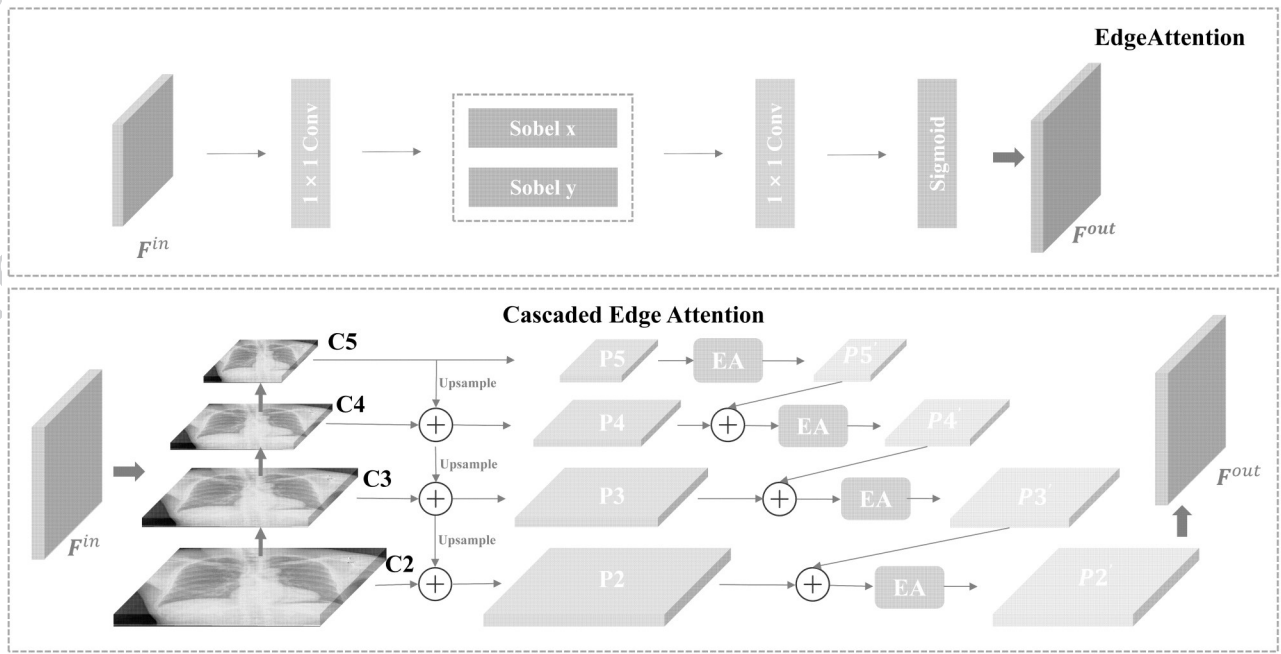


图2 Swin-EdgeNetd的级联边缘注意力结构

Fig. 2 Cascaded edge attention structure of Swin-EdgeNet

1.3 多尺度边缘密集模块

为了解决医学图像中“强语义、弱边界”的矛盾,本文提出了多尺度边缘密集模块(MS-EDM)。如图1所示,被部署在编码器与解码器的跳跃连接处,用于承接并增强多尺度特征。该模块在结构上由两个核心子组件构成:(1)级联边缘注意力模块(cascaded edge attention, CEA):负责利用Sobel算子从多尺度特征中提取边缘先验,解决边界模糊问题。(2)动态门控融合机制:负责自适应地整合上下文信息与边缘特征,替代传统的线性相加。与传统的密集空洞卷积(dense atrous convolution, DAC(Gu等, 2019))或特征金字塔(feature pyramid network, FPN)不同,MS-EDM不仅利用不同膨胀率捕获多尺度上下文,更通过上述两个组件的协同作用,实现了基于

边缘清晰度的特征重加权。

1.3.1 级联边缘注意力模块

级联边缘注意力模块通过改进的特征金字塔网络(FPN)架构,创新性地融入了级联边缘注意力机制。为了精准捕捉病灶轮廓,每个金字塔层级均引入了可学习的Sobel算子来生成边缘响应图,如图2所示。具体而言,水平与垂直方向的梯度特征通过平方和开方运算进行融合,并引入微小常数($1e-6$)确保数值稳定性。边缘响应经 1×1 卷积调整后,通过Sigmoid激活生成注意力权重,与原始特征进行点乘实现边缘增强。这种设计使网络能够自适应强化病灶边界特征,特别适用于医学图像中常见的弱边缘组织分割。

在多尺度特征融合过程中,模块采用级联式信
© 中国图象图形学报版权所有

息传递策略:高层级(如 p5)的边缘信息通过双线性插值下采样后与低层级(如 p4)特征相加,再输入下一级边缘注意力模块。这种设计保持了边缘特征在不同尺度间的连续性。各层级特征在融合前均经过 3×3 卷积平滑处理,以减轻上采样引入的混叠效应。最终输出通过 1×1 卷积将特征维度映射至目标类别数,确保网络适应不同分割任务的需求。

1.3.2 动态门控融合机制

在特征融合阶段,本文设计了动态门控机制以实现自适应特征融合。如图 3 所示,该机制摒弃了传统的串联或简单相加方式,转而采用基于注意力权重的智能融合策略。

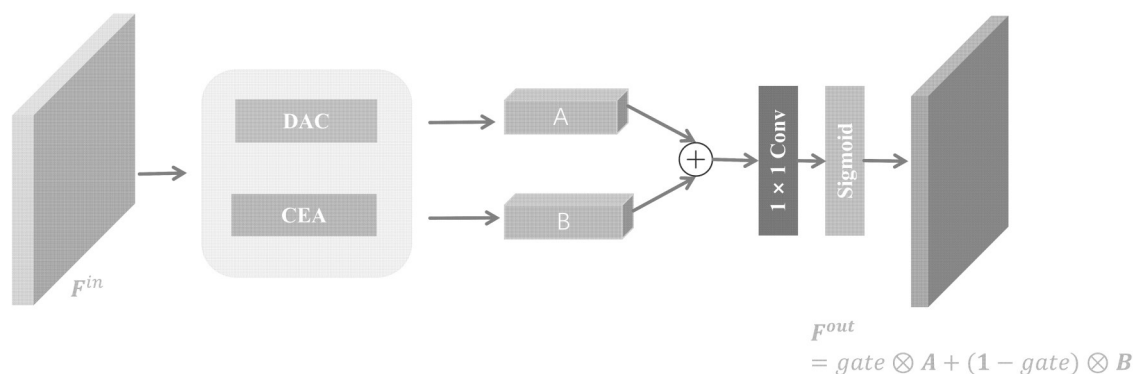


图 3 Swin-EdgeNet 的多尺度边缘密集模块

Fig. 3 Multi-scale edge dense module of Swin-EdgeNet

1×1 卷积的引入实现了跨通道特征交互,避免了简单加权导致的通道间独立性损失;最后,保持原始通道维度的设计确保模块可灵活嵌入各类分割网络,不会引起特征维度膨胀。

1.4 细节增强双重注意力

在医学图像分割任务中,细节信息的保留和关键区域的聚焦对于提高分割精度至关重要。经典的双重注意力模块以 convolutional block attention module (CBAM) (Woo 等, 2018) 为代表,其通过串联通道与空间注意力子模块,从输入特征中自适应地提取关键信息。尽管 CBAM 在自然图像中表现优异,但其注意力的计算仅依赖于当前层级的深层特征。在医学分割中,经过多次下采样后的深层特征往往丢失了微小病灶的边界纹理,导致 CBAM 难以在这些模糊区域生成准确的注意力响应。

为了进一步恢复在深层网络中丢失的细节,本文提出了细节增强双重注意力 (DEVA) 机制嵌入到解码器的上采样层级之后。与 CBAM 仅利用单输

具体而言,模块首先将 DAC 与 CEA 两个支路的输出特征在通道维度进行拼接,随后通过包含 1×1 卷积与 Sigmoid 激活的轻量级网络生成空间-通道自适应的融合权重(取值范围 0-1)。最终输出采用门控加权求和公式:

$$Fused = gate \otimes A + (1 - gate) \otimes B \quad (5)$$

其中 \otimes 表示逐元素乘法, A 和 B 分别代表 DAC 与 CEA 支路的特征图。这种融合设计具有三重优势:首先,门控机制使网络能够根据局部特征特性动态调整各支路的贡献比例,如在组织边界区域增强边缘支路权重,而在同质区域侧重上下文特征;其次

入特征不同,DEVA 的核心创新在于引入了“双输入”交互策略,旨在通过动态调整特征权重,增强模型对目标区域的聚焦能力,同时保留更多的细节信息。DEVA 机制如图 4 所示,是在保留传统的 CBAM 基础上,引入了浅层特征的细节信息,以增强模型对目标区域的细节捕捉能力。这种设计不仅增强了模型对复杂结构的聚焦能力,更从结构上解决了传统注意力机制在深层网络中细节丢失的痛点。具体而言,对于输入的特征图 $F \in R^{C \times H \times W}$ 和细节特征 F_d ,首先通过全局平均池化和全局最大池化分别获取每个通道的全局统计信息:

$$F_d = conv3 \times 3(F_{sh}) \quad (6)$$

$$F_{gap} = GAP(F + F_d) \in R^C, F_{gmp} = GMP(F + F_d) \in R^C \quad (7)$$

其中, F_{sh} 是细节特征,保留了浅层特征中的细节信息。空间注意力模块在 DEVA 中同样引入了细节特征 F_d 以增强空间注意力的计算。对于通道增强后的特征图 F_c 和细节特征 F_d ,首先通过卷积层提取其空间

特征:

$$F_{sh} = conv(F_c + F_d) \quad (8)$$

后续的操作与空间注意力相同。通过细节卷积模块提取浅层特征的细节信息,并将其融入通道注意力和空间注意力的计算中,可以增强模型对目标区域的细节捕捉能力,尤其是在处理具有复杂结构和模糊边界的医学图像时,能够显著提高分割精度。

2 实验

2.1 数据集

在本文的实验中,采用了四个公开的医学图像数据集来验证所提出方法的泛化能力,包括皮肤病变 ISIC2018、肺部 X 光 CHNCXR、眼底血管 DRIVE 和消化道息肉 Kvasir 数据集。由于医学图像具有显著的领域差异性,这种多数据集组合能够全面评估模型在不同解剖结构和病理特征下的表现[1]。ISIC2018、CHNCXR 和 Kvasir 数据集均提供了 512×512 像素的标准分辨率图像,分别包含 3664 张皮

肤镜图像、708 张胸部 X 光片和 1000 张内窥镜图像,这三个数据集的图像尺寸统一性有助于控制预处理阶段的复杂度(Codella 等, 2018)。特别的,DRIVE 眼底数据集原始图像分辨率为 960×960 像

素,考虑到该数据集仅有 40 张训练样本的实际情况,我们将其裁剪为 128×128 的小块,这一策略在不引入人工合成数据的前提下,有效地将样本数量扩充至原始规模的 16 倍(640 张),既缓解了数据稀缺问题,又保持了血管分割任务所需的局部细节特征。皮肤病变和眼底血管数据集均使用官方划分的训练/验证/测试集,而肺部和息肉数据集则按照 7:1.5:1.5 的比例进行随机划分以适配研究需求。

2.2 实验细节

本研究的实验基于统一的计算环境开展,硬件配置采用单块 NVIDIA RTX 4090(24GB 显存)显卡,搭载 Intel Xeon Platinum 8352V 处理器,系统内存 90GB,系统盘容量 30GB,所有实验均在该环境下执行以确保结果的可复现性。训练策略上,针对不同的医学图像分割任务采用差异化的参数配置。对于皮肤病变(ISIC2018)、肺部(CHNCXR)和息肉(Kvasir)分割任务,初始学习率设置为 1e-4,批次大小(batch size)为 4;而眼底血管(DRIVE)分割任务因数据特性采用更高的初始学习率 1e-3,批次大小为 8。所有任务均使用二元交叉熵损失(BCE)和 Dice 损失的加权组合作为优化目标。数据预处理与增强策略根据任务特点进行定制化设计。皮肤病变、肺部、息肉数据集采用随机水平/垂直翻转、±10度旋转以及色彩抖动作为训练集增强手段,并统一应用 Im-

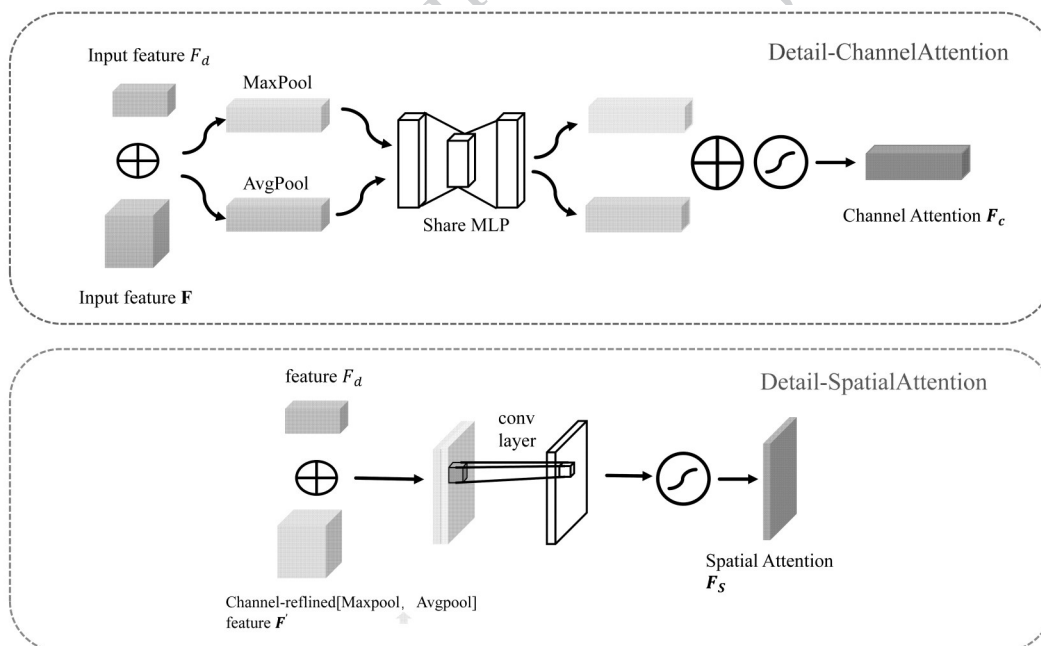


图 4 Swin-EdgeNet 的细节增强双重注意力模块图

Fig. 4 Detail-enhanced dual attention module of Swin-EdgeNet

geNet标准的归一化参数,验证集仅做归一化;而眼底数据集训练和验证集仅图像进行归一化处理。所有输入图像在进入网络前均经过标准化处理,确保数据分布的一致性。训练过程中采用Adam优化器进行参数更新,配合动态学习率调整策略以加速收敛并避免局部最优。测试阶段通过计算所有样本的指标均值作为最终性能评价基准,所有实验均重复三次取平均以消除随机性影响。

2.3 与先进方法的对比分析

本研究在多个医学图像分割任务上进行了全面的性能评估,与近年来最先进的方法进行了对比,包

括 U-Net、U-Net++、Swin-UNet、Trans-UNet、M2SNet (Zhao 等, 2023)、DCSAU-Net (Xu 等, 2023)、HTC-Net (Tang 等, 2024) 和 Rolling-Unet (Liu 等, 2024) 等。实验结果以定量(表 1-4)和定性(图 5-8)两种方式呈现。在定量指标方面,表 1(皮肤病变)、表 2(肺部 CT)、表 3(眼底血管)和表 4(内窥镜息肉)详细对比了本文方法与最新方法的核心指标;在图 5 至 8 中,我们进一步通过可视化结果展示不同方法在典型样本上的分割效果差异,验证了方法的先进性。以下分别对不同数据集的实验结果进行分析。

表 1 不同方法在 ISIC-2018 数据集上的分割性能

Table 1 Segmentation performance of different methods on the ISIC-2018 dataset

Method	Dic(%) ↑	IoU(%) ↑	Pre(%) ↑	F1(%) ↑	Rec(%) ↑	Acc(%) ↑
U-Net(2015)	85.07	77.64	80.92	79.59	86.14	89.63
U-Net++(2018)	86.90	79.85	82.02	82.66	89.28	90.87
Swin-UNet(2021)	88.16	81.02	81.94	84.76	92.15	91.78
Trans-UNet(2021)	89.14	81.73	85.71	86.00	93.34	92.89
XBound-Former(2022)	89.89	83.51	83.65	87.12	94.61	92.96
M2SNet(2023)	88.97	82.50	83.64	86.86	92.74	92.58
HTC-Net(2024)	89.52	83.61	86.34	87.06	91.41	92.87
Rolling-Unet(S)(2024)	89.64	83.70	86.62	87.45	92.33	92.98
Ours	89.87	83.97	86.11	87.24	93.44	93.02

注:加粗字体为每行最优值。

如表 1 所示,在 ISIC2018 皮肤病变数据集中本文的方法在多项关键指标上优于主流分割算法。与 Rolling-Unet 相比,在 Dice 和 IoU 上分别提升了 0.23% 和 0.27%。与 XBound-Former 的对比:尽管 XBound-Former 取得了最高的召回率,但在 IoU (83.97%) 和准确率 (93.02%) 上均优于 XBound-Former,且 F1 分数 (87.24%) 同样保持领先。这表明本文提出的模型在处理复杂病灶时,有效兼顾了病灶的完整检出与边界的精确分割,避免了 XBound-Former 存在的过分割问题,从而获得更

优的区域重叠度(IoU)。

在表 2 肺部分割任务中,本文的方法表现出卓越的分割能力。与当前先进的 Rolling-Unet 相比,本文方法将 Dice 系数和 IoU 分别提升至 97.62% 和 95.44%,达到了表中的最佳水平。在衡量轮廓重合度的 IoU 指标上,本文方法优于 XBound-Former

(95.44% vs 95.31%) 及 Rolling-Unet (95.44% vs 95.35%)。这表明本文提出的网络结构能更有效地提取边缘特征,从而在肺尖和肋膈角等细节区域生成更贴合真实标签的分割掩码。可视化结果如图 6 所示。

针对微细血管密集的 DRIVE 眼底数据集

(表 3),结果显示,本文方法在捕捉血管细节方面具有显著优势。本文方法在灵敏度上取得了 80.19% 的结果,显著优于当前 SOTA 方法 Rolling-Unet (77.21%) 和 DCSAU-Net (77.03%)。这一指标的大幅提升(约 3%)表明,本文模型极大地减少了微小血管的漏检率如图 7 所示。尽管 F1 分数略低于 Rolling-Unet,但考虑到灵敏度的大幅提升,这表明本文方法在临床应用中能提供完整的血管拓扑结构信息。

如表 4 所示,本文的方法在 Dice(93.80%)和召

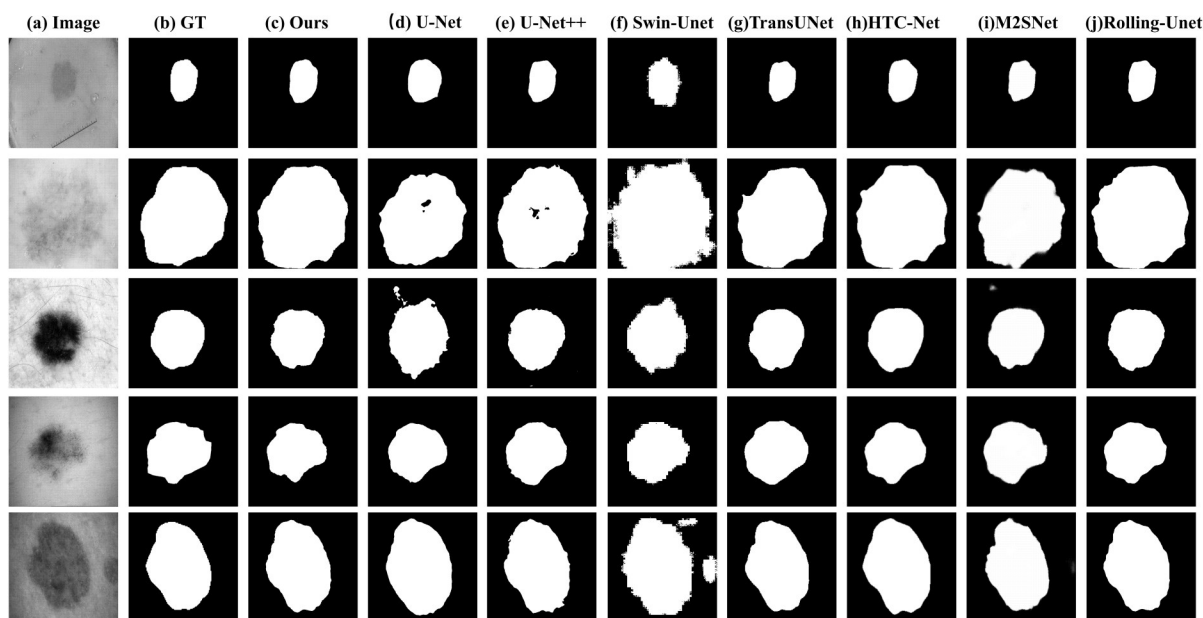


图5 皮肤病变分割结果的可视化对比。(a)原始图像、(b)真实标注(GT)、(c)本文方法 Swin-EdgeNet、(d) U-Net、(e)U-Net++ (f) Swin-Unet、(g) TransUNet、(h) HTC-Net、(i) M2SNet、(j) Rolling-Unet)

Fig. 5 Visualization Comparison of Skin Lesion Segmentation Results. ((a) original image、(b) ground truth (GT)、(c) our Swin-EdgeNet、(d) U-Net、(e)U-Net++ (f) Swin-Unet、(g) TransUNet、(h) HTC-Net、(i) M2SNet、(j) Rolling-Unet)

表2 不同方法在 CHNCXR 数据集上的分割性能

Table 2 Segmentation performance of different methods on the CHNCXR dataset

Method	Dic(%) ↑	IoU(%) ↑	Pre(%) ↑	F1(%) ↑	Rec(%) ↑	Acc(%) ↑
U-Net(2015)	96.71	93.75	96.28	94.11	95.07	97.57
U-Net++(2018)	97.17	94.60	96.72	94.99	95.75	97.91
Swin-Unet(2021)	95.94	92.33	94.76	93.27	93.91	97.00
Trans-Unet(2021)	97.04	94.36	96.79	94.54	95.56	97.83
XBound-Former(2022)	97.55	95.31	97.07	95.80	96.34	98.18
DCSAU-Net(2023)	97.56	95.35	97.02	95.58	96.35	98.21
HTC-Net(2024)	97.48	95.19	96.71	95.95	96.24	98.13
Rolling-Unet(S)(2024)	97.57	95.35	96.77	96.15	96.37	98.19
Ours	97.62	95.44	97.12	95.94	96.44	98.23

注:加粗字体为每行最优值。

回率(91.93%)上分别超越 Trans-Unet 2.12% 和

6.58%。图8的息肉分割案例进一步表明:传统 U-Net 系列(如 U-Net++)对扁平状息肉(图8第1行)的识别能力较弱(F1仅51.44%),而 Swin-EdgeNet 凭借双重注意力机制的聚焦能力,实现了对各类形态息肉的精准覆盖。

2.4 消融实验

为系统验证各模块的贡献,本研究选择肺部分割(CHNCXR)和视网膜血管分割(DRIVE)两个

具有显著解剖结构差异的数据集进行消融实验。

2.4.1 不同组件的消融实验

表5进一步探讨了在 CHNCXR 数据集上,不同组件对肺部分割性能的影响。从表中可以看出,随着各个组件的逐步加入,分割性能逐渐提升。例如加入密集空洞卷积后,Dice系数提升至97.54%,IoU提升至95.29%,这表明密集空洞卷积能够有效增强特征的表达能力。此外,边缘增强注意力(EEA)与门控融合机制的加入进一步优化了边

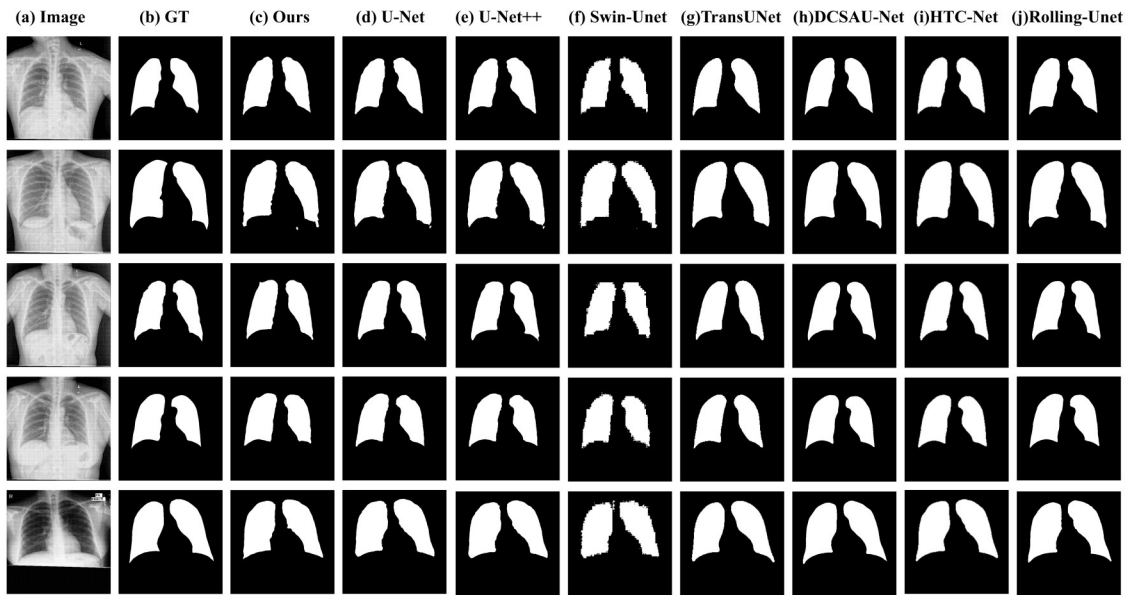


图6 肺部分割结果的可视化对比。(a)原始图像、(b)真实标注(GT)、(c)本文方法 Swin-EdgeNet、(d) U-Net、(e)U-Net++ (f) Swin-Unet、(g) TransUNet、(h) DCSAU-Net、(i) HTC-Net、(j) Rolling-Unet

Fig. 6 Visualization Comparison of Lung Lesion Segmentation Results. ((a) original image、(b) ground truth (GT)、(c) our Swin-EdgeNet、(d) U-Net、(e)U-Net++ (f) Swin-Unet、(g) TransUNet、(h) DCSAU-Net、(i) HTC-Net、(j) Rolling-Unet)

表3 不同方法在DRIVE数据集上的分割性能

Table 3 Segmentation performance of different methods on the DRIVE dataset

Method	Sen (%) ↑	Spe (%) ↑	F1(%) ↑	AUC (%) ↑
U-Net(2015)	64.16	96.12	66.23	80.64
U-Net++(2018)	61.56	96.85	71.00	80.20
MSNet(2023)	69.49	96.18	66.57	82.83
M2SNet(2023)	43.46	97.49	54.66	70.97
DCSAU-Net(2023)	77.03	97.67	78.36	87.58
Rolling-Unet(L)(2024)	77.21	97.69	78.60	87.72
Ours	80.19	97.76	77.77	88.63

注:加粗字体为每行最优值。

界细节,Dice系数进一步提升至97.59%,IoU提升至95.41%。最终,当所有组件全部集成时,Dice系数达到了97.62%,IoU为95.44%,精确率为97.12%,这表明各个组件之间具有良好的协同作用,能够显著提升分割性能。在DRIVE数据集上(表6)也观察到了相同的趋势:各组件的协同作用使灵敏度从基线的96.71%稳步提升至97.62%,验证了本文架构在处理细微结构时的有效性。

2.4.2 边缘注意力消融实验

边缘注意力机制在医学图像分割中具有重要作

用,能够有效突出目标边缘特征,从而提高分割精度。本研究中,我们设计了三种不同的边缘注意力模块(EEA、EEA1、EEA2和EEA3),并对其性能进行了评估,EEA是FPN结构中的p2输出结合

EA(Edge Attention),EEA1是FPN特征融合后加入EA,EEA2在FPN每一个阶段中加入EA,EEA3是级联边缘注意力。

表7展示了在CHNCXR数据集上,不同边缘注意力模块对肺部分割性能的影响。从表中可以看出

出,EEA3在所有指标上均取得了最佳性能,Dice系数达到了97.59%,IoU为95.41%,精确率为97.06%。这表明级联边缘注意力模块能够更有效地突出肺部边缘特征,从而提高分割精度。相比于仅在特定层级添加注意力的EEA1和EEA2,EEA3通过在特征金字塔的每一层级进行级联式的边缘强化,确保了边缘信息在多尺度特征融合过程中的连续性与完整性。

表8进一步探讨了在DRIVE数据集上,不同边缘注意力模块对眼底血管分割性能的影响。结果表明,EEA3同样在所有指标上取得了最佳性能,敏感性达到了78.59%,特异性为97.81%,F1分数为77.76%,AUC为88.14%。这充分证明了级联设计

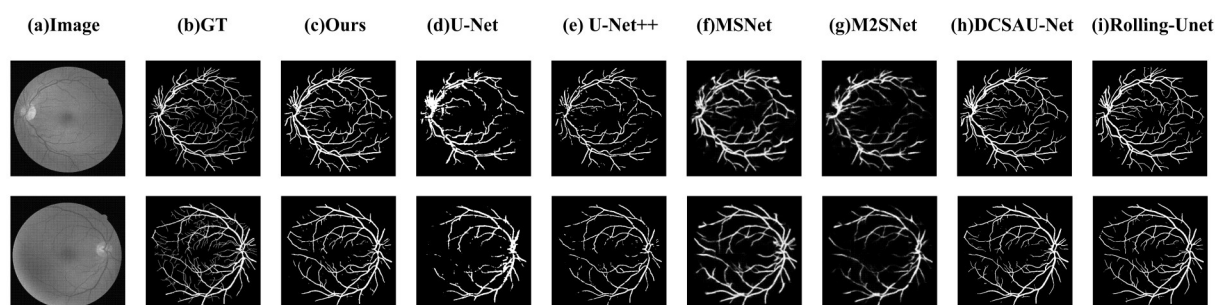


图7 眼底图像分割结果的可视化对比。(a)原始图像、(b)真实标注(GT)、(c)本文方法 Swin-EdgeNet、(d) U-Net、(e)U-Net++ (f) MSUnet、(g) M2SNet、(h) DCSAU-Net、(i) Rolling-Unet)

Fig. 7 Visualization of segmentation results for fundus images. ((a) original image, (b) ground truth (GT), (c) our Swin-EdgeNet, (d) U-Net, (e)U-Net++ (f) MSUnet, (g) M2SNet, (h) DCSAU-Net, (i) Rolling-Unet)

在捕捉复杂血管网络边缘方面的优越性。从而提高分割精度。

3 结论

本研究提出了一种融合 Swin Transformer 与 U-Net 架构的新型医学图像分割框架,通过引入多尺度边缘稠密模块与双重注意力机制,有效应对了不同医学影像中解剖结构多样性与病理特征差异所带来的挑战。系统实验在四个公共数据集 (ISIC2018、CHNCXR、DRIVE、Kvasir) 上展开,验证了本方法在

分割精度与模型鲁棒性方面相较于现有先进技术的优势,表明该框架具备良好的任务适应性与结构泛化能力。从方法学角度总结,本工作的主要特点包括:多尺度边缘信息与深层语义特征的协同提取机制,以及结合空间与通道维度的双重注意力融合策略,从而在复杂医学图像中实现更准确的结构边界识别与内部一致性表达。这些设计在多项分割任务中均体现出稳定且一致的性能提升。

然而,本研究仍存在若干局限:网络整体参数量较大,在计算资源受限场景下的部署效率尚需优化;此外,当前验证数据集的模式与病理类型仍有一定

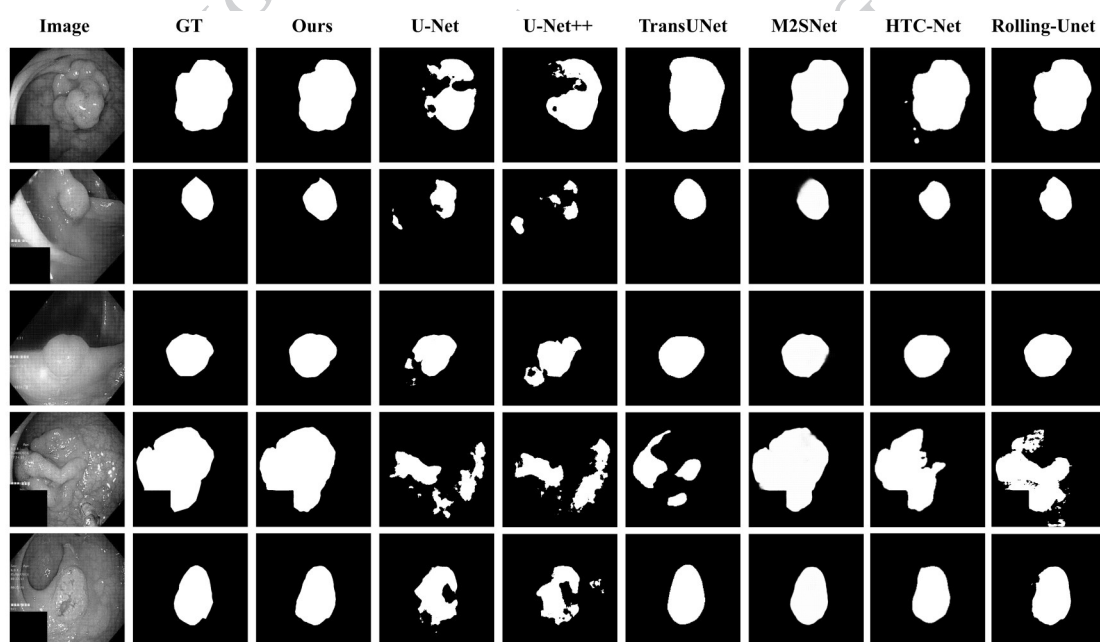


图8 息肉分割结果的可视化对比。(a)原始图像、(b)真实标注(GT)、(c)本文方法 Swin-EdgeNet、(d) U-Net、(e)U-Net++ (f) Swin-Unet、(g) TransUNet、(h) M2SNet、(i) HTC-Net、(j) Rolling-Unet)

Fig. 8 Visualization of segmentation results for fundus images. ((a) original image, (b) ground truth (GT), (c) our Swin-EdgeNet, (d) U-Net, (e)U-Net++ (f) Swin-Unet, (g) TransUNet, (h) M2SNet, (i) HTC-Net, (j) Rolling-Unet)

表 4 不同方法在 Kvasir 数据集上的分割性能

Table 4 Segmentation performance of different methods on the Kvasir dataset

Method	Dic (%) ↑	IoU (%) ↑	Pre (%) ↑	F1 (%) ↑	Rec (%) ↑	Acc (%) ↑
U-Net(2015)	73.70	64.71	71.91	53.50	54.46	88.67
U-Net++(2018)	72.06	63.12	70.08	49.68	51.44	88.18
Trans-UNet(2021)	91.68	86.32	90.48	86.04	86.11	95.84
M2SNet(2023)	87.04	81.25	93.45	73.57	78.22	93.57
HTC-Net(2024)	91.27	85.69	90.52	85.35	85.60	95.51
Rolling-Unet(L) (2024)	91.25	86.18	87.97	88.60	85.49	95.56
Ours	93.80	86.79	91.02	91.93	89.81	96.72

注:加粗字体为每行最优值。

局限,对极端罕见形态的泛化能力有待进一步考察。后续研究将着重于轻量化结构设计,扩展多中心、多

设备采集的实际临床数据验证,并探索结合自监督学习以降低标注依赖。同时,将该框架拓展至动态影像实时分割与多模态融合诊断,亦是未来值得探索的方向。

表 5 在 CHNCXR 数据集上不同组件的影响

Table 5 Effects of different components on the CHNCXR dataset

Models	Dic(%) ↑	IoU(%) ↑	Pre(%) ↑	Rec(%) ↑	F1(%) ↑	Acc(%) ↑
Backbone(U-Net)	96.71	93.75	96.28	94.11	95.07	97.57
UNet_Swin+DA	97.54	95.29	96.96	95.87	96.32	98.17
UNet_Swin+EEA	97.59	95.41	97.06	95.98	96.44	98.23
UNet_Swin+EB	97.60	95.40	97.16	95.84	96.42	98.22
UNet_Swin+MEFF	97.61	95.42	97.06	95.97	96.42	98.21
ALL	97.62	95.44	97.12	95.94	96.44	98.23

注:加粗字体为每行最优值。

表 6 在 DRIVE 数据集上不同组件的影响

Table 6 Effects of different components on the DRIVE dataset

Method	Sen (%) ↑	Spe (%) ↑	F1 (%) ↑	AUC (%) ↑
Backbone (U-Net)	96.71	93.75	96.28	94.11
UNet_Swin+DA	97.54	95.29	96.96	95.87
UNet_Swin+EEA	97.59	95.41	97.06	95.98
UNet_Swin+EB	97.60	95.40	97.16	95.84
UNet_Swin+MEFF	97.61	95.42	97.06	95.97
ALL	97.62	95.44	97.12	95.94

注:加粗字体为每行最优值。

表 7 在 CHNCXR 数据集上不同边缘注意力的影响

Table 7 Effects of different Edge Attention on the CHNCXR dataset

Method	Dic (%) ↑	IoU (%) ↑	Pre (%) ↑	Rec (%) ↑	F1 (%) ↑	Acc (%) ↑
EEA	97.60	95.41	97.04	95.97	96.42	98.22
EEA1	97.48	95.15	96.67	95.97	96.24	98.13
EEA2	97.59	95.42	97.26	95.78	96.43	98.23
EEA3	97.59	95.41	97.06	95.98	96.44	98.23

注:加粗字体为每行最优值。

表 8 在 DRIVE 数据集上不同边缘注意力的影响

Table 8 Effects of different Edge Attention on the DRIVE dataset

Method	Sen (%) ↑	Spe (%) ↑	F1 (%) ↑	AUC (%) ↑
EEA	81.48	97.17	77.49	89.33
EEA1	78.76	97.62	77.52	88.19
EEA2	78.60	97.66	77.59	88.13
EEA3	78.59	97.81	77.76	88.14

注:加粗字体为每行最优值。

参考文献 (References)

Adams R and Bischof L. Seeded region growing [J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 1994, 16(6): 641-647. [DOI: 10.1109/34.295913]

Canny J. A computational approach to edge detection [J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 1986, 8(6): 679-698. [DOI: 10.1109/TPAMI.1986.4767851]

Cao H, Wang Y, Chen J, Jiang D, Zhang X, Tian Q, et al. Swin-unet: unet-like pure transformer for medical image segmentation[C]//Karlinsky L, Michaeli T, Nishino K, eds. Computer Vision - ECCV 2022 Workshops. Cham: Springer, 2023: 205-218. [DOI: 10.1007/978-3-031-25066-8_9]

Chen J, Lu Y, Yu Q, Luo X, Adeli E, Wang Y, et al. Transunet: transformers make strong encoders for medical image segmentation. arXiv preprint arXiv: 2102.04306, 2021. [DOI: 10.48550/arXiv.2102.04306]

Çiçek Ö, Abdulkadir A, Lienkamp S S, Brox T and Ronneberger O. 3D U-net: learning dense volumetric segmentation from sparse annotation [C]. In: Medical Image Computing and Computer-Assisted Intervention - MICCAI 2016. Lecture Notes in Computer Science, vol 9901. Springer, Cham, 2016: 424-432. [DOI: 10.1007/978-3-319-46723-8_49]

Codella N C F, Gutman D, Celebi M E, Helba B, Marchetti M A, Dusza S W, et al. Skin lesion analysis toward melanoma detection: a challenge at the 2017 international symposium on biomedical imaging (ISBI), hosted by the international skin imaging collaboration (ISIC) [C]. In: 2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018). Washington, DC, 2018: 168-172. [DOI: 10.1109/ISBI.2018.8363547]

Ghafoorian M, Mehrtash A, Kapur T, Karssemeijer N, Marchiori E, Pesteie M, et al. Transfer learning for domain adaptation in mri: application in brain lesion segmentation [C]. In: Medical Image Computing and Computer Assisted Intervention - MICCAI 2017. Lecture Notes in Computer Science, vol 10435. Springer, Cham, 2017: 516-524. [DOI: 10.1007/978-3-319-66170-6_60]

Gu Z, Cheng J, Fu H, Zhou K, Hao H, Zhao Y, et al. Ce-net: context encoder network for 2D medical image segmentation [J]. IEEE Transactions on Medical Imaging, 2019, 38(10): 2281-2292. [DOI: 10.1109/TMI.2019.2903562]

Hatamizadeh, Ali, Dong Yang, Holger R. Roth and Daguang Xu. "Unetr: transformers for 3D medical image segmentation." 2022 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV) (2021): 1748-1758. [DOI: 10.1109/WACV51458.2022.00181]

Jiang Ting, Li Xiaoning. 2024. Segmentation of abdominal CT and cardiac MR images with multi scale visual attention. Journal of Image and Graphics, 29(01): 0268-0279 DOI: 10.11834/jig.221032 (蒋婷, 李晓宁. 2024. 采用多尺度视觉注意力分割腹部CT和心脏MR图像. 中国图象图形学报, 29(01): 0268-0279) DOI: 10.11834/jig.221032

Lecun Y, Bengio Y and Hinton . Deep learning [J]. Nature, 2015, 521(7553): 436-444. [DOI: 10.1038/nature14539]

Litjens G, Kooi T, Bejnordi B E, Setio A A A, Ciompi F, Ghafoorian M, et al. A survey on deep learning in medical image analysis [J]. Medical Image Analysis, 2017, 42: 60-88. [DOI: 10.1016/j.media.2017.07.005]

Liu Y, Zhu H, Liu M, Chen J, Li C and Zhang D. Rolling-unet: revitalizing mlp's ability to efficiently extract long-distance dependencies for medical image segmentation [C]. In: Proceedings of the AAAI Conference on Artificial Intelligence, 2024, 38(4): 3819-3827. [DOI: 10.1609/aaai.v38i4.28174]

Liu Z, Lin Y, Cao Y, Hu H, Wei Y, Zhang Z, et al. Swin transformer: hierarchical vision transformer using shifted windows [C]. In: 2021 IEEE/CVF International Conference on Computer Vision (ICCV). Montreal, QC, Canada, 2021: 9992-10002. [DOI: 10.1109/ICCV48922.2021.00986]

Liu Z, Hu H, Lin Y, Yao Z, Xie Z, Wei Y, et al. Swin transformer V2: scaling up capacity and resolution [C]. In: 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). New Orleans, LA, USA, 2022: 11999-12009. [DOI: 10.1109/CVPR52688.2022.01170]

- Oktay O, Schlemper J, Folgoc L L, Lee M, Heinrich M, Misawa K, et al. Attention u-net: learning where to look for the pancreas [J]. arXiv preprint arXiv: 1804.03999, 2018. [DOI: 10.48550/arXiv.1804.03999]
- Otsu N. A threshold selection method from Gray-Level histograms [J]. IEEE Transactions on Systems, Man, and Cybernetics, 1979, 9(1): 62-66. [DOI: 10.1109/TSMC.1979.4310076]
- Ronneberger O, Fischer P and Brox T. U-net: convolutional networks for biomedical image segmentation [C]. In: Medical Image Computing and Computer-Assisted Intervention - MICCAI 2015. Lecture Notes in Computer Science, vol 9351. Springer, Cham, 2015: 234-241. [DOI: 10.1007/978-3-319-24574-4_28]
- Shi Jun, Wang Tiantong, Zhu Ziqi, Zhao Minfan, Wang Bingxun, An Hong. 2025. Deep learning-based medical image segmentation methods. Journal of Image and Graphics, 30(6): 2161-2186 DOI: 10.11834/jig.240467 (石军, 王天同, 朱子琦, 赵敏帆, 王炳勋, 安虹. 2025. 基于深度学习的医学图像分割方法综述. 中国图象图形学报, 30(6): 2161-2186 DOI: 10.11834/jig.240467)
- Tang H, Chen Y, Wang T, Liu M, Zhao X and Ouyang X. Htc-net: a hybrid cnn-transformer framework for medical image segmentation [J]. Biomedical Signal Processing and Control, 2024, 88: 105605. [DOI: 10.1016/j.bspc.2023.105605]
- Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez A N, et al. Attention is all you need [C]. In: Advances in Neural Information Processing Systems, 2017: 5998-6008. [DOI: 10.48550/arXiv.1706.03762]
- Jiacheng Wang, Fei Chen, Yuxi Ma, Li Wang, Zhaodong Fei, Jianfeng Shuai, et al. "Xbound-former: toward cross-scale boundary modeling in transformers." IEEE Transactions on Medical Imaging. 2022. [DOI: 10.1109/TMI.2023.3236037]
- Wang W, Xie E, Li X, Fan D P, Song K, Liang D, et al. Pyramid vision transformer: a versatile backbone for dense prediction without convolutions [C]. In: 2021 IEEE/CVF International Conference on Computer Vision (ICCV). Montreal, QC, Canada, 2021: 568-578. [DOI: 10.1109/ICCV48922.2021.00061]
- Woo S, Park J, Lee J Y and Kweon I S. Cham: convolutional block attention module [C]. In: Computer Vision - ECCV 2018. Lecture Notes in Computer Science, vol 11211. Springer, Cham, 2018: 3-19. [DOI: 10.1007/978-3-030-01234-2_1]
- Xu Q, Ma Z and Duan W. DCSAU-Net: A deeper and more compact split-attention U-Net for medical image segmentation [J]. Computers in Biology and Medicine, 2023, 154: 106626. [DOI: 10.1016/j.compbiomed.2023.106626]
- Zhou HY, Guo J, Zhang Y, Yu L, Wang L, Yu Y. nnFormer: interleaved transformer for volumetric segmentation [J]. 2021. [DOI: 10.1109/TIP.2023.3293771]
- Zhao X, Jia H, Pang Y, Li C, Chen G and Zhang D. M2snet: multi-scale in multi-scale subtraction network for medical image segmentation [J]. arXiv preprint arXiv: 2303.10894, 2023. [DOI: 10.48550/arXiv.2303.10894]
- Zhou Z, Siddiquee M M R, Tajbakhsh N and Liang J. Unet++: a nested u-net architecture for medical image segmentation [J]. IEEE Transactions on Medical Imaging, 2020, 39(6): 1856-1867. [DOI: 10.1007/978-3-030-00889-5_1]

作者简介

刘善洁, 女, 硕士研究生, 主要研究方向为医学图像分割和深度学习。E-mail: 1084585774@qq.com