

中图法分类号: TP391 文献标识码: A 文章编号: 1006-8961(XXXX)XX-0001-18

论文引用格式: Chen Xiaolei, Zhong Zhihua, Shen Yujie. XXXX. A dual-modal salient object detection network for 360° omnidirectional images. Journal of Image and Graphics, XX(XX):0001-0018(陈晓雷, 钟智华, 申玉杰. XXXX. 双模态360度全景图像显著目标检测网络. 中国图象图形学报, XX(XX):0001-0018)[DOI:10.11834/jig.250564]

双模态360度全景图像显著目标检测网络

陈晓雷, 钟智华, 申玉杰

兰州理工大学微电子现代产业学院, 兰州 730050

摘要: 目的 显著目标检测(SOD)旨在模拟人类视觉注意力机制,从图像或视频中识别并分割最显著的物体。尽管基于深度学习的2D SOD已取得显著进展,但面向360°全景图像的SOD因球面投影畸变和边界不连续性问题,面临独特挑战。与此同时,虽然深度信息在2D RGB-D SOD中被证明能增强几何推理能力,但由于模态对齐困难、噪声敏感性以及缺乏针对ERP畸变的融合框架,其在360°SOD中的应用仍探索不足。**方法** 本文提出一种新颖的非对称双分支U-Net网络用于RGB-D 360°SOD。该网络包含全景感知感受野模块(panoramic-aware receptive field module, PA-RFM)、注意力引导融合模块(attention-guided fusion module, AFM)和跨模态引导协同解码策略。PA-RFM通过经度-纬度-全局三重注意力机制缓解投影畸变,AFM实现自适应跨模态特征融合,有效利用深度信息,跨模态引导协同解码策略利用高分辨率RGB解码细节提升边界恢复精度。**结果** 在两个基准数据集(360-SOD和360-SSOD)上的大量主客观实验表明,本文方法性能优于现有10种代表性RGB-D 2D SOD先进方法和7种代表性RGB 360°SOD先进方法,在360-SOD数据集中,相比于性能第2的模型,MAE降低了13.7%,max-F提升4.86%,mean-F提升3.24%,S_m提升1.96%,同时也在360-SSOD数据集中展示竞争优势。**结论** 本文提出的面向360°全景图像显著性检测的RGB-D网络,通过PA-RFM、AFM和跨模态协同解码三大模块协同优化,显著提升了检测精度与鲁棒性,同时验证了深度信息在360°SOD中的有效性与潜力。

关键词: 显著目标检测;360°全景图像;RGB-D;跨模态融合;注意力机制

A dual-modal salient object detection network for 360° omnidirectional images

Chen Xiaolei, Zhong Zhihua, Shen Yujie

School of Microelectronics Industry-education Integration, Lanzhou University of Technology, Lanzhou 730050, China

Abstract: **Objective** Salient object detection (SOD) seeks to emulate human visual attention by locating and segmenting the most conspicuous objects in visual scenes. While deep-learning-based 2D RGB and RGB-D SOD methods have advanced substantially, 360° omnidirectional (panoramic) images introduce unique challenges—most notably spherical-to-plane projection distortions (e. g., ERP artifacts), severe polar-region deformation, and boundary discontinuities across the left-right wraparound. At the same time, depth cues are known to improve geometric reasoning in conventional RGB-D SOD, but their effective exploitation for 360° SOD is underexplored due to modality-alignment difficulties, noise in estimated depth, and the lack of fusion architectures tailored to panoramic projection effects. The objective of this work is to investigate and demonstrate how carefully designed cross-modal modeling and projection-aware feature processing can harness depth information to substantially improve salient object detection on ERP panoramic images, while remaining robust to depth estimation noise and projection artifacts. **Method** We propose an asymmetric dual-branch U-Net architecture for

收稿日期:2025-11-07;修回日期:2026-01-21

基金项目:国家自然科学基金(项目编号:61967012)

Supported by:Project supported by the National Natural Science Foundation of China (Grant No. 61967012)

©中国图象图形学报版权所有

RGB-D 360° SOD that explicitly accounts for ERP projection properties and for cross-modal alignment. The model contains three novel components: (1) a panoramic-aware receptive field module (PA-RFM) that enhances direction-sensitive context modeling via a longitude–latitude–global triplet attention scheme and a ring-padding/cropping strategy to respect horizontal periodicity; (2) an attention-guided fusion module (AFM) that performs adaptive, dynamic weighting of RGB and depth channels using combined channel and spatial attention plus a light-weight modality weight generator to suppress noisy depth responses while preserving boundary cues; and (3) a cross-modal guided collaborative decoder in which a fusion-centric decoding path is explicitly guided by an independent high-resolution RGB decoding path via attentional gating, yielding superior edge and texture recovery. The RGB encoder is a pre-trained SAM2 Hierarchical backbone with lightweight adapters to limit finetuning cost; the depth encoder is a ResNet-34 modified for single-channel input. PA-RFM employs a four-branch multi-scale dilated design with a panoramic context aggregation module (PCAM) at each branch to merge longitude attention (via circular padding + horizontal convolutions), latitude attention (via latitude-weighted vertical convolutions that attenuate polar exaggeration), and global semantic attention (via latitude-weighted GAP and channel recalibration). Training is performed in PyTorch on NVIDIA RTX 3090 hardware using AdamW (lr=1e-3, weight decay=5e-4), cosine annealing to 1e-7, batch size 4, up to 50 epochs, input resolution 512×1024, and a structured loss composed of weighted BCE + weighted IoU with multi-layer deep supervision. **Result** We evaluate the proposed model on two public panoramic SOD benchmarks augmented with depth maps produced by a state-of-the-art 360° depth estimator: 360-SOD (500 ERP images; 400 train / 100 test) and 360-SSOD (1,105 ERP images; 850 train / 255 test). Quantitatively, our method achieves MAE = 0.0151, max-F = 0.8388, mean-F = 0.8150, max-E = 0.9331, mean-E = 0.9240 and S_m = 0.8736 on 360-SOD (same scores on 360-SSOD), outperforming ten representative RGB-D 2D SOD methods and seven representative RGB 360° SOD methods in both objective metrics and visual quality. Relative to the strongest competing baseline, our MAE is reduced by 13.7%, max-F increases by 4.86%, mean-F by 3.24%, and S_m by 1.96%, demonstrating substantial gains particularly in boundary accuracy and structural consistency. Extensive qualitative comparisons show our model is markedly better at detecting objects in severely distorted polar regions, maintaining left–right boundary continuity, and recovering fine edges in cluttered or low-contrast scenes. Ablation studies confirm the contribution of each module: PA-RFM yields the largest single-module improvement (vs. RFB and a no-PA-RFM baseline), PCAM’s three attention branches (longitude, latitude, global) act synergistically, and a ring-padding width of k=2 provides the best tradeoff between boundary continuity and redundancy. AFM outperforms simple fusion strategies (Add, Concat, Multiply) and several SOTA fusion blocks (MobileSal, HIDANet, CPNet variants). The cross-modal guided decoder meaningfully improves edge/detail recovery over a variant that omits the RGB decoding path. We also test robustness to different depth estimators (DA, EGFormer, CRF360D, Joint_360depth) and find consistent improvements over RGB-only baselines; performance varies modestly with depth source, with Joint_360depth yielding the best overall results but not causing catastrophic degradation when replaced. **Conclusion** This study demonstrates that depth information can be effectively and robustly exploited for 360° omnidirectional SOD when combined with projection-aware receptive field design and modality-aware fusion and decoding strategies. The proposed PA-RFM addresses ERP-specific distortions (horizontal periodicity and polar exaggeration) via circular padding, latitude weighting and global recalibration; AFM adaptively reconciles noisy depth with RGB structure; and the cross-modal guided collaborative decoder preserves high-resolution RGB priors to restore fine boundaries and textures lost in fused representations. Together, these innovations lead to consistent and significant gains across standard panoramic SOD benchmarks, especially in polar regions and at projection boundaries. The approach is general and can be extended to other panoramic vision tasks and to setups with sensor-captured depth; future work will explore domain adaptation to real depth sensors, multi-projection fusion (ERP + CMP), and model compression for mobile/AR/VR deployment. **Keywords:** 360° omnidirectional image, salient object detection, RGB-D fusion, projection-aware attention, panoramic context aggregation.

Key words: salient object detection; 360° omnidirectional images; RGB-D; cross-modal fusion; attention mechanism

0 引言

随着全景视觉在 VR/AR 等应用中的日益普及, 360° 全景图像显著目标检测 (360° salient object detection, 后文简称 360° SOD) 逐渐成为视觉理解的重要方向。与常规 2D 图像不同, 360° 全景图像本质上是球形图像, 需要通过等矩柱状投影 (equal rectangular projection, ERP)、立方体投影 (cube map projection, CMP) 等技术转换为平面图像, 才能进行后续处理, 无论哪种投影方式都不可避免的会引起图像畸变, 这种投影畸变会严重影响 SOD 的效果。此外, 360° 全景图像还存在极区畸变、经纬向采样不均匀和边界不连续等复杂特性, 导致显著区域的语义连贯性和空间一致性难以保持。

图像的深度信息可以反映图像的几何结构, 具有内部一致性和光照不变性, 能够为 RGB 图像提供辅助信息从而帮助 SOD 模型有效区分前景和背景, 在光照条件低、目标与背景对比度弱或纹理信息易混淆的情形下, 深度信息仍能保持清晰的空间结构线索, 有效辅助模型区分前景与背景并抑制由纹理导致的误检。然而, 针对 360° 全景图像, 将深度信息有效地与 RGB 信息融合并非易事: 一方面, 现有全景相机无法同时采集 RGB 信息和深度信息, 已有的两个 360° SOD 公开数据集也只是 RGB 单模态数据集, 只能采用深度估计方法为现有数据集增加深度信息。另一方面模态间语义对齐与尺度一致性在 360° 全景图像上更难保证。尽管近年来已有大量 RGB-D 2D SOD 和若干 RGB 360° SOD 研究, 但系统性地将深度信息引入 360° SOD, 并同时考虑 ERP 特性与跨模态鲁棒融合的研究未见文献报道。

针对以上问题, 本文提出一种 RGB-D 360° SOD 网络。网络总体为非对称 RGB-D 双分支 U-Net 结构: 为保证对全景图像畸变区域和大范围显著区域的表征能力, RGB 分支采用了具备强全局建模能力与长程依赖捕获能力的 SAM2-Hiera 作为编码器。深度信息分支采用 ResNet-34 专门提取单通道几何边界特征。为适配 360 度全景图像的几何特性, 提出了全景感知感受野模块 (panoramic-aware receptive field module, PA-RFM) 以增强方向感知上下文建模能力并缓解极区畸变与边界不连续性。为实现跨模态融合, 设计注意力引导融合模块 (attention-

guided fusion module, AFM) 自适应地抑制深度噪声同时保留边界补偿能力; 解码端采用以融合分支为主, 独立 RGB 解码路径作为高分辨率细节引导的协同解码策略, 通过多尺度注意力加权显著提升边缘与纹理恢复。

本文主要创新点和工作总结如下: 1) 提出了一个非对称双分支 U-Net 结构的 RGB-D 360° SOD 网络, 该网络同时考虑了 360° 全景图像自身特性带来的挑战和深度信息对显著性目标检测性能的提高作用, 进一步提高了 360° SOD 性能。2) 设计了全景感知感受野模块 PA-RFM, 通过经度注意、纬度加权与全局语义注意的联合建模, 有效缓解边界不连续与极区畸变问题。3) 设计了注意力引导融合模块 AFM, 以注意力机制为核心驱动, 结合通道与空间注意力、动态权重分配策略及多尺度特征建模, 实现 RGB 特征与深度信息特征的高效对齐与融合。4) 提出跨模态引导协同解码策略, 通过独立的 RGB 解码路径为融合分支提供高分辨率结构先验, 显著提升边缘与细节恢复效果。

5) 在两个公开的 360° SOD 数据集上评估了本文方法, 并将其与现有的 2D RGB-D SOD 和 360° SOD 方法进行比较。实验结果表明, 本文方法的性能超过了现有代表性先进方法。

本文其余部分组织如下: 第 2 节回顾相关工作 (包括 RGB-D SOD 与 360° SOD 的最新进展), 第 3 节详细介绍所提网络结构与各模块设计, 第 4 节给出实验设置、本文方法与现有代表性先进方法的对比结果及消融分析, 最后在第 5 节对本文工作进行总结并讨论未来研究方向。

1 相关工作

1.1 RGB-D 2D SOD 方法

近年来, RGB 2D SOD 虽已取得较好性能, 但在复杂背景与弱光场景下仍存在局限, 因此研究者开始引入深度信息增强显著性检测。随着深度传感器的发展, 早期 RGB-D SOD 工作 (如 Lang 等 2012; Ciptadi 等 2013; Ren 等 2015; Cong 等 2019) 主要依赖手工特征, 性能有限。深度学习出现后, 大量基于卷积神经网络 (convolutional neural networks, CNN) 的 RGB-D 2D SOD 方法取得明显进展, 例如 Wu 等人 (2021) 利用轻量级 CNN 与隐式深度恢复实现高

效多尺度融合。但受限于卷积的局部感受野,CNN在建模长程依赖方面仍存在不足。

随着 Vision Transformer(ViT)的兴起,更强的跨区域建模能力被引入 RGB-D SOD。部分方法采用纯 Transformer 框架,如 Liu 等人(2021)通过三元组 Transformer 嵌入与三分支解码器实现跨模态融合。也有方法结合 CNN 与 Transformer,如 Liu 等人(2021)使用双流 Swin Transformer 进行多模态编码;Lee 等人(2022)将 Transformer 与图卷积神经网络(graph convolutional networks, GCN)结合以实现原型采样与鲁棒融合;Cong 等人(2023)提出 CNN 辅助的 Transformer 架构,通过注意力触发的跨模态点感知交互提升特征融合;Yin 等人(2023)在编码器中加入 RGB-D block 以加强模态交互;Hu 等人(2024)基于 Swin Transformer 设计跨模态融合与渐进解码结构,在无需额外增强模块的情况下获得良好性能。

1.2 RGB 360°SOD 方法

近年来,已有部分研究开始探索 360°SOD,但整体仍较有限。DDS(Li 等,2019)首次提出 RGB 360°SOD 网络,并引入失真自适应模块处理 ERP 畸变;LDNet(Huang 等,2023)基于失真感知与深度可分离卷积缓解投影失真;DATFormer(Zhao 等,2023)利用 Transformer 建模 ERP 特性;ACoNet(Chen 等,2024)通过多分支结构实现跨尺度特征交互;DSANet(Chen 等,2025)结合畸变自适应卷积与多尺度注意增强语义融合;DPNet(Chen 等,2025)采用 ViT+CNN 双编码器并加入可变形卷积以适应 ERP 几何畸变。在互补模态利用方面,FANet(Huang 等,2020)同时利用 ERP 与 CMP 信息实现自适应特征融合,MPFR-Net(Cong 等,2023)采用多 CMP 图像与动态加权策略保持目标结构完整。

2 本文模型

2.1 网络总体结构下

本文提出的模型结构如 1 所示,支持 RGB 图像与深度图并行输入。RGB 分支编码器采用 SAM2(Ravi 等,2024)预训练的 Hiera(Ryali 等,2024)主干,输出多尺度特征 $F_i(i=1,2,3,4)$ 。Hiera 的结构如图 1(a)所示,原始 Hiera 的大参数使得完全微调的计算代价很高,为了提高参数效率,本文采用了 SAM2-UNET(Xiong 等,2024)的适配器设计,冻结了 Hiera

的部分参数,在每个多尺度块之前插入了轻量级适配器。输出多尺度特征后通过全景感知感受野模块(PA-RFM)增强上下文建模以缓解极区畸变;由于深度图本身已具备良好的结构边界和空间位置信息,特征较为简洁直接,深度分支使用轻量化 ResNet-34 提取深度特征,无需复杂增强模块以避免冗余。RGB 与深度特征通过注意力引导融合模块(AFM)自适应整合模态互补信息并抑制噪声,生成融合特征用于多级解码。

解码阶段,RGB 与融合分支分别采用 U-Net 结构的上采样模块进行逐级上采样重建。RGB 分支的解码侧重于纹理细节与色彩信息还原,其结构如图 1(b)所示,最终输出三层辅助显著性预测图,为减少计算量与复杂度,只对拥有最精细特征的输出预测 S_4 进行监督。对于融合特征的解码,本文设计了一个跨模态引导协同解码策略,利用注意力机制将 RGB 解码特征引导融合至融合分支,突出显著区域与边缘信息,强化解码分辨能力。最终,融合分支的输出包含主显著图 S_1 及两层多尺度侧输出 S_2, S_3 , 实现主-辅协同监督策略,有效提升模型的显著性区域聚焦能力与整体检测鲁棒性。

2.2 全景感知感受野模块 PA-RFM

本文为解决 ERP 格式下 360 度全景图像存在的空间几何畸变与上下文不均衡问题,设计了全景感知感受野模块 PA-RFM,其结构如图 2(a)所示,旨在同时建模方向敏感的局部上下文、多尺度语义依赖与结构一致性的全局感知能力。PA-RFM 为四路并行结构,通过不同卷积核尺寸与空洞率的空洞卷积感知多尺度信息,每一路都利用本文设计的全景上下文聚合模块(panoramic context aggregation module, PCAM)提高全景感知能力,最终拼接融合并加残差以增强信息保留,从而兼顾多尺度显著目标与 ERP 图像几何特性。

其中全景上下文聚合模块 PCAM 作为 PA-RFM 的核心组件,专为 ERP 格式下的 360 度全景图像特性量身设计,结构图 2(b)所示,旨在从方向敏感性、空间均衡性与语义一致性三个维度对特征进行增强建模。该模块引入三重上下文建模机制:经度方向注意、纬度方向注意与全局语义注意,以充分刻画 ERP 投影下非均匀几何结构对显著性建模的影响,并融合为一个统一的上下文增强图以调节原始特征表示。

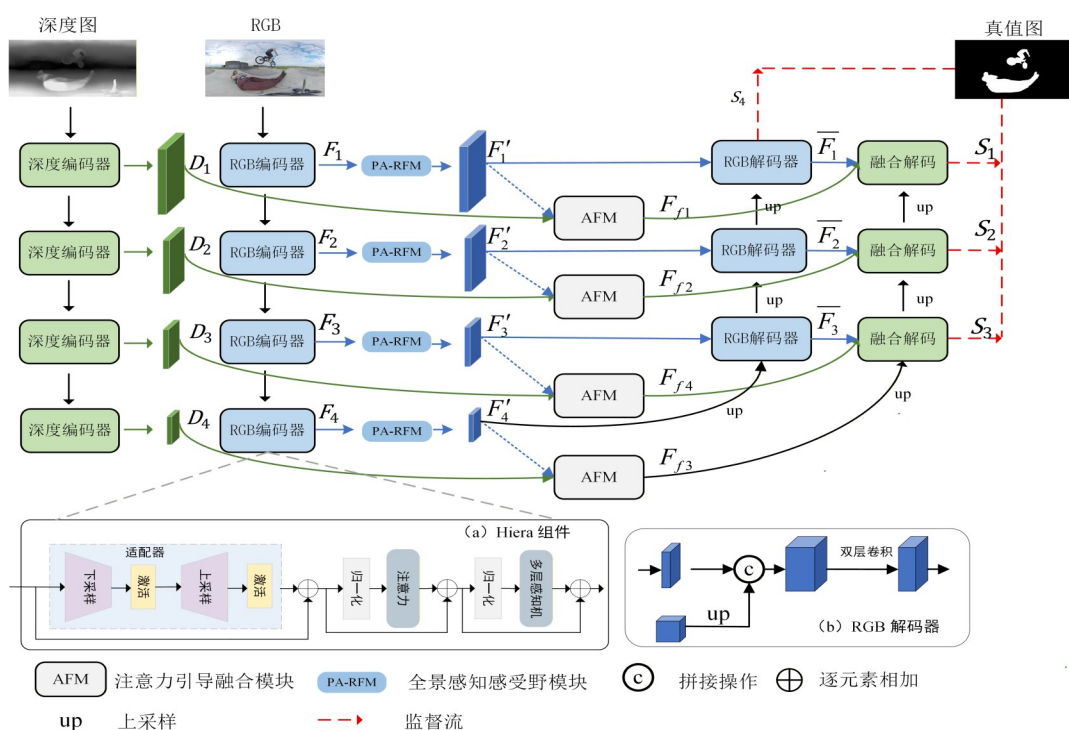


图1 本文模型网络结构

Fig. 1 The network architecture of the proposed model

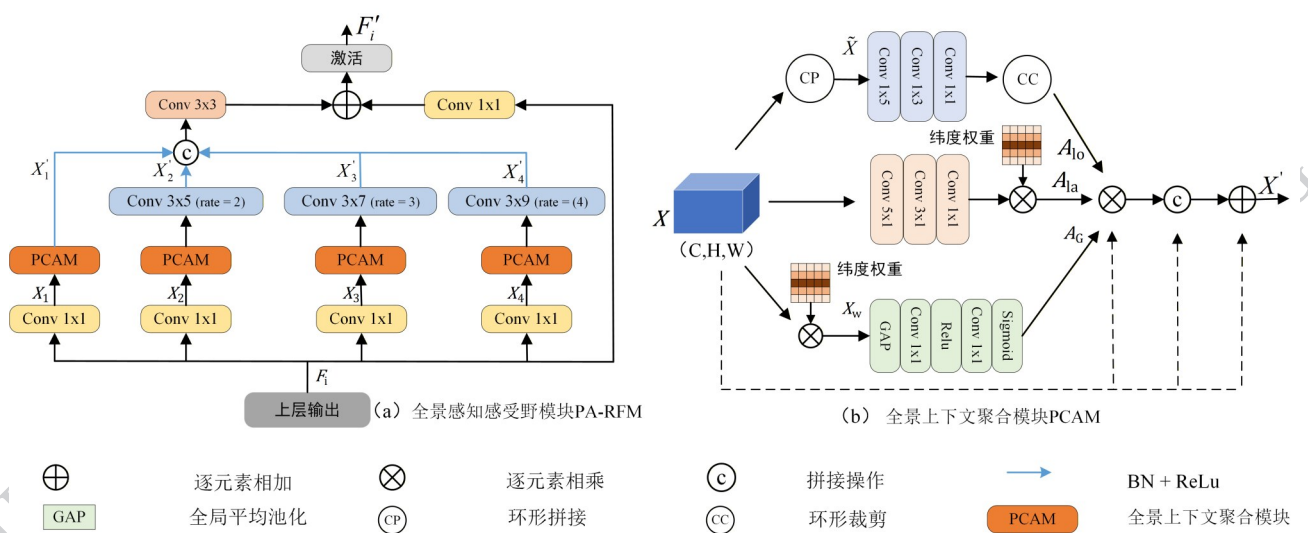


图2 全景感知感受野模块PA-RFM结构图

Fig. 2 Structure of the Panoramic-Aware Receptive Field Module (PA-RFM)

经度方向建模如图2(b)中的上部分支所示,首先针对360度全景图像在水平方向具备周期性边界的拓扑属性,引入环形填充与裁剪策略,以弥补常规卷积在边界处感受野断裂的问题,其效果如图3所示,设给定输入特征 $X \in \mathbf{R}^{B \times C \times H \times W}$, $\mathbf{R}^{B \times C \times H \times W}$ 表示特征的维度。首先进行环形填充,基本原理如图3上部所示,从完整特征的最右侧和最左侧各提取 k

列特征(本文取 $k=2$,既能保证卷积在边界处仍能感知完整的局部窗口并与全景图像的周期性拓扑对齐,同时避免不必要的填充带来的冗余计算,下文环形裁剪操作也取 $k=2$),再分别拼接到左侧和右侧,完成环形扩展:

$$\tilde{X} = \text{concat}\left(X_{[\dots, w-k:w]}, X, X_{[\dots, 0:k]}\right) \in \mathbf{R}^{B \times C \times H \times (W+2k)} \quad (1)$$

式中, \tilde{X} 表示进行环形裁剪后得到的特征图, $concat$ 表示拼接操作, X 表示输入的原始特征。该操作使模型在处理边缘区域时具备连续空间感知能力, 有效提升水平结构一致性建模性能。对于填充后的特征 \tilde{X} , 设计了一个经度方向注意力机制, 以增强其在水平方向的上下文感知能力。该模块依次采用核为 1×5 和 1×3 的卷积聚合横向信息, 随后通过 1×1 卷积恢复通道数, 并经 Sigmoid 得到经度注意力权重图。

随后进行环形裁剪, 基本原理如图 3 下部所示, 对 A'_o 在宽度维度裁剪去除左右各 k 列, 恢复为与原始输入宽度相同的经度注意力图:

$$A_{lo} = A'_o[:, :, k:W+k] \in \mathbf{R}^{B \times C \times H \times W} \# (2)$$

式中, A_{lo} 表示进行环形裁剪后得到的特征图, A'_o 表示进行卷积操作后的中间结果, k 表示裁剪宽度。该机制可有效建模跨水平方向的上下文依赖关系, 从而提升模型对全景图像水平连续区域的理解能力。

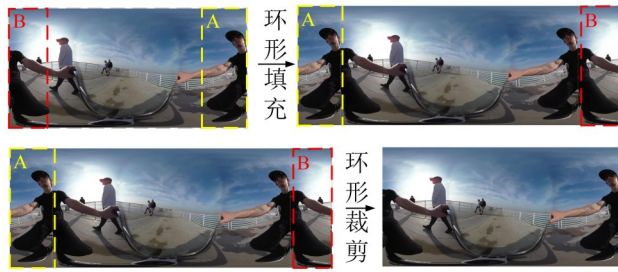


图 3 环形填充与裁剪示意图

Fig. 3 Illustration of Ring Padding and Cropping

在以上经度方向建模中, 主要解决了全景图像中边界不连续的问题, 但是全景图像还有一个严重问题, 就是在极区的严重畸变与失真, 为解决这一问题, 如图 2(b) 中的中部分支所示, 本文首先根据 ERP 图像的球面投影特性, 对每一行像素计算对应的纬度角:

$$\theta_i = -\frac{\pi}{2} + \frac{i\pi}{H-1}, i = 0, \dots, H-1 \# (3)$$

式中, θ_i 表示表示计算得到的纬度角, i 代表特征图高度方向上的像素行索引, H 代表 ERP 格式全景图像特征图的高度。以余弦函数生成初始权重向量, 再乘以一个可学习标量 ∂ :

$$w_{la}(i) = \partial \cos \theta_i \# (4)$$

式中, w_{la} 表示特征对应的纬度权重, ∂ 表示可学习标

量。这一权重能够衰减极区过度聚集的投影效应, 强化赤道附近区域的特征响应。接着, 对原始特征 $X \in \mathbf{R}^{B \times C \times H \times W}$ 逐像素乘以 w_{la} , 先后施加三层竖向卷积以捕获多尺度南北上下文, 具体方法如下: 首先由 5×1 卷积及后续的批归一化和 ReLU 激活提取宽度范围的垂直信息; 紧接着用 3×1 卷积加 BN (batch normalization, BN) + ReLU 对该信息进行局部细化; 最后以 1×1 卷积恢复原始通道数并通过 Sigmoid 激活生成归一化的纬度注意力图。

将该注意力映射与预先计算的纬度权重相乘, 可得最终的纬度注意力。

$$A_{la} = \tilde{A}_{la} \otimes w_{la} \in \mathbf{R}^{B \times C \times H \times W} \# (5)$$

式中, A_{la} 表示纬度注意力图, \tilde{A}_{la} 表示进行三层竖向卷积后得到的特征。

在完成经度与纬度方向的局部上下文建模后, 局部感受野虽然能有效捕捉跨边界与方向感知信息, 但是依然难以覆盖整幅图像的语义全貌, 尤其对于 360 度全景图像中跨区域的长距离依赖关系。为此, PCAM 模块进一步引入全局上下文建模机制, 其结构如图 2(b) 下部分支所示。具体而言, 首先利用纬度权重对原始特征进行加权:

$$X_w = X \otimes w_{la} \# (6)$$

式中, X_w 表示经过纬度加权的特征, \otimes 表示逐元素相乘。随后, 对 X_w 执行全局平均池化, 将空间维度压缩为 1×1 :

$$Z = GAP(X_w) \in \mathbf{R}^{B \times C \times 1 \times 1} \# (7)$$

式中, Z 表示压缩后的特征, GAP 表示全局平均池化。但在进入全局平均池化之前, 先对特征图按纬度加权, 这是因为在 ERP 投影下, 两极像素代表的球面区域远小于赤道像素, 若直接做全局平均池化, 极区冗余信息会过度主导通道统计。接下来, 通过两层 1×1 卷积与 ReLU 激活, 构建通道级的语义重标定, 再经 Sigmoid 归一化生成全局语义注意力向量 $A_g \in [0, 1]^{B \times C \times 1 \times 1}$ 。最后将 A_g 沿空间维度广播回 (H, W) , 得到与原始特征尺寸对齐的全局上下文图 $A_c \in \mathbf{R}^{B \times C \times H \times W}$ 。以上处理过程可以表示如下:

$$A_c = \text{expand} \left(\text{sigmoid} \left(\text{Conv}_{1 \times 1} \left(\text{ReLu} \left(\text{Conv}_{1 \times 1} Z \right) \right) \right) \right) \# (8)$$

式中, A_c 表示计算得到的全局上下文图, $\text{expand}(\cdot)$ 表示将 $(B, C, 1, 1)$ 的注意力向量广播到 (B, C, H, W) , ReLu 表示激活函数, $\text{Conv}_{1 \times 1}$ 表示 1×1 的卷积,

sigmoid 表示 sigmoid 激活函数。

将三种上下文注意力图进行逐像素相乘融合, 形成统一的上下文增强图 A_c :

$$A_c = A_{lo} \otimes A_{la} \otimes A_c \# (9)$$

式中, A_{lo} 、 A_{la} 、 A_c 分别表示上文中计算得到的经度注意力图、纬度注意力图以及全局上下文图。

随后将其与原始特征进行逐通道加权相乘并拼接:

$$X' = \text{concat}(X \otimes A_c, X) \# (10)$$

式中, X' 表示经过全景上下文聚合模块最后得到特征。

2.3 注意力引导融合模块AFM

本文为了实现RGB与深度模态之间的高效融

合, 设计了注意力引导融合模块AFM, 结构如图4所示, 该模块结合通道-空间注意力机制以及动态权重分配策略, 以实现跨模态特征的有效增强与融合。具体而言, 对于输入的RGB特征 F' 和深度特征 D , 首先通过通道注意力与空间注意力机制增强其表示能力。通道注意力机制建模了特征图不同通道在全局语义上的重要性, 同时对质量较差或噪声较多的深度通道进行自适应抑制, 通道注意力图 M_c 的计算过程如下:

$$M_c = \text{sigmoid}(\text{Conv}_{1 \times 1} c(\text{Conv}_{1 \times 1} \text{GAP}(F'))) \# (11)$$

$$\check{F} = F' \otimes M_c \# (12)$$

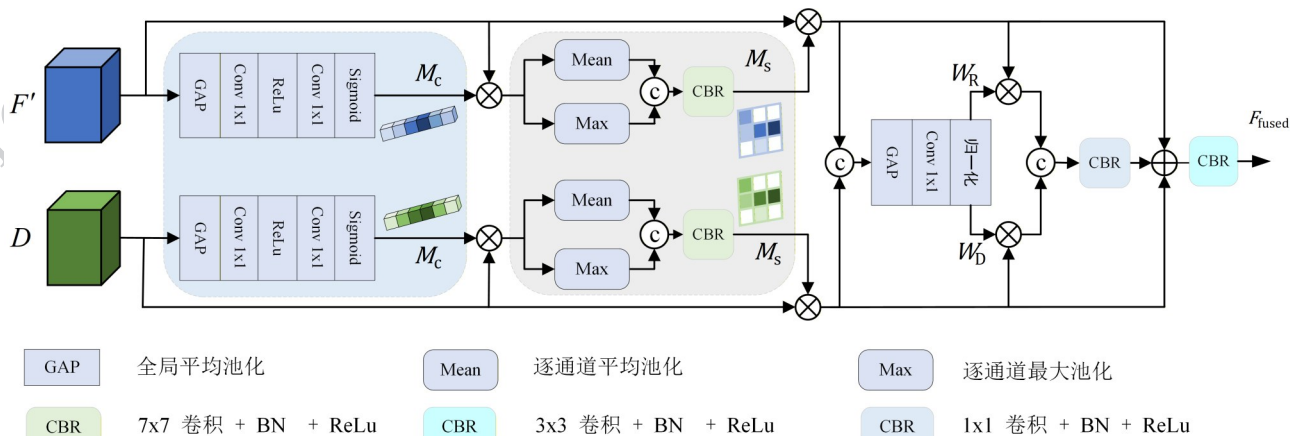


图4 注意力引导融合模块AFM结构图

Fig. 4 Structure of the Attention-Guided Fusion Module (AFM)

式中, M_c 表示得到的通道注意力图, F' 表示经过全景感知感受野模块模块处理后的特征, \check{F} 表示进行将 F' 与通道注意力图进行逐通道相乘得到的特征,接着,为进一步强调空间上的显著区域,引入空间注意力机制,计算空间注意力图 M_s :

$$M_s = \text{Conv}_{7 \times 7}(\text{concat}[\text{Avg}(\check{F}), \text{Max}(\check{F})]) \# (13)$$

$$\check{F} = F' \otimes M_s \# (14)$$

式中, M_s 表示得到的空间注意力图, $\text{Conv}_{7 \times 7}$ 表示7x7的卷积, Avg 表示逐通道平均池化操作, Max 表示逐通道池化操作, \check{F} 表示通过通过与空间注意力图进行逐元素相乘得到的特征,为实现模态间的信息选择性融合,本文设计了一个轻量级动态模态权重生成器,用于生成模态级别的权重向量:

$$\check{F}_c = \text{concat}(\check{F}_R, \check{F}_D) \# (15)$$

$$W = \text{Softmax}(\text{Conv}_{1 \times 1}(\text{GAP}(\check{F}_c))) \# (16)$$

式中, \check{F}_c 表示将经过通道和空间注意力得到的RGB和深度信息特征进行融合后的特征, \check{F}_R 表示RGB分支的注意力图, \check{F}_D 表示深度分支得到的注意力图。

$$F_{out} = \text{concat}(W_R \otimes \check{F}_R, W_D \otimes \check{F}_D) \# (17)$$

式中, W_R 和 W_D 表示 W 经通道拆分得到模态特定权重, F_{out} 表示两个模态的最终融合结果。

在融合特征生成后,为进一步提升特征表达能力,本文引入了残差连接与特征增强机制。将 F_{out} 送入一个1x1卷积层,得到增强特征。随后,与两个原始模态的注意力增强特征进行残差融合,并通过3x3卷积输出最终结果。

综上所述,AFM模块通过引入双重注意力机制
© 中国图象图形学报版权所有

与动态模态融合策略,在保持结构简洁的同时有效增强了多模态特征的表达与融合能力。

2.4 跨模态引导协同解码策略

在完成全景感知与跨模态融合设计后,核心挑战在于如何高效将多尺度特征从编码端引导至解码端以精细还原显著区域。为此,本文提出一种跨模态协同解码策略(图5),以融合分支为主干,动态汲取RGB分支的高分辨率细节信息。该策略通过逐

层注意力加权实现信息重塑与多尺度一致性,同时引入独立的RGB解码路径缓解融合特征中RGB语义弱化问题。具体而言,每层融合特征先与上层上采样特征拼接并经过两次Conv-BN-ReLU生成初级解码特征 X_i ,再与RGB解码特征 \bar{F}_i 拼接,通过轻量注意力生成通道加权系数 A ,得到高级解码特征 Y_i 并进行特征融合,最后经 3×3 Conv-BN-ReLU生成该层最终输出 S_i 。该设计有效增强了RGB特征与融

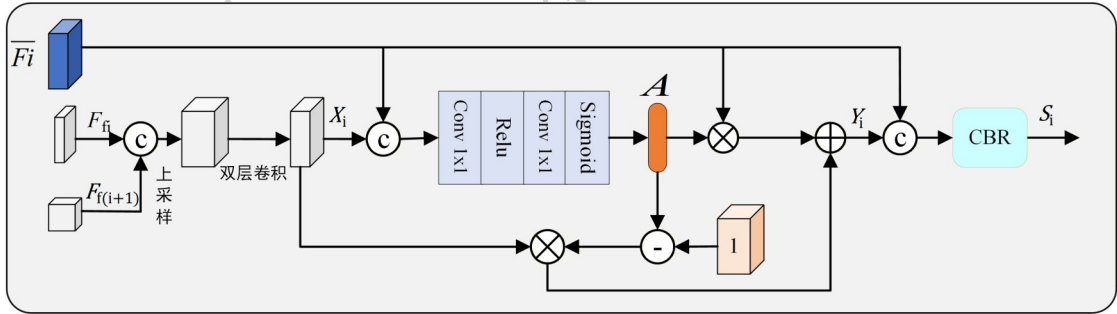


图5 跨模态引导协同解码策略结构图

Fig. 5 Structure of the Cross-Modal Guided Collaborative Decoding Strategy

合特征间的互补性,提升了显著区域的结构保留与边缘细节恢复能力,以上处理过程可以表示如下:

$$X_i = \sigma \left(\text{concat} \left(F_{f_i}, \text{upsampling} \left(F_{f(i+1)} \right) \right) \right), i = 1, 2, 3 \quad (18)$$

$$A = \text{sigmoid} \left(\text{Conv}_{1 \times 1} \text{ReLU} \left(\text{Conv}_{1 \times 1} \text{concat} \left(X_i, \bar{F}_i \right) \right) \right) \# \quad (19)$$

$$Y_i = A \otimes \bar{F}_i + (1 - A) \otimes X_i \# \quad (20)$$

$$S_i = \varnothing \left(\text{Conv}_{3 \times 3} \left(\text{concat} \left[Y_i, \bar{F}_i \right] \right) \right) \# \quad (21)$$

式中, F_{f_i} 表示经过融合模块得到的特征, $F_{f(i+1)}$ 表示上一层得到的融合特征, A 表示通道加权系数, Y_i 与 X_i 表示计算过程中的中间结果, \bar{F}_i 表示RGB解码特征, $i = 1, 2, 3$, σ 表示进行两次的 $\text{Conv}_{3 \times 3} - \text{Batch Normalization} - \text{ReLU}$ 操作, upsampling 表示双线性插值上采样。

最后,为RGB最底层特征 S_4 和每一层融合解码特征 S_1, S_2, S_3 附加侧输出监督,采用 1×1 卷积层进行显著图生成,并通过双线性插值上采样操作恢复至原图尺度以便于训练与评估。最终输出由最底层融合解码特征 S_1 生成,表示最终的预测结果。

2.5 损失函数

本文的损失函数为加权的联合交集(IOU)损失 L_{wIoU} 和加权的二进制交叉熵(BCE)损失 L_{wBCE} 组合,定义为:

$$L = L_{wIoU} + L_{wBCE} \# \# \quad (22)$$

为了增强训练有效性,对RGB分支解码的顶层输出和融合解码的所有输出应用深度监督,从而得到最终的总损失函数:

$$L_{total} = L_{RGB}(\mathbf{G}, \mathbf{S}) + \sum_{i=1}^3 L_f(\mathbf{G}, \mathbf{S}_i) \# \quad (23)$$

式中, \mathbf{G} 表示真实标签, \mathbf{S} 表示RGB分支解码的顶层输出, \mathbf{S}_i 表示由融合解码生成的多级分割输出, L_{total} 表示最终损失。

3 实验与分析

3.1 数据集及评价指标

本文使用360-SOD(Li等,2019)和360-SSOD(Ma等,2020)两个公开的 360° 全景图像数据集来测试本文模型的性能,将Yun等人(2022)提出的 360° 全景图像深度估计方法为现有的2个 360° 全景图像数据集生成了深度图子集。扩充后的360-SOD包含500张高分辨率ERP全景图像及其对应的深度图,

360-SSOD包含1105张高分辨率ERP全景图像及其对应的深度图。本文使用以下几个评价指标来对所有模型的性能进行评估:平均绝对误差(mean absolute error, MAE) (Perazzi 等, 2012)、F-measure (Achanta 等, 2009)、E-measure (Fan 等, 2009)和结构测度(structure-measure,) (Fan 等, 2017)。

3.2 实验设置

本文使用Pytorch来训练所提出的模型,所有实验均在NVIDIA RTX 3090 GPU服务器上完成。编码器采用经过预训练的SAM2 Hiera-L模型和ResNet-34,为避免大型预训练骨干导致算法对比时的不公平性,本文冻结了Hiera的部分参数,在每个多尺度块之前插入了轻量级适配器(Xiong 等, 2024)。输入包括RGB图像与对应的深度图像,图像分辨率为512×1024。本文使用AdamW优化器进行训练,最大学习率设置为1e-3,权重衰减设置为5e-4,为防止过拟合提供正则化约束。训练过程中,引入了余弦退火策略,将学习率逐步衰减至最小值1e-7,衰减周期设置为整个训练周期数(即每50个epoch衰减一次),batch size设置为4,最大训练周期数设为50。训练过程中采用第3.5节定义的结构化损失(加权BCE+加权IoU)并结合多层深度监督,以提升模型收敛与特征学习效果。实验中数据集的划分如下:360-SOD数据集中400张RGB图像及对应深度图用于训练,100张图像及对应深度图用于测试;360-SSOD中850张图像及对应深度图用于训练,255张图像及对应深度图用于测试。训练过程中,使用随机垂直和水平翻转来对数据集进行增强。

3.3 实验对比

本文将所提出的模型与10种代表性的RGB-D 2D SOD方法和7种代表性的RGB 360° SOD方法进行了对比。RGB-D 2D SOD方法包括VST(Liu 等, 2021)、TriTransNet(Liu 等, 2021)、BTS-Net(Zhang 等, 2021)、SwinNet(Liu 等, 2021)、MobileSal(Wu 等, 2021)、SPSN(Lee 等, 2022)、PICR-Net(Cong 等, 2023)、HIDANet(Wu 等, 2023)、DFormer-B(Yin 等, 2023)、CPNet(Hu 等, (2024)),上述10种方法使用360-SOD和360-SSOD数据集中的RGB图像和本文生成的深度估计图像进行训练。RGB 360° SOD方法包括FANet(Huang 等, 2020)、MPFR-Net(Cong 等, 2023)、LDNet(Huang 等, 2023)、DATFormer(Zhao 等, 2023)、ACoNet Chen 等, 2024)、DPNet(Chen 等,

2025), DSANet(Chen 等, 2025),这些方法只使用360-SOD和360-SSOD数据集中的RGB图像进行训练。为了公平比较,所有模型均使用官方代码或者作者提供的显著图在同一软硬件环境中测试,并且都做了微调来获得最好结果。

所有方法在360-SOD数据集上的客观指标对比如表1所示,可以看到本文方法相较于现有的RGB-D 2D SOD方法和RGB 360° SOD方法取得了最好的性能,尤其相对于RGB 360° SOD方法,本文方法展示出了较大程度的提高。相对于所有对比方法中的次优值,本文方法的MAE指标降低了13.7%,max-F提高了4.86%,mean-F提高了3.24%,max-Em提高了0.42%,mean-Em提高了0.39%,Sm提高了1.96%

在360-SOD数据集上的主观结果比较如图6所示,本文选取了若干具有挑战性和代表性的场景进行比较。图6是本文方法与10种代表性RGB-D 2D SOD方法和5种代表性RGB 360° SOD方法的对比,从图中可以看出,本文方法生成的显著性目标图比其他方法更接近真值图。例如:第1行和第7行的场景中,对于存在严重畸变的极区,只有本文方法正确检测到了极区物体,反映出本文方法对于严重畸变区域的检测效果是优秀的;在具有复杂背景的场景中,例如第3、4、5、8行,本文方法能够准确识别目标主体和边界并完成分割,优于其他方法;在投影边界不连续场景中,例如第2行,ERP投影导致图像边界处的显著目标出现结构断裂,本文方法在该场景下的检测结果优于现有方法;另外,当图像中存在多个检测目标时,例如第6行场景,本文方法能够精准定位所有显著目标并实现精细分割,而其他方法则存在部分显著目标遗漏的情况。

所有方法在360-SSOD数据集的客观指标对比如表2所示。整体来看,本文方法综合性能最好,各项指标上的表现同样具有竞争力。

360-SSOD数据集上的主观结果比较如图7所示,从图7中可以看出与目前代表性RGB-D 2D SOD先进方法和RGB 360° SOD先进方法相比,本文方法在畸变极区的(例如图7第2行图像)检测效果优于其他方法,另外对低对比度场景(例如图7第1、5行中的图像)和具有复杂背景的场景(例如图7第3、6、7行中的图像),本文方法能够有效抑制背景或突出前景,实现了更精细的显著目标检测。

表1 不同方法在360-SOD数据集上的客观指标对比

Table 1 Comparison of Objective Metrics of Different Methods on the 360-SOD Dataset

方法	MAE ↓	max-F ↑	mean-F ↑	max-Em ↑	mean-Em ↑	Sm ↑
VST(2021)	0.0229	0.7353	0.6919	0.8846	0.8537	0.8147
TriTransNet(2021)	0.0217	0.7740	0.7321	0.9026	0.8970	0.8204
BTS-Net(2021)	0.0240	0.7230	0.6853	0.8780	0.8032	0.7972
SwinNet(2022)	0.0283	0.7063	0.6661	0.8640	0.7932	0.7783
MobileSal(2022)	0.0268	0.6136	0.5973	0.8222	0.7827	0.7420
SPSN(2022)	0.0193	0.7748	0.7584	0.8938	0.8779	0.8314
PICR-Net(2023)	0.0204	0.7734	0.7583	0.9005	0.8943	0.8306
HIDANet(2023)	0.0191	0.7862	0.7730	0.9086	0.9039	0.8473
DFormer-B(2024)	0.0175	0.7999	0.7758	0.9292	0.9204	0.8568
CPNet(2024)	0.0187	0.7942	0.7894	0.9039	0.8946	0.8405
FANet(2020)	0.0249	0.6846	0.6628	0.8630	0.8225	0.7697
MPFRNet(2023)	0.0191	0.7651	0.7549	0.8848	0.8744	0.8418
LDNet(2023)	0.0289	0.6562	0.6391	0.8655	0.8414	0.7679
DATFormer(2023)	0.0186	0.7647	0.7490	0.8878	0.8773	0.8408
ACoNet(2024)	0.0181	0.7893	0.7815	0.9141	0.9043	0.8493
DPNet(2025)	0.0190	0.8016	0.7864	0.9207	0.9096	0.8486
DSANet(2025)	0.0198	0.7842	0.7703	0.9081	0.8988	0.8463
RGB-D 2D SOD方法						
RGB 360° SOD方法						
RGB-D 360° SOD方法						
本文方法	0.0151	0.8388	0.8150	0.9331	0.9240	0.8736

注:黑色加粗字体表示最优结果,↑(↓)表示客观指标数值越大越好(越小越好)

图8和图9展示了不同模型在360-SOD和360-SSOD数据集上的P-R曲线和F-measure曲线,从图中可以看到,本文模型(红实线)在两个数据集上的表现均为最优,进一步证明了本文模型的优越性能。

3.4 消融实验

为本节将通过删除或者替换本文模型的不同组件来研究它们的有效性,主要验证方法与组件为:不同深度估计方法的影响分析、全景感知感受野模块PA-RFM的有效性、注意力引导融合模块AFM的有效性以及跨模态引导协同解码策略的有效性。本文在360-SOD数据集上分别对以上方法和组件做了消融实验。

3.4.1 不同深度估计方法的影响分析

在相同的网络结构(均使用SAM2-Hiera作为编码器)、训练配置与数据预处理条件下,对四种主流的360°深度估计方法(Joint_360depth(Yun等,

2022),本文采用的方法)、DA(Wang等,2024)、EGFormer(Yun等,2023)和CRF360D(Cao等,2024))分别生成深度图并重新训练模型;同时,以仅使用RGB的w/o depth方案作为对照。表3的结果显示:无论采用哪一种深度估计方法,将深度模态引入模型后,各项指标均显著优于单模态(RGB-only),并且全部优于表1中未使用深度信息的RGB 360° SOD方法。这说明深度特征能够稳定辅助全景显著目标检测,从整体上提升检测精度和边缘完整度。同时值得注意的是,在RGB-only情况下,本文模型的性能仅处于现有RGB 360° SOD方法的中等偏上水平,说明SAM2-Hiera的表征能力并不足以在缺乏深度模态时解决全景投影畸变和结构歧义等核心难题,这证明本文模型的优越性主要来自于网络结构和各子模块设计,而非由大型预训练骨干网络带来。

3.4.2 全景感知感受野模块PA-RFM的有效性

为验证本文提出的全景感知感受野模块PA-

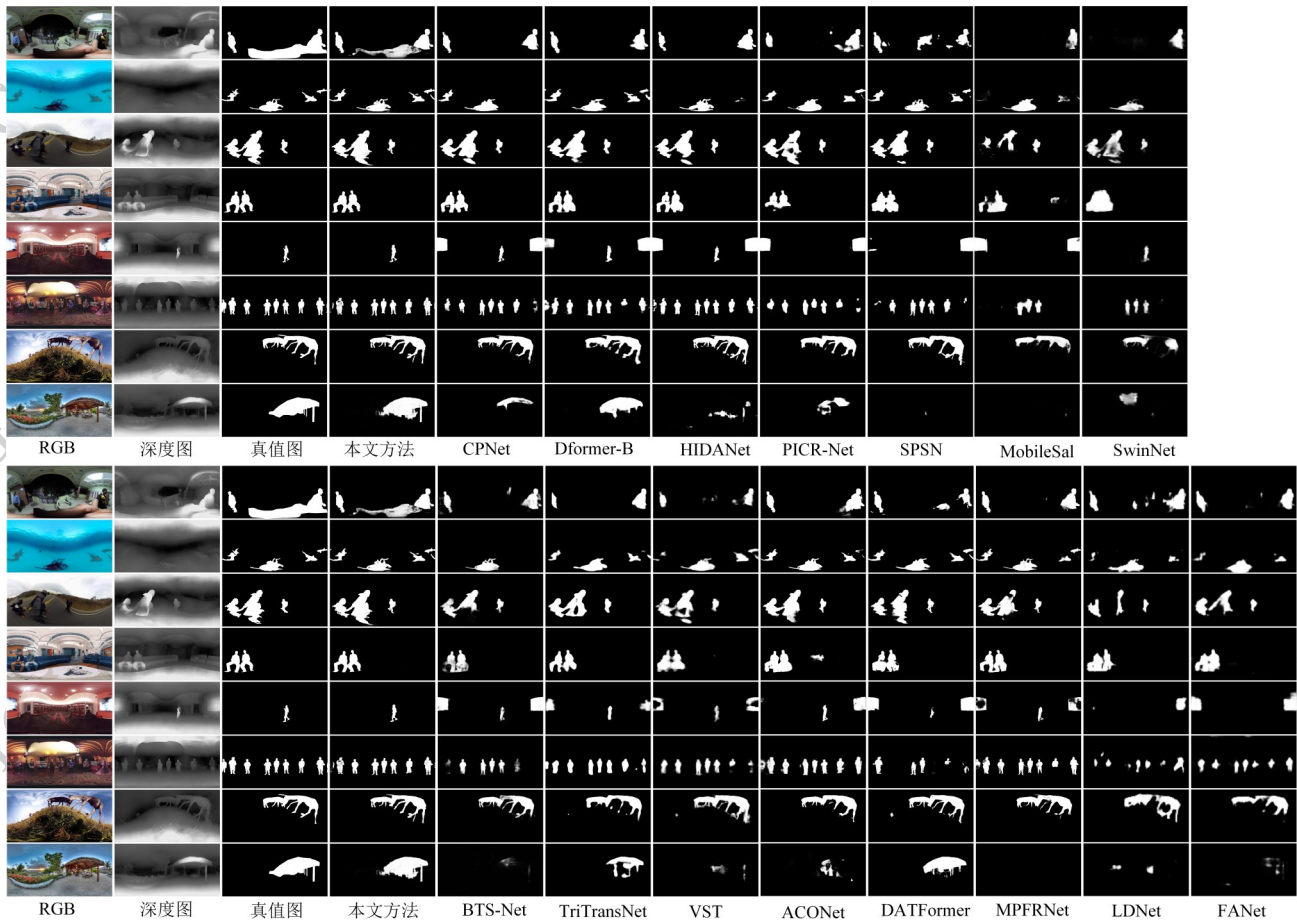


图6 不同方法在360-SOD数据集上的主观结果比较

Fig. 6 Qualitative comparison of different methods on the 360-SOD dataset

RFM 的有效性, 设计了三种对比方案: 1) with RFB: 采用经典感受野模块 RFB (Liu 等, 2018) 替换 PA-RFM; 2) without PA-RFM: 以 1×1 卷积替代 PA-RFM, 即完全去除该模块; 3) w/ PA-RFM (ours): 使用本文提出的 PA-RFM 模块。

客观指标对比表 4 可见, 加入 PA-RFM 后各项指标均达到最优, 较其他替代方案表现更佳。主观可视化结果 (图 10) 进一步验证了该模块的有效性: 第 2 行显示其在极区畸变区域具有更准确的显著性定位; 第 3 行则表明其在处理跨边界连续物体时能更好保持目标完整性。同时对全景上下文聚合模块 (PCAM) 及其三类注意力机制 (经度注意 Lon、纬度注意 Lat、全局注意 Glob) 进行了消融实验。以去除 PCAM 的完整模型 (No. 1) 为基线, 逐步加入不同注意力分支, 其结果如表 5 所示。单独加入任一注意力机制 (No. 2 - No. 4) 均带来显著提升; 同时使用两种机制 (No. 5 - No. 7) 可获得更优表现, 说明它们具备互补性; 三种机制同时使用时 (No. 8) 在全部指标

上达到最佳, 证明三类注意力存在协同作用, 可有效增强全景特征建模与显著性检测性能。

3.4.3 注意力引导融合模块 AFM 的有效性

为验证所提出的注意力引导融合模块 AFM 的有效性, 本文在统一网络架构下设计了五组对比实验, 结果如表 6 所示: 表中 No. 1 代表单独使用 RGB 模态 (RGB); No. 2 代表 RGB 模态与深度模态逐元素相加 (Add) 融合; No. 3 代表 RGB 模态与深度模态通道拼接 (Concat) 融合; No. 4 代表 RGB 模态与深度模态逐元素乘法 (Multi) 融合; No. 5 代表使用本文提出 AFM 模块。由表 6 实验结果可以看出, 相较于 RGB 单模态, 加入深度模态后, Add、Concat、Multi、AFM 四种双模态模型性能均有显著提升, 这充分证明了深度信息对 360 RGB-D SOD 检测具有重要价值, 即便在采用最简单的融合策略时亦能显著提升性能。进一步地, 采用本文提出的 AFM 模块融合 RGB 与深度特征时, 性能在所有指标上均优于其他融合方式, 说明 AFM 能够更充分地挖掘跨模态的互补性, 从而

表2 360-SSOD数据集上各模型客观指标对比

Table 2 Comparison of Objective Metrics of Different Models on the 360-SSOD Dataset

方法	MAE ↓	max-F ↑	mean-F ↑	max-Em ↑	mean-Em ↑	Sm ↑	
VST(2021)	0.0280	0.6189	0.5900	0.8511	0.7979	0.7583	
TriTransNet(2021)	0.0283	0.6263	0.6225	0.8578	0.8372	0.7451	
BTS-Net(2021)	0.0340	0.6190	0.5643	0.8289	0.7556	0.7447	
SwinNet(2022)	0.0292	0.6504	0.6217	0.8622	0.8160	0.7689	
RGB-D 2D SOD 方法	MobileSal(2022)	0.0324	0.5078	0.4891	0.7652	0.7016	0.6873
DFM-Net(2022)	0.0407	0.05747	0.5084	0.8555	0.7433	0.7190	
SPSN(2022)	0.0283	0.6379	0.6288	0.8446	0.8109	0.7383	
PICR-Net(2023)	0.0299	0.6099	0.6013	0.8509	0.8247	0.7320	
HIDANet(2023)	0.0281	0.6613	0.6409	0.8623	0.8547	0.7782	
DFormer-B(2024)	0.0273	0.6762	0.6603	0.8758	0.8482	0.7846	
CPNet(2024)	0.0278	0.6752	0.6676	0.8778	0.8649	0.7748	
FANet(2020)	0.0415	0.5619	0.5335	0.8027	0.7703	0.7119	
MPFRNet(2023)	-	-	-	-	-	-	
RGB 360° SOD方法	LDNet(2023)	0.0342	0.5862	0.5672	0.8390	0.8187	0.7245
DATFormer(2023)	0.0252	0.6196	0.6050	0.8325	0.7920	0.7527	
ACoNet(2024)	0.0288	0.6641	0.6564	0.8695	0.8632	0.7796	
RGB-D 360° SOD方法	本文方法	0.0273	0.6938	0.6706	0.8759	0.8564	0.7871

注:黑色加粗字体表示最优结果,↑(↓)表示客观指标数值越大越好(越小越好),“-”表示因没有源代码而无法测试的值。

获得更高的检测精度与鲁棒性。

为了进一步证明AFM的优越性,将AFM替换为其他SOTA方法的特征融合模块,包括MobileSal(Wu等,2021),HIDANet(Wu等,2023)和CPNet(Hu等,(2024))。如表7所示,AFM模块比其他方法中使用的模块性能更好。

3.4.4 跨模态引导协同解码策略的有效性

为验证所提出的跨模态引导协同解码策略的有效性,将完整模型(Ours)与移除RGB解码路径的变体(w/o RGB)进行对比,结果如表8所示。可以看到,去掉RGB解码路径后,实验结果显示模型在多项评测指标上均有轻微下降,说明在缺少RGB引导时模型的边缘/细节恢复以及整体结构保持能力均受到影响。且由图12所展示的主观对比可以看出,RGB解码路径能够在融合分支解码过程中提供高分辨率的结构先验和细节补偿,有效缓解融合特征中RGB语义弱化问题,显著提升模型的边缘区域与细节纹理检测能力。

4 结论

本文针对360°SOD面临的投影畸变、边界不连续及跨模态融合等问题,提出了一种RGB-D 360°SOD方法。该方法包含三个协同工作的核心模块:PA-RFM模块通过显式编码经纬度信息并结合全局注意力,缓解极区畸变与边缘断裂,增强几何感知;AFM模块利用动态权重与双重注意力机制,抑制深度噪声,强化RGB-D互补表征;跨模态协同解码策略引入独立RGB路径作为高分辨率先验,在融合解码中有效恢复细节,克服传统方法的细节弱化问题。在两个数据集上的系统性实验与消融结果表明,相比现有代表性先进方法,本文方法在检测精度和鲁棒性方面均取得了显著提升,证明了深度模态在360°SOD中的应用潜力。需要说明的是,使用深度估计方法生成的深度图与使用深度相机拍摄的深度图还有一定的差距,且不同深度估计方法之间的性

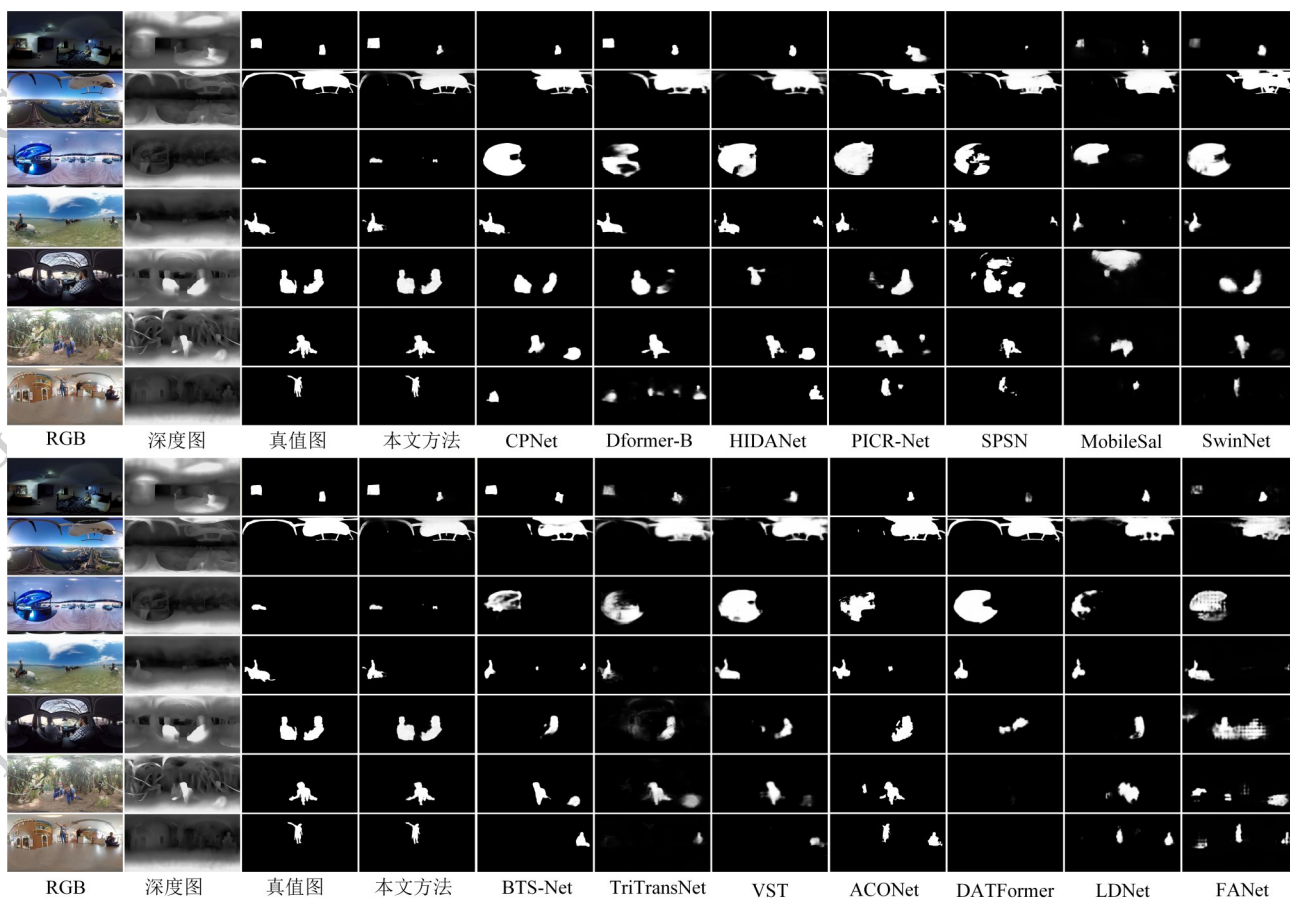


图7 不同方法在360-SSOD数据集上的主观结果比较

Fig. 7 Qualitative comparison of different methods on the 360-SOD dataset

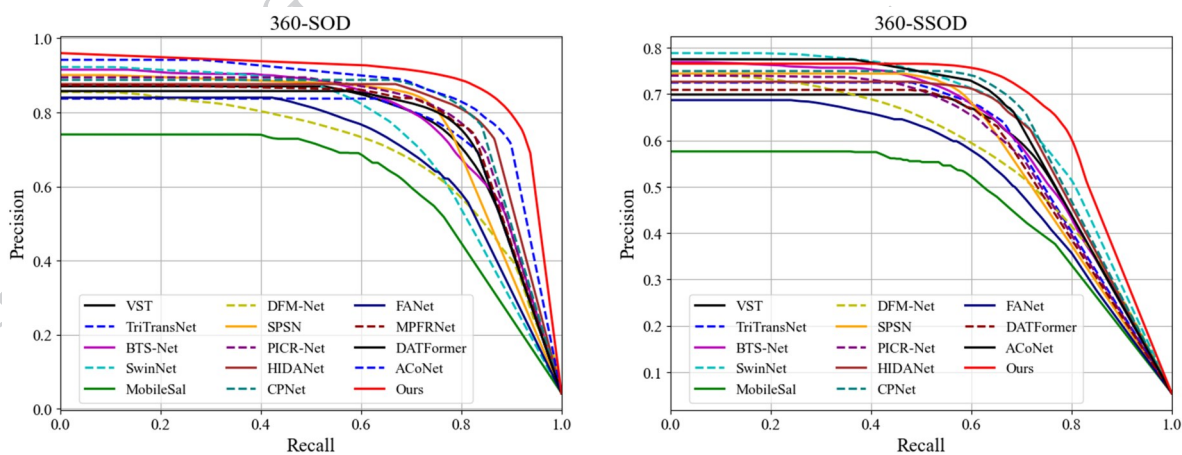


图8 不同模型在360-SOD和360-SSOD数据集上的P-R曲线比较

Fig. 8 Comparison of P-R Curves of Different Models on the 360-SOD and 360-SSOD Datasets

能存在差异,故本文模型的效果可能会因采用的深度图不同而有所变化。但是从大量实验来看,这种变化在可控范围内,不会造成较大程度的检测性能变化。未来研究将从以下几个方向展开:第一,结合真实深度传感器数据开展跨域自适应学习;第二,探

索多投影视角联合建模,以进一步缓解极区畸变问题;第三,研究模型蒸馏与网络压缩策略,实现模型轻量化与移动端友好部署。

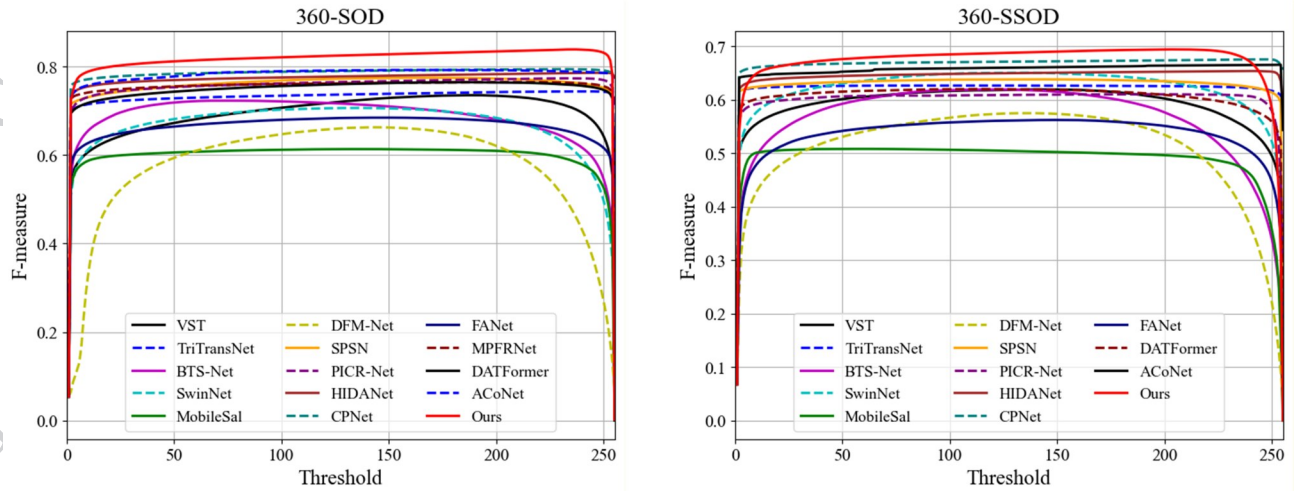


图9 不同模型在360-SOD和360-SSOD数据集上的F-measure曲线比较

Fig. 8 Comparison of P-R Curves and F-measure Curves of Different Models on the 360-SOD and 360-SSOD Datasets

表3 使用不同深度估计方法的客观指标对比

Table 3 Comparison of Objective Metrics Using Different Depth Estimation Methods

模型方法	深度估计方法	MAE ↓	max-F ↑	mean-F ↑	max-Em ↑	mean-Em ↑	Sm ↑
	w/o depth (RGB only)	0.0220	0.7304	0.7021	0.8526	0.8414	0.8096
	DA(2023)	0.0160	0.8300	0.8064	0.9211	0.9154	0.8731
本文模型	EGFormer(2024)	0.0163	0.8243	0.8079	0.9205	0.9095	0.8643
	CRF360D(2024)	0.0171	0.8179	0.7980	0.9204	0.9101	0.8725
	Joint_360depth(2022)	0.0151	0.8388	0.8150	0.9331	0.9240	0.8736
	(本文使用方法)						

注:黑色字体表示最优结果,↑(↓)表示客观指标数值越大越好(越小越好),w/o depth表示不使用深度信息

表4 全景感知感受野模块PA-RFM有效性的客观指标对比

Table 4 Comparison of Objective Metrics for the Effectiveness of the Panoramic-Aware Receptive Field Module (PA-RFM)

Method	MAE ↓	max-F ↑	mean-F ↑	max-E ↑	mean-E ↑	Sm ↑
with RFB	0.0162	0.8173	0.8009	0.9177	0.9073	0.8703
w/o PA-RFM	0.0166	0.8140	0.8011	0.9201	0.9008	0.8610
本文	0.0151	0.8388	0.8150	0.9331	0.9240	0.8736

注:黑色加粗字体表示最优结果,↑(↓)表示客观指标数值越大越好(越小越好),with RFB表示使用RFB模块,w/o PA-RFM表示不适用PA-RFM模块。

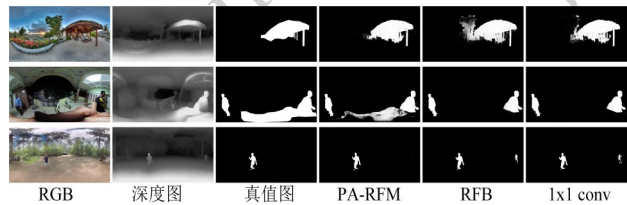


图10 全景感知感受野模块PA-RFM有效性的主观结果对比

Fig. 10 Qualitative comparison demonstrating the effectiveness of the Panoramic-Aware Receptive Field Module

表5 PCAM模块的有效性及其三种注意力不同组合的客观指标对比

Table 5 comparison of Objective Metrics for the Effectiveness of the PCAM Module and Different Combinations of Three Attention Mechanisms

NO.	Lon	Lat	Glob	MAE ↓	max-F ↑	mean-F ↑	max-E ↑	mean-E ↑	Sm ↑
1				0.0215	0.7471	0.7197	0.8720	0.8597	0.8201
2	■			0.0193	0.7966	0.7792	0.8883	0.8763	0.8431
3		■		0.0181	0.8010	0.7839	0.9147	0.9085	0.8568
4			■	0.0186	0.8096	0.7996	0.9274	0.9142	0.8393
5	■	■		0.0166	0.8164	0.7800	0.9248	0.9093	0.8683
6	■		■	0.0171	0.8299	0.8013	0.9340	0.9199	0.8547
7		■	■	0.0174	0.8327	0.8114	0.9293	0.9184	0.8553
8	■	■	■	0.0151	0.8388	0.8150	0.9331	0.9240	0.8736

注:黑色加粗字体表示最优结果, ↑(↓)表示客观指标数值越大越好(越小越好), ■表示使用对应注意力机制。

表6 注意力引导融合模块AFM有效性的客观指标对比

Table 6 Comparison of Objective Metrics for the Effectiveness of the Attention-Guided Fusion Module (AFM)

NO.	RGB	Depth	Fusion method	MAE ↓	max-F ↑	mean-F ↑	max-E ↑	mean-E ↑	Sm ↑
1	■			0.0220	0.7304	0.7021	0.8526	0.8414	0.8096
2	■	■	Add	0.0208	0.7445	0.7128	0.8799	0.8624	0.8185
3	■	■	Concat	0.0157	0.8217	0.8066	0.9185	0.9114	0.8314
4	■	■	Multi	0.0166	0.8145	0.8095	0.9019	0.8814	0.8450
5	■	■	AFM	0.0151	0.8388	0.8150	0.9331	0.9240	0.8736

注:黑色加粗字体表示最优结果, ↑(↓)表示客观指标数值越大越好(越小越好), ■表示使用。

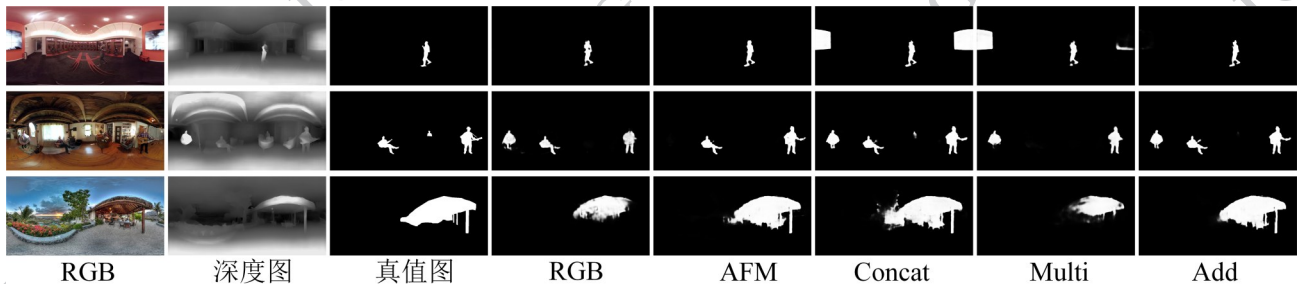


图11 注意力引导融合模块AFM有效性的主观结果对比

Fig. 11 Qualitative comparison demonstrating the effectiveness of the Attention-Guided Fusion Module (AFM)

表7 AFM 模块与其他 SOTA 方法的特征融合模块客观指标对比

Table 7 Comparison of Objective Metrics between the AFM Module and Feature Fusion Modules of Other SOTA Methods

Method	MAE ↓	max-F ↑	mean-F ↑	max-E ↑	mean-E ↑	Sm ↑
MobileSal	0.0159	0.8330	0.8134	0.9260	0.9164	0.8634
HIDANet	0.0163	0.8373	0.8112	0.9323	0.9231	0.8734
CPNet	0.0165	0.8355	0.8100	0.9282	0.9163	0.8663
本文	0.0151	0.8388	0.8150	0.9331	0.9240	0.8736

注:黑色加粗字体表示最优结果,↑(↓)表示客观指标数值越大越好(越小越好)。

表8 跨模态引导协同解码策略有效性的客观指标对比

Table 8 Comparison of Objective Metrics for the Effectiveness of the Cross-Modal Guided Collaborative Decoding Strategy

Method	MAE ↓	max-F ↑	mean-F ↑	max-E ↑	mean-E ↑	Sm ↑
w/o RGB	0.0169	0.8328	0.8051	0.9255	0.9157	0.8724
本文	0.0151	0.8388	0.8150	0.9331	0.9240	0.8736

注:黑色加粗字体表示最优结果,↑(↓)表示客观指标数值越大越好(越小越好),w/o表示移除RGB解码。

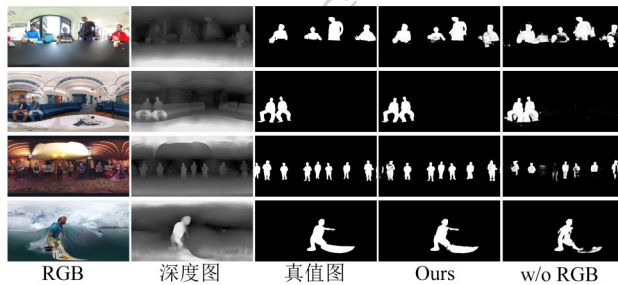


图12 跨模态引导协同解码策略的有效性主观对比

Fig. 11 Qualitative comparison demonstrating the effectiveness of the Cross-Modal Guided Collaborative Decoding Strategy

参考文献 (References)

- Achanta R, Hemami S, Estrada F and Sussstrunk S. 2009. Frequency-tuned salient region detection//Proceedings of the 2009 IEEE Conference on Computer Vision and Pattern Recognition. Miami, USA: IEEE: 1597-1604[DOI:10.1109/CVPR.2009.5206596]
- Cao Z and Wang L. 2025. Monocular 360 depth estimation via spherical fully-connected CRFs. IEEE Robotics and Automation Letters, 10

(2):1409-1416[DOI:10.1109/LRA.2024.3518109]

Chen X L, Zhang X G, Du Z L and Wang X. 2025. Distortion semantic aggregation network for 360° omnidirectional image salient object detection. Journal of Image and Graphics, 30(7):2451-2467 (陈晓雷, 张学功, 杜泽龙, 王兴. 2025. 面向360°全景图像显著目标检测的畸变语义聚合网络. 中国图象图形学报, 30(7):2451-2467)[DOI:10.11834/jig.240371]

Chen X L, Zhang X G, Du Z L and Wang X. 2025. Distortion-adaptive and position-aware 360° omnidirectional image salient object detection network. Journal of Image and Graphics, 30(8):2758-2774 (陈晓雷, 杜泽龙, 张学功, 王兴. 2025. 畸变自适应与位置感知的360°全景图像显著目标检测网络. 中国图象图形学报, 30(8):2758-2774)[DOI:10.11834/jig.240592]

Chen X, Wang X, Zhang X and Du Z. 2024. Adjacent coordination network for salient object detection in 360 degree omnidirectional images. Journal of Electronics and Information Technology, 46(12):4529-4541[DOI:10.11999/JEIT240502]

Ciptadi A, Hermans T and Rehg J M. 2013. An in depth view of saliency//Proceedings of the British Machine Vision Conference. Bristol, United Kingdom: BMVC: 1-11[DOI:10.5244/C.27.76]

Cong R, Huang K, Lei J, Zhao Y, Huang Q and Kwong S. 2024. Multi-projection fusion and refinement network for salient object detection in 360° omnidirectional image. IEEE Transactions on Neural Networks and Learning Systems, 35(7):9495-9507[DOI:10.1109/TNNLS.2022.3233883]

Cong R, Lei J, Fu H, Hou J, Huang Q and Kwong S. 2020. Going from RGB to RGBD saliency: a depth-guided transformation model. IEEE Transactions on Cybernetics, 50(8):3627-3639[DOI:10.1109/TCYB.2019.2932005]

Cong R, Liu H, Zhang C, Zhang W, Zheng F, Song R and Kwong S. 2023. Point-aware interaction and CNN-induced refinement network for RGB-D salient object detection//Proceedings of the 31st ACM International Conference on Multimedia. New York, USA: ACM: 406-416[DOI:10.1145/3581783.3611982]

Donoser M, Urschler M, Hirzer M and Bischof H. 2009. Saliency driven total variation segmentation//Proceedings of the 2009 IEEE 12th International Conference on Computer Vision. Kyoto, Japan: IEEE: 817-824[DOI:10.1109/ICCV.2009.5459296]

Dosovitskiy A, Beyer L, Kolesnikov A, Weissenborn D and Houlsby N. 2020. An image is worth 16x16 words: transformers for image recognition at scale. arXiv preprint arXiv:2010.11929[DOI:10.48550/arXiv.2010.11929]

Fan D P, Gong C, Cao Y, Ren B, Cheng M M and Borji A. 2018. Enhanced-alignment measure for binary foreground map evaluation. arXiv preprint arXiv:1805.10421 [DOI:10.48550/arXiv.1805.10421]

Fan D P, Gong C, Cao Y, Ren B, Cheng M M and Borji A. 2018. Enhanced-alignment measure for binary foreground map evaluation//Proceedings of the 27th International Joint Conference on Arti-

- ficial Intelligence. Stockholm, Sweden: AAAI Press: 698-704 [DOI:10.5555/3304415.3304515]
- He K, Zhang X, Ren S and Sun J. 2016. Deep residual learning for image recognition//Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition. Las Vegas, USA: IEEE: 770-778[DOI:10.1109/CVPR.2016.90]
- Hu X, Sun F, Sun J, Wang F and Li H. 2024. Cross-modal fusion and progressive decoding network for RGB-D salient object detection. *International Journal of Computer Vision*, 132: 3067-3085 [DOI: 10.1007/s11263-024-02020-y]
- Huang M, Li G, Liu Z and Zhu L. 2023. Lightweight distortion-aware network for salient object detection in omnidirectional images. *IEEE Transactions on Circuits and Systems for Video Technology*, 33(10):6191-6197[DOI:10.1109/TCSVT.2023.3253685]
- Huang M, Liu Z, Li G, Zhou X and Le Meur O. 2020. FANet: features adaptation network for 360° omnidirectional salient object detection. *IEEE Signal Processing Letters*, 27: 1819-1823 [DOI: 10.1109/LSP.2020.3028192]
- Läng C, Nguyen T V, Katti H, Yadati K, Kankanhalli M and Yan S. 2012. Depth matters: influence of depth cues on visual saliency//Proceedings of the 12th European Conference on Computer Vision. Berlin, Heidelberg: Springer: 101-115 [DOI: 10.1007/978-3-642-33709-3_8]
- Lee M, Park C, Cho S and Lee S. 2022. SPSN: superpixel prototype sampling network for RGB-D salient object detection//Proceedings of the European Conference on Computer Vision. Cham: Springer: 630-647[DOI:10.1007/978-3-031-19818-2_36]
- Li J, Su J, Xia C and Tian Y. 2020. Distortion-adaptive salient object detection in 360° omnidirectional images. *IEEE Journal of Selected Topics in Signal Processing*, 14(1): 38-48 [DOI: 10.1109/JSTSP.2019.2957982]
- Li J, Su J, Xia C and Tian Y. 2020. Distortion-adaptive salient object detection in 360° omnidirectional images. *IEEE Journal of Selected Topics in Signal Processing*, 14(1): 38-48 [DOI: 10.1109/JSTSP.2019.2957982]
- Liu N, Zhang N, Wan K, Han J and Shao L. 2021. Visual saliency transformer//Proceedings of the IEEE/CVF International Conference on Computer Vision. Montreal, Canada: IEEE: 4722-4732 [DOI:10.1109/ICCV48922.2021.00468]
- Liu S T, Huang D and Wang Y H. 2018. Receptive field block net for accurate and fast object detection//Proceedings of the 15th European Conference on Computer Vision. Munich, Germany: Springer: 404-419[DOI:10.1007/978-3-030-01252-6_24]
- Liu Z, Lin Y, Cao Y, Hu H, Wei Y, Zhang Z, Lin S and Guo B. 2021. Swin transformer: hierarchical vision transformer using shifted windows//Proceedings of the IEEE/CVF International Conference on Computer Vision. Montreal, Canada: IEEE: 9992-10002[DOI:10.1109/ICCV48922.2021.00986]
- Liu Z, Tan Y, He Q and Xiao Y. 2022. SwinNet: swin transformer drives edge-aware RGB-D and RGB-T salient object detection. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(7):4486-4497[DOI:10.1109/TCSVT.2021.3127149]
- Liu Z, Wang Y, Tu Z, Xiao Y and Tang B. 2021. TriTransNet: RGB-D salient object detection with a triplet transformer embedding network//Proceedings of the 29th ACM International Conference on Multimedia. New York, USA: ACM: 4481-4490 [DOI: 10.1145/3474085.3475601]
- Perazzi F, Krähenbühl P, Pritch Y and Hornung A. 2012. Saliency filters: contrast based filtering for salient region detection//Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition. Providence, USA: IEEE: 733-740 [DOI: 10.1109/CVPR.2012.6247743]
- Ravi N, Gabeur V, Hu Y T, Bolya D, Wei C, Fan H, Huang P Y and Feichtenhofer C. 2024. Sam 2: segment anything in images and videos. arXiv preprint arXiv:2408.00714 [DOI: 10.48550/arXiv.2408.00714]
- Ren J, Gong X, Yu L, Zhou W and Yang M Y. 2015. Exploiting global priors for RGB-D saliency detection//Proceedings of the 2015 IEEE Conference on Computer Vision and Pattern Recognition Workshops. Boston, USA: IEEE: 25-32 [DOI: 10.1109/CVPRW.2015.7301391]
- Ryali C, Hu Y T, Bolya D, Wei C, Fan H, Huang P Y and Feichtenhofer C. 2023. Hiera: a hierarchical vision transformer without the bells-and-whistles//Proceedings of the 40th International Conference on Machine Learning. Honolulu, USA: PMLR: 29441-29454
- Wang N H and Liu Y L. 2025. Depth anywhere: enhancing 360 monocular depth estimation via perspective distillation and unlabeled data augmentation//Proceedings of the 38th International Conference on Neural Information Processing Systems. New York, USA: Curran Associates Inc: 127739-127764 [DOI:10.5555/3737916.3741972]
- Wang W, Shen J, Sun H and Shao L. 2018. Video co-saliency guided co-segmentation. *IEEE Transactions on Circuits and Systems for Video Technology*, 28(8):1727-1736 [DOI:10.1109/TCSVT.2017.2701279]
- Wu Y H, Liu Y, Xu J, Bian J W, Gu Y C and Cheng M M. 2022. MobileSal: extremely efficient RGB-D salient object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(12):10261-10269 [DOI:10.1109/TPAMI.2021.3134684]
- Wu Z, Allibert G, Meriaudeau F, Ma C and Demonceaux C. 2023. HiDAnet: RGB-D salient object detection via hierarchical depth awareness. *IEEE Transactions on Image Processing*, 32:2160-2173 [DOI:10.1109/TIP.2023.3263111]
- Xiong X, Wu Z, Tan S, Liu Y and Wang X. 2024. Sam2-unet: segment anything 2 makes strong encoder for natural and medical image segmentation. arXiv preprint arXiv:2408.08870 [DOI:10.48550/arXiv.2408.08870]
- Yin B, Zhang X, Li Z, Liu L, Cheng M M and Hou Q. 2023. Dformer: rethinking RGBD representation learning for semantic segmenta-

- tion. arXiv preprint arXiv: 2309.09668 [DOI: 10.48550/arXiv.2309.09668]
- Yun I, Lee H J and Rhee C E. 2022. Improving 360 monocular depth estimation via non-local dense prediction transformer and joint supervised and self-supervised learning//Proceedings of the 36th AAAI Conference on Artificial Intelligence. Palo Alto, USA: AAAI Press: 3224-3233 [DOI: 10.1609/aaai.v36i3.20231]
- Yun I, Shin C, Lee H, Lee H J and Rhee C E. 2023. EGformer: equirectangular geometry-biased transformer for 360 depth estimation//Proceedings of the 2023 IEEE/CVF International Conference on Computer Vision. Paris, France: IEEE: 6078-6089 [DOI: 10.1109/ICCV51070.2023.00561]
- Zhang W, Jiang Y, Fu K and Zhao Q. 2021. BTS-Net: bi-directional transfer-and-selection network for RGB-D salient object detection//Proceedings of the 2021 IEEE International Conference on Multimedia and Expo. Shenzhen, China: IEEE: 1-6 [DOI: 10.1109/ICME51207.2021.9428263]
- Zhang Y, Qian X, Tan X, Han J and Tang Y. 2016. Sketch-based image retrieval by salient contour reinforcement. IEEE Transactions on Multimedia, 18 (8) : 1604-1615 [DOI: 10.1109/TMM. 2016. 2568138]
- Zhao Y, Zhao L, Yu Q, Sheng L, Zhang J and Xu D. 2023. Distortion-aware transformer in 360° salient object detection//Proceedings of the 31st ACM International Conference on Multimedia. New York, USA: ACM: 499-508 [DOI: 10.1145/3581783.3612025]
- Zhu S, Liu C and Xu Z. 2020. High-definition video compression system based on perception guidance of salient information of a convolutional neural network and HEVC compression domain. IEEE Transactions on Circuits and Systems for Video Technology, 30 (7) : 1946-1959 [DOI: 10.1109/TCSVT.2019.2911396]

作者简介

陈晓雷, 通讯作者, 男, 教授, 硕士生导师, 主要研究方向为人工智能与计算机视觉。E-mail: chenxl703@lut.edu.cn

钟智华, 男, 硕士研究生, 主要研究方向为全景图像显著性目标检测。E-mail: 1603524098@qq.com

申玉杰, 女, 硕士研究生, 主要研究方向为全景图像显著性目标检测。E-mail: 1455922810@qq.com