

中图法分类号: TP391 文献标识码: A 文章编号: 1006-8961(XXXX)XX-0001-15

论文引用格式: Zhao Gongbo, Yu Ke, Wang Feng. Local-aware state-space model for remote sensing image super-resolution[J/OL]. Journal of Image and Graphics, XXXX:1-15. DOI: 10.11834/jig.250580. (赵公博, 于珂, 王峰. 面向遥感图像超分辨率的局部感知状态空间模型[J/OL]. 中国图象图形学报, XXXX:1-15. DOI: 10.11834/jig.250580.) [DOI:10.11834/jig.250580]

面向遥感图像超分辨率的局部感知状态空间模型

赵公博, 于珂, 王峰

复旦大学未来信息创新学院, 上海 200433

摘要: 目的 遥感图像超分辨率(remote sensing image super-resolution, RSISR)旨在从低分辨率观测中重建高分辨率影像,实现细节纹理恢复与全局结构保持。然而,现有方法难以兼顾全局依赖建模与局部特征表达。尽管状态空间模型(state space model, SSM)在长程依赖建模方面具有线性复杂度优势,但其一维展开机制易导致局部语义破坏与像素次序扰动,并在性能与计算效率之间陷入权衡。**方法** 为应对上述局限,本文提出基于局部感知的状态空间重建框架 LAMA-SR(local-aware Mamba for super-resolution)。首先,本文对输入遥感影像显式施加二维旋转位置编码(two-dimensional rotary position embedding, RoPE-2D),以建立像素之间的相对空间约束与跨区域关联性,从而缓解一维扫描导致的空间邻接关系丢失问题。随后,编码后的特征通过二维选择性扫描进行全局依赖建模,以在线性计算复杂度下刻画长程上下文。同时,本文引入轻量卷积式局部信息聚合器(local information aggregator, LIA),以增强局部上下文聚合与细粒度纹理恢复,实现对细节区域的补偿性重建。该设计在全局建模能力与局部细节保真度之间形成协同。**结果** 在多个遥感图像基准数据集上,LAMA-SR 相较于主流卷积神经网络(convolutional neural network, CNN)和 Transformer 模型,在峰值信噪比(peak signal-to-noise ratio, PSNR)与结构相似性指数(structural similarity index, SSIM)等指标上均取得显著提升。同时,该模型保持了较低的参数规模与计算开销,展现出优越的效率-效果平衡性。**结论** 实验结果验证了所提出 LAMA-SR 框架在遥感图像超分辨率任务中的有效性与普适性。其在保留 Mamba 长程依赖建模优势的同时,显著改善了局部语义一致性与细节重建能力。

关键词: 深度学习;超分辨率重建;Mamba;位置编码;注意力机制

Local-aware state-space model for remote sensing image super-resolution

Zhao Gongbo, Yu Ke, Wang Feng

School of Information Science and Technology, Fudan University, Shanghai, 200433, China

Abstract: Objective Remote sensing image super-resolution (RSISR) seeks to reconstruct high-resolution (HR) images from low-resolution (LR) observations, thereby enhancing spatial detail while preserving structural integrity across large-scale scenes. Despite remarkable progress achieved by convolutional neural networks (CNNs) and Transformer-based architectures, existing approaches still face an intrinsic trade-off between reconstruction fidelity and computational efficiency. CNNs are adept at capturing fine-grained local textures but suffer from limited receptive fields, whereas Transformers model long-range dependencies effectively but incur quadratic complexity, making them computationally expensive for high-resolution remote-sensing imagery. The recently proposed Selective State Space Model (Mamba) introduces linear-time sequence modeling with strong global dependency representation. However, when extended to vision tasks, Mamba

收稿日期: 2025-11-17; 修回日期: 2026-02-13

* 通信作者: 王峰 fengwang@fudan.edu.cn

基金项目: 上海市协同创新项目(项目编号: XTCX-KJE001-2025-08)

Supported by: the Shanghai Collaborative Innovation Program (XTCX-KJE001-2025-08).

exhibits several limitations, including unidirectional information flow, disrupted pixel ordering due to one-dimensional unfolding, and weak local semantic integrity. These challenges motivate the design of an architecture that can jointly enhance local detail representation and global spatial consistency under constrained computation. **Method** To this end, we propose LAMA-SR (Local-Aware Mamba for Super-Resolution), a novel state-space-based reconstruction framework tailored for RSISR. Instead of embedding positional cues within the scanning module, LAMA-SR first imposes an explicit two-dimensional rotary position embedding (RoPE-2D) on input feature maps to encode relative spatial relationships and inter-pixel geometric dependencies across both axes. This design preserves structural continuity across neighboring regions and mitigates spatial distortion caused by one-dimensional scanning. The position-enriched representation is subsequently processed through a Two-Dimensional Selective Scan (SS2D) mechanism, enabling efficient long-range dependency modeling with linear computational complexity. To enhance the model's capability in capturing fine-grained spatial structures, a Multi-Scale Mixture-of-Experts (MoE) Local Information Aggregator is introduced prior to the Mamba module. The aggregator partitions the input features into small-, medium-, and large-receptive-field branches, enabling adaptive extraction of local patterns across different spatial scales. Unlike conventional MoE architectures that rely on explicit routing, a lightweight channel-attention-based soft routing strategy is employed to dynamically modulate the contribution of each expert according to the input content. This design allows the network to emphasize appropriate local receptive fields—favoring small-scale experts for texture details while leveraging larger-scale experts for structural continuity—thereby enriching the local representation quality and providing stronger inputs for subsequent state-space propagation. Collectively, these components form a cooperative mechanism that unifies global modeling efficiency and local fidelity within a compact Mamba backbone. **Result** Extensive experiments were conducted on four widely used benchmark datasets—UCMerced, RSSCN7, AID, and WHU-RS19—under upscaling factors of $\times 2$, $\times 3$, and $\times 4$. Quantitative evaluations using peak signal-to-noise ratio (PSNR) and structural similarity index (SSIM) demonstrate that LAMA-SR consistently outperforms representative CNN-based and Transformer-based methods, as well as Mamba-based baselines. For instance, on the UCMerced dataset with a $\times 4$ scale, LAMA-SR achieves 29.26 dB PSNR with only 8.10 M parameters and 20.23 G FLOPs, surpassing MambaIR (20.42 M parameters, 79.70 G FLOPs) while maintaining superior image fidelity. Ablation studies further verify that the joint incorporation of RoPE-2D and LIA yields complementary improvements—RoPE-2D enhances spatial coherence and relative position modeling, whereas LIA strengthens texture recovery and channel selectivity. Visual and class-wise analyses confirm that LAMA-SR achieves sharper object boundaries and more natural textures, especially in complex landscapes such as airports, industrial areas, and vegetation regions, validating its robustness and generalization capability across diverse remote-sensing scenarios. **Conclusion** In summary, LAMA-SR provides a state-space-based framework for remote sensing image super-resolution that improves both reconstruction fidelity and computational efficiency. By promoting more consistent spatial information propagation and targeted refinement of local structures, the model better preserves neighborhood continuity and scene layout during reconstruction compared with conventional Mamba-style baselines. Experiments on multiple benchmark datasets show that LAMA-SR yields sharper textures, clearer object boundaries, and a more favorable accuracy-efficiency balance across different upscaling factors. These results indicate that state-space models can be adapted to high-resolution remote sensing imagery while maintaining fine-grained spatial semantics, offering a stable and scalable direction for further research on lightweight super-resolution models.

Key words: deep learning; super-resolution reconstruction; Mamba; position embedding; attention mechanism

0 引言

遥感图像超分辨率重建 (Remote Sensing Image Super-Resolution, RSISR) 旨在从低分辨率遥感影像中恢复出具有更高空间分辨率的图像, 从而提升下游任务如地物识别 (Helber 等, 2020)、目标检测、环

境监测 (Zhao 等, 2020)、灾害管理 (Ghaffarian 等, 2018) 和目标检测 (Dong 等, 2019)。与自然图像不同, 遥感图像往往覆盖大范围场景, 包含尺度差异极大的地物目标 (道路网、建筑群、农田纹理与水系边界等), 并受空气散射、成像平台抖动、光照和传感器物理极限等因素影响而存在细节缺失与噪声污染 (Li 等, 2024)。这使得 RSISR 不仅需要恢复高频纹

© 中国图象图形学报版权所有

理,还必须保持大尺度空间结构和几何一致性,即既不能把边缘“磨平”,也不能让目标结构出现形变或错位(Zhang等,2018)。因此,如何在有限计算预算下同时实现局部细节还原和全局结构保持,成为当前国内外该领域的核心科学问题(Liang等,2021)。

传统方法主要依赖卷积神经网络(convolutional neural network, CNN)(He等,2016)进行单幅图像的超分辨率重建(Dong等,2016;Kim等,2016;Lei等,2022;Li等,2022;Lim等,2017;Wang等,2022;Zhang等,2018)。这些CNN类方法的优势在于:在局部纹理恢复方面具有很强的精细建模能力,尤其适合恢复房屋边缘、道路骨架等高频结构(Lim等,2017)。但它们也有共性局限:卷积的感受野本质上是局部的,即便通过堆叠加深网络或引入空洞卷积,也很难在长距离范围内保持全局一致性,这在遥感场景中表现为“大范围结构容易断裂或不连续”,例如道路网络变形、农田块边界被打散、城市街区轮廓不连贯(Zhang等,2018)。因此,单纯依靠卷积堆叠很难同时建模遥感图像中的长程依赖关系和跨区域一致性(Liang等,2021)。

Transformer结构引入的自注意力机制为长距离依赖建模提供了一种显式的全局建模能力(Vaswani等,2017)。在超分辨率场景中,基于自注意力的模型能够同时观察到广域上下文,避免仅根据局部邻域进行“盲目锐化”,从而在恢复大尺度地物形态和语义连贯性方面优于纯CNN方法(Cai等,2023;Lei等,2022;Liang等,2021;Lu等,2022;Xiao等,2024)。尤其是窗口化或分层式视觉Transformer结构通过分区注意力和跨窗口交流,实现了对大范围纹理模式的一致建模,同时一定程度上控制了注意力计算的规模。不过,自注意力的计算与显存开销依然随空间分辨率近似平方增长。在遥感场景中,原图往往是超大幅面切块推理,高分辨率补偿本身就非常昂贵,这使得Transformer类RSISR方法在边缘设备、机载/星载实时判读、地面快速应急监测等场景下难以直接部署(Wu等,2025)。

为同时获得长程建模能力与较低复杂度,近年来兴起的状态空间模型(state space model, SSM),尤其是Mamba结构,被引入到遥感图像超分辨率中并迅速成为新的研究焦点(Gu等,2021)。Mamba通过选择性状态传递实现跨长距离的依赖建模,其推理复杂度在序列长度上接近线性规模,相比自注意力

的全局两两交互更加高效(Gu等,2023)。这为RSISR带来了一个新的可能路径:不再依赖纯卷积的局部堆叠,也不完全依赖全注意力的二次复杂度,而是用线性时间的全局建模去维持大尺度地物结构的连贯性。但由于Mamba最初是从序列建模场景发展而来,它天然依赖对输入进行“扫描式”的状态更新,这在视觉任务中通常意味着把二维影像按行、列或块展平成一维序列,再按序推进状态。这种展平会破坏原始空间邻接关系,导致局部几何细节(例如屋顶边缘的转角、道路的交叉口)在重建时出现模糊、粘连或错位(Jiang等,2025)。

针对Mamba在遥感超分中的局限,不同研究工作从结构融合、注意力增强到频域建模等方向展开了改进探索,形成了一系列具有代表性的方法。MambaFormerSR提出将卷积分支、Transformer注意力分支与Mamba分支并行融合,分别负责局部纹理恢复、显著区域加权与全局一致性保持(Zhi等,2024)。该结构能较好兼顾细节与语义,但多分支耦合导致模型推理路径冗长,计算与内存开销增大,不利于轻量化部署。Efficient Mamba-Attention Network更注重效率,通过显著性引导与多尺度轻量特征提取模块,将计算聚焦于关键地物区域,并在Mamba状态传播中嵌入通道-空间注意力以强化长程依赖(Wu等,2025)。其优势在于在保持近线性复杂度的同时提升判读精度,但依赖注意力补偿使模型仍存在人工权衡;当场景纹理规则或重复时,显著性机制易压制真实边界,造成局部平滑。MFEM(multiscale frequency-enhanced mamba)在Mamba主干内部引入多尺度卷积与频域增强模块,直接提升模型的局部与高频感知能力(Chen等,2025)。该方法能有效缓解“全局一致而局部模糊”

的问题,使Mamba既具长程建模又能保持纹理锐利;但多尺度卷积与频域分支增加了结构复杂性与推理负担。Rep-Mamba则从空间结构保持性出发,提出跨尺度状态传播与重参数化卷积机制,在推理阶段折叠多分支卷积以降低运行成本(Jiang等,2025)。该方法强化了全局-局部协同,对遥感图像中复杂几何边界表现优异;但其补偿仍发生于Mamba外部,未从根本上解决二维影像在一维扫描中空间关系受损的问题。

综上所述,现有基于Mamba的遥感超分辨率方法虽然在不同层面提升了模型的适应性,但普遍依

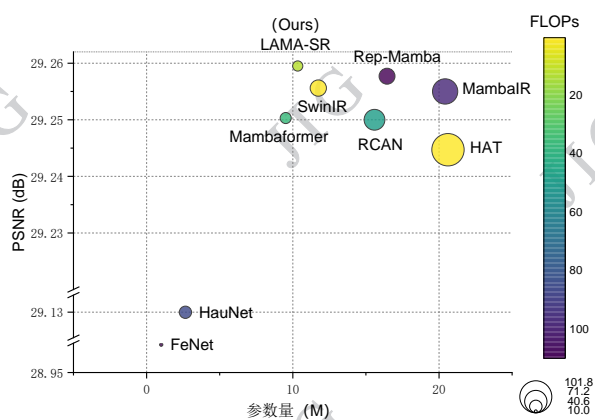


图1 模型复杂度与重建性能对比图

Fig. 1 Comparison of model complexity and reconstruction performance

赖外部卷积、注意力或频域分支进行修正,结构复杂、推理开销大,且缺乏一种能够在 Mamba 内部同时实现二维空间建模与局部几何保持的统一机制。为此,本文提出一种基于二维旋转位置编码(RoPE-2D)(Su等,2024)与轻量局部信息聚合器(Local Information Aggregator, LIA)的 Mamba 遥感超分辨率模型。该方法在状态传播过程中显式注入空间相对位置信息,使 Mamba 具备局部几何感知与边缘连续性;同时, LIA 模块在进入 Mamba 前以混合专家(Shazeer等,2017)(Mixture of Experts, MoE)结构自适应聚合不同尺度的局部邻域信息。具体而言,输入特征按通道划分为多组并分别经小、

中、大感受野的专家卷积分支提取特征,再由通道注意力机制(Hu等,2018)实现轻量“软路由”加权,从而在低计算开销下自适应选择最合适的局部感受野。该机制显式建模多尺度局部依赖与边缘连续性,为后续的 Mamba 状态传播提供结构保持性更强的输入表征。整体上, RoPE-2D 提升了 Mamba 的二维空间建模能力,而 LIA 则弥补其局部几何感知的不足,使模型能够在保持线性复杂度的同时兼顾细节纹理与全局结构一致性。如图 1 所示,本文的方法在 UC Merced $\times 4$ 数据集上对多个典型超分辨率模型进行了 PSNR、参数量与 FLOPs 对比。实验结果表明,本文提出的方法在重建质量和计算效率之间实现了良好平衡:在显著降低参数数量与计算开销的同时,取得了当前最优的重建性能,充分展示了其在性能与效率上的综合优势。

本文的主要贡献可概括如下:

(1)提出了一种融合 RoPE-2D 与状态空间建模的遥感图像超分方法,从结构层面提升 Mamba 的二维空间感知能力;

(2)基于混合专家模型设计轻量局部聚合模块(LIA),在保持线性复杂度的同时,增强了模型对地物边缘与纹理的恢复效果;

(3)设计了完整的实验方案,在多个遥感数据集上验证了所提模型在性能-效率平衡方面的优越性,为轻量化高保真 RSISR 提供了一种可推广的解决方案。

1 本文方法

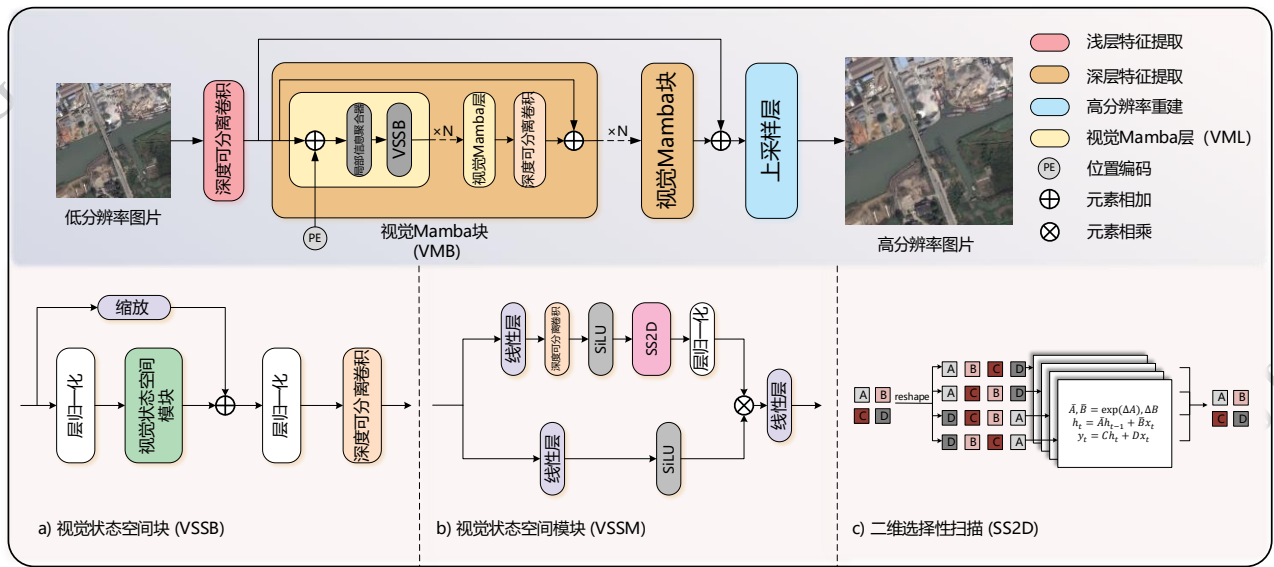
本节将介绍本文提出的超分辨率(SR)框架,其目标是从低分辨率输入中重建高分辨率图像。该网络主要由三个模块组成:浅层特征提取器、基于视觉 Mamba 模块(vision mamba blocks, VMBs)构建的深层特征提取主干网络,以及图像重建模块。为了进一步增强特征表示能力,框架中集成了三个关键组件:局部信息聚合器(LIA)、旋转位置编码(RoPE)以及视觉状态空间模块(visual state space block, VSSB)。各组件的具体细节将在下文中逐一介绍。

1.1 整体框架

如图 2 所示,本文提出的网络专为遥感图像超分辨率任务设计,主要由三个阶段组成:浅层特征提取、深层特征提取和高分辨率图像重建。网络的输入为一张低分辨率图像 $I_{LR} \in \mathbb{R}^{H \times W \times 3}$,其中 H 和 W 分别表示图像的高度和宽度,3 表示 RGB 通道数。首先,该图像经过一个深度可分离卷积(depthwise separable convolution, DWConv)层处理,用于提取低级或浅层特征,记为 $F_s = DWConv(I_{LR})$,输出表示为 $F_s \in \mathbb{R}^{H \times W \times C}$,以及 C 表示特征通道的数量。

这些浅层特征随后被送入一系列 VMB(vision mamba blocks)中进行深层特征提取。每个 VMB 包含多个 VML(Vision Mamba Layers),而每个 VML

由三个关键组件组成:RoPE、LIA 和 VSSB。首先,将位置编码添加到输入特征中,以向特征注入空间位置信息;然后,特征图通过 LIA 以捕捉短程依赖关系,从而增强局部上下文信息;接着, VSSB 模块对整张图像建模长程空间依赖,实现高效的全局特征表示。第 i 个 VMB 中第 l 个 VML 的完整计算过程



(a) VSSB; (b) VSSM; (c) SS2D

图2 LAMA-SR整体框架

Fig. 2 Framework of the LAMA-SR

可表示为:

$$\mathbf{F}^{(i,l)} = \text{VSSB}(\text{LIA}(\mathbf{F}^{(i,l-1)} + \text{PE})) \quad (1)$$

在第 i 个 VMB 中经过所有 VML 处理后,会在 VMB 末尾应用一个 DWConv 层,以进一步优化局部特征信息。最终,为了保持训练的稳定性,将输入与 VMB 的输出通过残差连接进行相加:

$$\mathbf{F}^{(i)} = \text{DWConv}(\mathbf{F}^{(i,L)}) + \mathbf{F}^{(i-1)} \quad (2)$$

其中 L 表示每个 VMB 中 VML 的数量。上述过程重复执行 N 次,得到最终的深层特征图 \mathbf{F}_D 。整个深层特征提取阶段可总结为:

$$\mathbf{F}_D = \text{VMB}_N(\dots(\text{VMB}_1(\mathbf{F}_S))) \quad (3)$$

然后,将浅层特征与深层特征通过逐元素相加进行融合:

$$\mathbf{F}_{\text{fused}} = \mathbf{F}_D + \mathbf{F}_S \quad (4)$$

最后,通过上采样层对融合后的特征图进行重建,生成高分辨率图像:

$$\mathbf{R}_{\theta,m} = \begin{pmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \cdots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos m\theta_{d/2} & -\sin m\theta_{d/2} \\ 0 & 0 & 0 & 0 & \cdots & \sin m\theta_{d/2} & \cos m\theta_{d/2} \end{pmatrix} \quad (8)$$

该变换在嵌入空间的二维子空间中施加一个逆时针旋转,旋转角度为 θ_i ,定义如下:

$$\theta_i = 10000^{-2(i-1)/d_{\text{model}}}, i \in [1, \dots, d_{\text{model}}/2] \quad (9)$$

其中 i 表示嵌入维度的索引, d 是嵌入的总维度。该

公式确保不同的嵌入维度以不同的速率进行旋转：较低索引的维度编码更细粒度的位置差异，而较高索引的维度则用于捕捉更宽泛、粗粒度的位置关系。

旋转矩阵 $R_{\theta,m}$ 在复向量空间中起作用，其中嵌入的每个二维子空间都会根据其位置索引执行一个受控旋转。该变换可重写为复指数形式如下：

$$R_{\theta,m} = e^{im\theta} \quad (10)$$

该形式将位置变换表示为复平面中的纯相位偏移，这意味着在改变角度关系的同时保持向量模不变。因此，依赖于嵌入向量之间内积的自注意力机制，能通过这种旋转相位偏移自然地捕捉相对位置信息。

一种处理图像中位置编码的自然方法是将二维

$$R_{\theta(p',p'')} = \begin{pmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos n\theta_1 & -\sin n\theta_1 & \cdots & 0 & 0 \\ 0 & 0 & \sin n\theta_1 & \cos n\theta_1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos n\theta_{d/4} & -\sin n\theta_{d/4} \\ 0 & 0 & 0 & 0 & \cdots & \sin n\theta_{d/4} & \cos n\theta_{d/4} \end{pmatrix} \quad (12)$$

由于 RoPE 由一维序列扩展到二维坐标后，单轴位置索引的尺度通常从长度 L 变为约 \sqrt{L} 量级，为保持与原始一维 RoPE 相近的相位变化范围，需要对式 (9) 中的常数系数将取平方根，从而得到更新后的定义：

$$\theta_i = 100^{-4(i-1)/d_{\text{model}}}, i \in [1, \dots, d_{\text{model}}/4] \quad (13)$$

二维 RoPE 保证了不同嵌入维度在适当比例下进行旋转，使得模型能够同时捕捉局部与全局的空间关系。由于旋转变换分别作用于 x 和 y 坐标，模型在两个维度上都能保持等距的空间关系，避免了展平成一维序列所引入的空间扭曲 (Heo 等, 2024)。

1.3 局部信息聚合器

如图 3 所示，本文提出了一种混合专家 (Mixture of Experts, MOE) 驱动局部信息聚合器，用于在特征输入 Mamba 模块之前充分整合其邻域像素信息，从而弥补状态空间模型在局部几何结构建模方面的不足。给定输入特征 $X \in R^{C \times H \times W}$ ，首先按照通道维度将其划分为三组 $\{X^S, X^M, X^L\}$ ，分别对应小、中、大三类局部感受野。随后，这三组特征分别送入卷积核大小为 3、5、7 的多尺度专家，以捕获不同空间尺度下的局部结构与纹理细节：

$$Z_3 = f_{3 \times 3}(X^S), Z_5 = f_{5 \times 5}(X^M), Z_7 = f_{7 \times 7}(X^L) \quad (14)$$

特征图展平为一维序列，并应用标准的一维位置编码。然而，这种方法会引入相对空间关系的不一致性。为了解决这一问题，RoPE 被扩展到了二维空间 (Heo 等, 2024)。在一维 RoPE 中，位置信息通过乘法旋转变换进行编码，从而使相对位置信息自然融入嵌入向量中。将这一概念扩展到二维时，需要在两个空间维度上分别施加独立的旋转变换，同时保留相对位置结构。

对于具有网格表示的图像，特征图中每个嵌入向量 $x_{(m,n)}$ 对应于位置 (m,n) 。由此，二维 RoPE 的变换形式为：

$$f(x_{(m,n)}, m, n) = R_{\theta(m,n)} x_{(m,n)} \quad (11)$$

式中， $R_{\theta(m,n)}$ 为二维旋转矩阵，其定义为：

与传统 MoE 结构依赖显式路由 (hard routing) 不同，传统方法通常通过 Top-k 等离散门控对每个位置、特征进行“硬选择”，仅将输入分发给少数专家参与

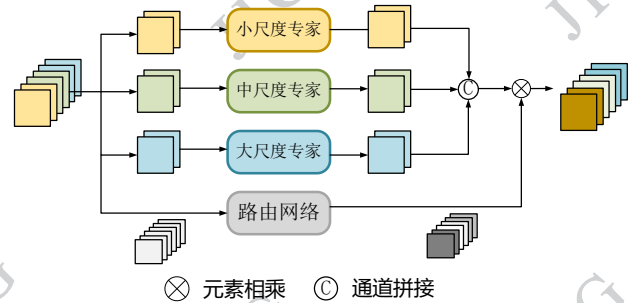


图 3 通道注意力驱动的多尺度局部信息聚合器

Fig. 3 Multi-scale local information aggregator

计算，并需要额外的负载均衡约束来缓解专家塌缩等训练不稳定问题；并且，在超分重建中，同一空间位置往往需要同时借助小感受野分支恢复细纹理、并借助大感受野分支维持结构轮廓，而离散硬路由的“单一尺度选择”容易丢失多尺度互补信息。相比之下，本文在 LIA 中以通道注意力作为门控，生成对各尺度 (3/5/7 卷积) 分支的连续权重并进行加权融合，形成无需离散分发、梯度可同时回传至所有分

支的“软路由”机制,使多尺度信息能够按内容自适应地平滑混合,等效于动态选择有效感受野并保留互补细节;从而在保持轻量化的同时获得更平滑的专家选择。具体而言,首先对输入特征进行全局通道汇聚,得到压缩后的通道描述子:

$$r = \text{GAP}(X) \in R^c \quad (15)$$

然后通过两层瓶颈式线性变换生成对应于三个尺度专家的通道级注意力权重:

$$[a^s, a^m, a^l] = \text{Sigmoid}(W_2 \varphi(W_1 r)) \quad (16)$$

其中 GAP 表示全局池化操作, $a^{(c)} \in R^{c^{(c)}}$ 分别对应小、中、大尺度专家的通道权重, $\varphi(\cdot)$ 为 relu 激活函数。得到的小、中、大尺度专家权重在空间维度上进行广播,并分别与各尺度专家的输出进行逐通道加权,最后沿通道维拼接形成最终的聚合特征:

$$Y = \text{Concat}(a^{(s)} \odot Z_3, a^{(m)} \odot Z_5, a^{(l)} \odot Z_7) \quad (17)$$

通过这种通道注意力驱动的多尺度特征汇聚机制,局部信息聚合器具备以下显著优势:

1) 自适应选择局部感受野:在几乎不增加计算开销的情况下,模型能够根据输入内容自动选择合适的尺度分支。对于纹理细密区域,更多权重将分配给小尺度专家;在结构跨度较大的区域,中、大尺度专家将获得更高权重,这使模型能够更灵活地处理不同类型的局部结构。

2) 增强多尺度局部依赖建模能力:通过显式融合多尺度感受野,模型能够更好地保持边缘连续性和局部几何一致性,提升对复杂空间的表达能力。

3) 为状态空间传播提供更优特征表示:经过多尺度聚合后得到的特征在局部纹理、边缘信息及空间结构上更为丰富与稳定,为后续 Mamba 模块的长距离状态传播提供了高质量的输入。

2 实验与分析

2.1 数据集与实验设置

2.1.1 数据集

为评估所提出的超分辨率方法的性能,本文选取了四个广泛使用的遥感图像数据集。这些数据集包含多种地物类型和场景变化的高分辨率航拍图像,为遥感图像超分辨率模型提供了全面的测试平台:

1) UCMerced(Yang等,2010):该数据集图像来

源于美国地质调查局的国家地图。该数据集涵盖 21 类地物使用类型,例如农业用地、飞机、海滩、建筑、森林、高速公路和停车场等,每类包含 100 张图像。数据集被分为训练集和测试集,各包含 1050 张图像,其中训练集的 20% 被进一步用作验证集。

表 1 基于 UCMerced×4 的不同方法的参数, FLOPs, 内存占用以及 PSNR 对比

Table 1 Comparison of parameters, FLOPs, memory usage, and PSNR on UCMerced ×4

方法	参数量	FLOPs	PSNR	SSIM
FeNet	0.37 M	1.44 G	28.9975	0.7832
HAT	20.62 M	101.72 G	29.2587	0.7957
HauNet	2.66 M	38.95 G	29.1298	0.7895
MambaIR	20.42 M	79.70 G	29.2550	0.7956
RCAN	15.59 M	65.25 G	29.2500	0.7942
SwinIR	11.75 M	51.33 G	29.2556	0.7949
SR-Mambaformer	9.52 M	34.82 G	29.2542	0.7933
Rep-Mamba	16.45 M	50.04G	29.2617	0.7964
Ours	8.10 M	20.23 G	29.2695	0.7980

2) RSSCN7(Zou等,2015):包含 2800 张遥感图像,涵盖 7 类场景类型,分别为:草地、森林、农田、工业区、居民区、河流、湖泊和停车场。每类场景包含 400 张图像,图像采集于不同季节和大气条件下,以增强数据的多样性。所有图像的空间分辨率为 400×400 像素。该数据集被平均划分为训练集和测试集,各包含 1400 张图像。

3) AID(Xia等,2017):包含 10000 张图像,图像采集自 Google Earth,空间分辨率为 600×600 像素。该数据集涵盖 30 类地物类型,包括机场、港口、公园、教堂、高架桥和广场等。其中 20% 的图像被用作测试集,剩余 80% 用于训练,并从每个类别的训练集中随机选取 5 张图像作为验证集。

4) WHU-RS19(Dai等,2011):由 1005 张分辨率为 600×600 像素的航拍图像组成,涵盖 19 类地物类型,包括机场、港口、沙漠、草地、河流和农田等。该数据集具有高分辨率和多样化的地貌场景,为评估超分辨率模型在遥感图像中重建细节能力提供了良好的测试基础。该数据集用作跨数据集评估。

2.1.2 实验设置

在上述数据集中,超分辨率(SR)任务采用三种

表2 不同方法在 UCMerced 和 RSSCN7 数据集上的性能比较

Table 2 PSNR and SSIM comparisons of different methods on the UCMerced and RSSCN7 datasets

Methods	Metric	UCMerced ×2	UCMerced ×3	UCMerced ×4	RSSCN7 ×2	RSSCN7 ×3	RSSCN7 ×4
Bicubic	PSNR	30.6585	27.3787	25.5859	29.3264	27.4110	26.2983
	SSIM	0.8893	0.7811	0.6892	0.7975	0.6952	0.6193
FeNet	PSNR	35.8796	31.3746	28.9975	31.6946	29.5600	28.3596
	SSIM	0.9419	0.8592	0.7832	0.8282	0.7362	0.6741
RCAN	PSNR	36.1224	31.7282	29.2500	31.8503	29.6913	28.2911
	SSIM	0.9440	0.8656	0.7942	0.8339	0.7419	0.6711
HauNet	PSNR	36.0486	31.5285	29.1298	31.7428	29.6283	28.4546
	SSIM	0.9431	0.8649	0.7895	0.8305	0.7401	0.6808
SwinIR	PSNR	36.1641	31.7632	29.2556	31.8216	29.6724	28.4485
	SSIM	0.9449	0.8671	0.7949	0.8330	0.7422	0.6803
HAT	PSNR	36.1721	31.7704	29.2587	31.8627	29.7008	28.4886
	SSIM	0.9451	0.8671	0.7957	0.8340	0.7416	0.6825
MambaIR	PSNR	36.1376	31.6984	29.2550	31.8545	29.6946	28.4930
	SSIM	0.9437	0.8670	0.7956	0.8336	0.7431	0.6820
SR-Mambaformer	PSNR	36.1432	31.7321	29.2542	31.8487	29.6941	28.4832
	SSIM	0.9443	0.8660	0.7933	0.8325	0.7419	0.6811
Rep-Mamba	PSNR	36.1659	31.7787	29.2617	31.8743	29.7070	28.4993
	SSIM	0.9438	0.8669	0.7964	0.8333	0.7420	0.6813
LAMA-SR(本文)	PSNR	36.1750	31.7916	29.2695	31.9076	29.7068	28.5095
	SSIM	0.9450	0.8684	0.7980	0.8351	0.7436	0.6816

不同的放大倍率:×2、×3和×4。低分辨率(LR)图像通过对原始高分辨率(HR)图像进行双三次插值降采样生成。训练后的模型将这些低分辨率图像重建为高分辨率图像,其性能通过(peak signal-to-noise ratio, PSNR) (Gao 等, 2009)与结构相似性指数(structural similarity index, SSIM) (Wang 等, 2004), 进行评估,从而能够在多个数据集和放大倍率下对不同的超分辨率方法进行全面比较。模型采

用 Adam 优化器(Kingma 等, 2014)并使用 L1 损失函数进行优化,超参数设置为 $\beta_1 = 0.9$ 以及 $\beta_2 = 0.99$ 。在训练过程中,应用了水平翻转和随机旋转等数据增强技术,以提升模型的泛化能力。所有实验基于 PyTorch 框架实现,并在四张 NVIDIA

RTX 4090 GPU 上运行。展示的实验结果均为更换随机种子训练三次后取平均值所得。

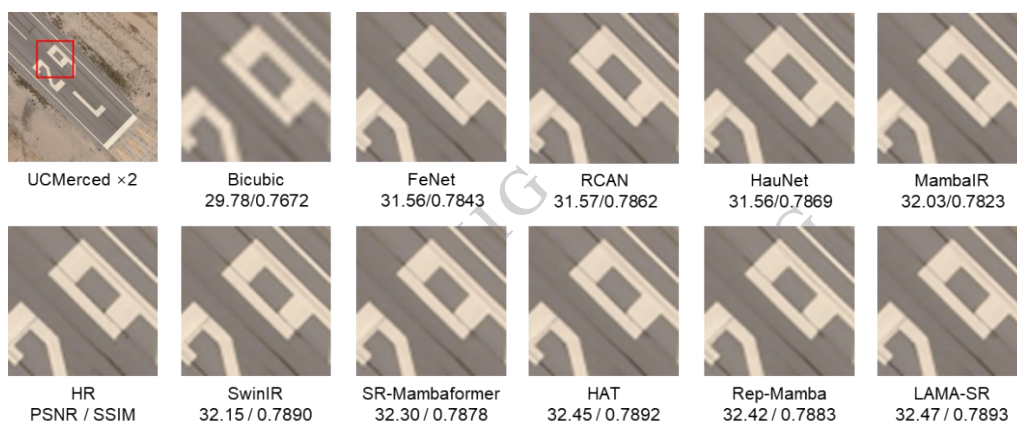
2.2 对比试验

为评估所提出方法在性能与效率方面相较于主流超分辨率模型的表现,实验选取了多种具有代表性的基线模型进行对比,包括 FeNet (feature enhancement network) (Wang 等, 2022), HauNet (hybrid attention based U-shaped network) (Wang 等, 2023), MambaIR (Mamba based image restoration) (Guo 等, 2024), RCAN (Zhang 等, 2018), HAT (hybrid attention Transformer) (Chen 等, 2023), SwinIR (image restoration using Swin Transformer)

表3 不同方法在AID和WHU-RS19数据集上的性能比较

Table 3 Quantitative results of different models on AID and WHU-RS19 datasets

Methods	Metric	AID ×2	AID ×3	AID ×4	WHU-RS19 ×2	WHU-RS19 ×3	WHU-RS19 ×4
Bicubic	PSNR	32.3762	29.0681	27.2769	30.0867	27.5652	26.0896
	SSIM	0.9009	0.8033	0.7183	0.8401	0.7364	0.6543
FeNet	PSNR	36.8855	32.7916	30.5350	31.9460	29.5212	28.0519
	SSIM	0.9440	0.8711	0.8033	0.8463	0.7611	0.6968
RCAN	PSNR	37.0335	33.0194	30.5921	31.9359	29.5107	28.0713
	SSIM	0.9451	0.8761	0.8048	0.8470	0.7612	0.6978
HauNet	PSNR	37.0055	32.9760	30.7812	31.9414	29.575	28.0749
	SSIM	0.9452	0.8751	0.8110	0.8475	0.7617	0.6979
SwinIR	PSNR	37.0807	33.0515	30.8544	31.9512	29.5305	28.0577
	SSIM	0.9463	0.8774	0.8128	0.8480	0.7610	0.6972
HAT	PSNR	37.0904	33.0721	30.8229	31.9534	29.5437	28.0563
	SSIM	0.9464	0.8780	0.8130	0.8481	0.7617	0.6977
MambaIR	PSNR	37.0605	33.0397	30.8059	31.9530	29.5310	28.0752
	SSIM	0.9456	0.8760	0.8113	0.8473	0.7615	0.6980
SR-Mambaformer	PSNR	37.0736	33.0587	30.8148	31.9502	29.5364	28.0693
	SSIM	0.9453	0.8755	0.8107	0.8472	0.7613	0.6975
Rep-Mamba	PSNR	37.0910	33.0679	30.8616	31.9531	29.5432	28.0778
	SSIM	0.9459	0.8765	0.8096	0.8478	0.7615	0.6977
LAMA-SR(本文)	PSNR	37.1002	33.0711	30.8652	31.9610	29.5583	28.0869
	SSIM	0.9466	0.8772	0.8143	0.8481	0.7621	0.6983



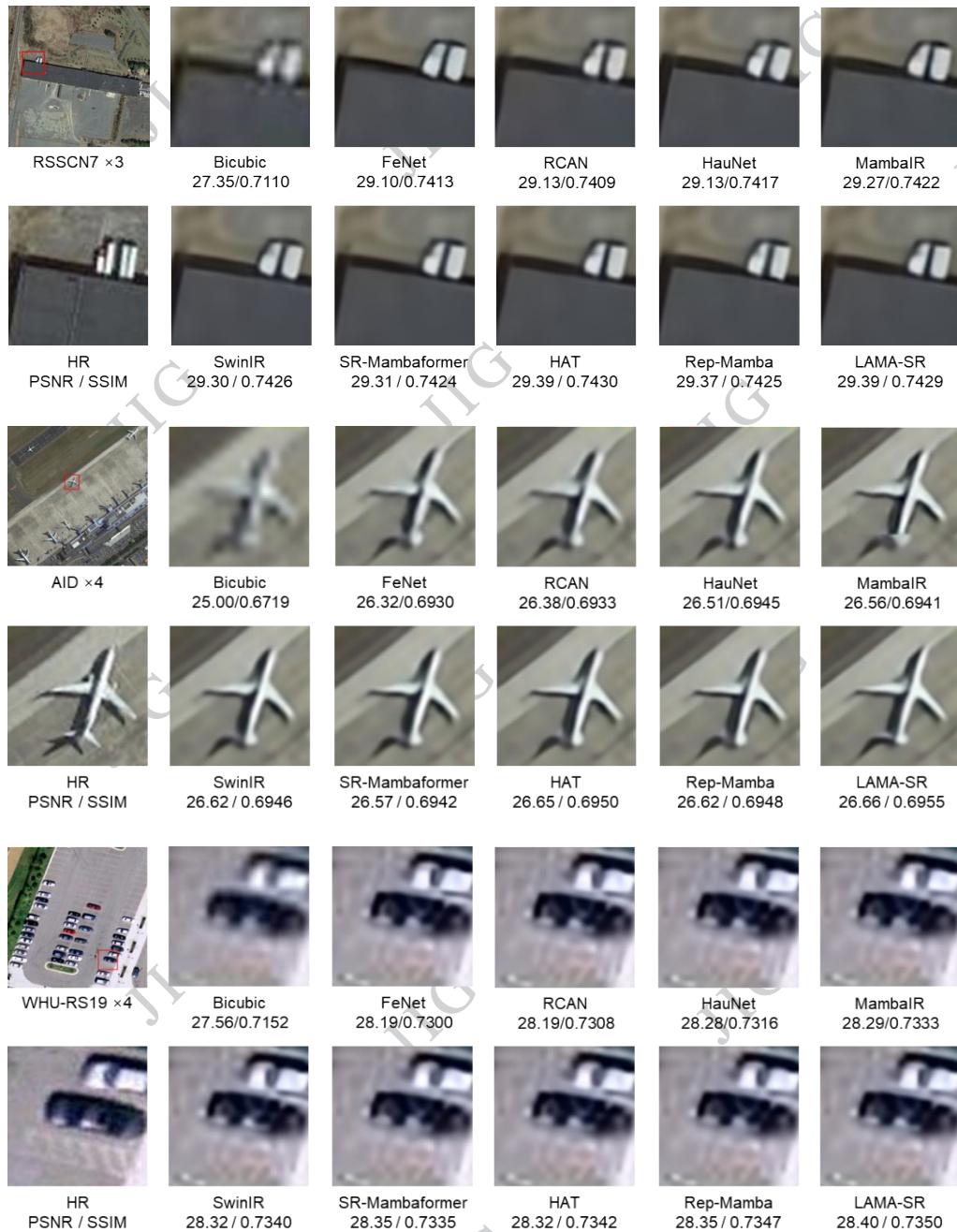


图4 不同方法在四个数据集上的视觉对比

Fig. 4 Visual comparisons of different methods on four datasets

(Liang 等, 2021), MambaformerSR (A Lightweight Model for Remote-Sensing Image Super-Resolution) (Zhi 等, 2024), Rep-Mamba (reparameterization in vision mamba for lightweight remote sensing image super-resolution) (Jiang 等, 2025)。

表 1 给出了 UC Merced 数据集 ×4 放大下各方法的参数量与 FLOPs 以及对应的 PSNR/SSIM。本文的方法在 PSNR 上取得最高值 29.2695 dB, 同时

FLOPs 仅为 20.23G、参数量 8.10M, 在达到 SOTA 级别精度的同时显著降低计算开销。对比来看, SwinIR (51.33G) 与 MambaIR (79.70G) 虽在 PSNR 上接近 (分别为 29.2556 与 29.2550 dB), 但计算量显著更高; 而 FeNet 与 HauNet 等轻量模型的 PSNR 分别为 28.9975 与 29.1298 dB, 均未达到 29.15 dB。在 SSIM 指标上, 本文的方法取得 0.7980 的最高值。总体而言, 本文的方法在保持低计算复杂度的前提下, 展示了更优的精度-效率权衡。

此外,本文在四个遥感数据集—UCMerced、RSSCN7、AID 和 WHU-RS19 上进行了全面评估,涵盖多个放大倍率($\times 2$ 、 $\times 3$ 和 $\times 4$),以测试所提出方法的性能与泛化能力。与体量较小的 UCMerced 和 RSSCN7 相比,AID 和 WHU-RS19 包含更大规模的场景与更多类别的多样性。特别地,WHU-RS19 用于跨数据集评估,即所有模型在 AID 数据集上训练后,直接在 WHU-RS19 上进行测试。相关结果展示于表 II 和表 III。图 4 提供了在四个数据集上进行不同网络定性对比的结果,以补充定量实验结果。在这些基准测试中,本文的方法展现出稳健且一致的性能表现。在 UCMerced 和 RSSCN7 数据集上,模型在所有放大倍率下始终获得最高或接近最高的 PSNR 和 SSIM 分数。在更大规模的 AID 数据集上,尽管本文的模型相比于 HAT 和 SwinIR 等基于 Transformer 的方法具有明显更少的参数量,但仍取得了具有竞争力的性能。特别地,与参数量远超本文模型的 HAT 模型相比,本文的方法在 PSNR 上仍能取得优势。值得注意的是,尽管 LAMA-SR 与 MambaSR、Rep-Mamba、SR-Mambaformer 都基于 Mamba 架构,但在 PSNR 指标上,LAMA-SR 始终优于其他基于 Mamba 架构的模型,进一步验证了所提出设计的有效性。此外,在用于跨数据集评估的 WHU-RS19 数据集上,LAMA-SR 在 $\times 3$ 和 $\times 4$ 放大倍率下取得了最高的 PSNR 和 SSIM,表明其在训练数据分布之外的出色泛化能力。

2.3 消融试验

2.3.1 Mamba block 数量消融

在网络深度(Mamba Block 数量)的设定上,本文遵循“任务表征能力—计算复杂度—跨场景泛化”的三方权衡。首先,遥感超分常面对大幅面高分辨率影像,训练与推理通常采用滑窗或切块方式,网络深度增加会带来近线性的参数与 FLOPs 增长,并显著提升显存占用与推理时延;因此本文对主干深度设置上限,以满足实际应用中对效率与可扩展性的约束。其次,从任务适配性出发,RSISR 不仅需要长程依赖以保证道路、建筑群等大尺度结构的连续性,也需要多尺度局部建模以恢复纹理与边缘细节。本文在单个基本单元中通过 RoPE-2D 显式注入二维相对位置约束,通过 LIA 完成多尺度局部细节聚合,并结合 Mamba 进行全局依赖建模,使得每个 Block 本身已具备“局部增强+全局一致”的协同表达

能力,从而降低了单纯依赖堆叠更深网络来扩展感受野的必要性。基于上述原则,本文在合理复杂度范围内对 进行验证:如图 5 所示,当 从较小值增加至 6 时,模型表征能力增强、重建质量显著提升;但继续加深会带来收益递减并增加过拟合风险——在规模较小、场景相对集中的 UCMerced 上性能明显下降,而在规模更大、场景更复杂的 AID 上增益趋于饱和并出现轻微回落。综合精度提升、复杂度约束与跨数据集泛化表现,本文最终采用 $=6$ 作为默认配置,以在重建质量与计算开销之间取得更优折中。

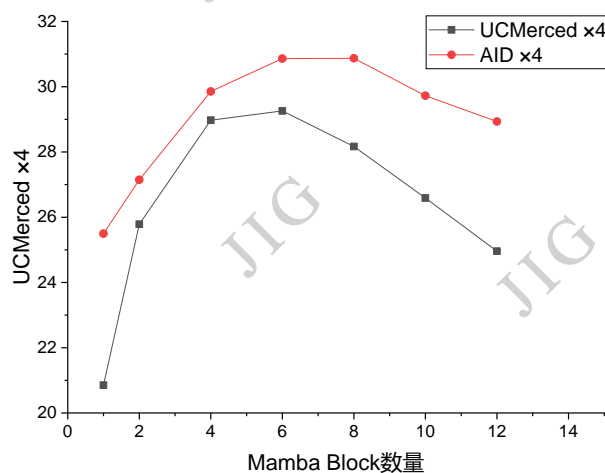


图 5 Mamba block 对性能的影响

Fig. 5 Effect of Mamba Block Number on Model Performance

2.3.2 不同位置编码方法对比

为验证不同位置编码策略对模型性能的影响,本文在相同网络结构与训练配置下在 UCMerced 数据集上进行 $\times 2$ 实验以对多种位置编码方式进行系统对比实验。模型中的位置编码分别替换为 sinusoidal-1D、sinusoidal-2D 和 RoPE-2D,实验结

果如表 4 所示。结果显示,引入 sinusoidal-1D 编码后,相较于未使用位置编码的模型,性能得到一定提升;当采用 sinusoidal-2D 编码时,模型的 PSNR

表 4 不同位置编码方法比较

Table 4 Comparison of different positional encoding methods

Network	PSNR
Baseline with sinusoidal-1D	36.0746
Baseline with sinusoidal-2D	36.1173
Baseline with RoPE-2D	36.1233

进一步提高,说明在需要空间一致性的任务中,相较一维编码,捕捉二维空间位置信息(横向与纵向)更加有效。值得注意的是, RoPE-2D 编码在所有配置中表现最为优异,取得了最高的 PSNR。该性能提升归因于 RoPE-2D 能够建模更复杂的空间关系与旋转不变特征,相较于传统正弦位置编码,在保持空间一致性和捕捉高阶位置依赖方面具有更显著优势。

2.3.3 RoPE 和 LIA 的有效性分析

为评估 RoPE-2D 与所提出的 LIA 模块的有效性,在基线模型的基础上分别集成 RoPE-2D 和 LIA,进行了系列消融实验。

在 UCMerced 数据集 $\times 2$ 超分辨率任务的消融结果如表 5 所示。相较于基线模型(36.0539 dB),

仅引入 RoPE-2D 后 PSNR 提升至 36.1233 dB,且参数量保持不变(7.85M),FLOPs 仅增加约 0.01G,表明二维相对位置建模能够以极低开销增强空间结构一致性与几何约束。仅引入 LIA 时 PSNR 提升至 36.1674 dB,同时参数量增加约 0.25M、FLOPs 增加约 2.83G,说明多尺度局部信息聚合对纹理与边缘细节恢复贡献更显著,但计算开销主要由该模块带来。进一步地, RoPE-2D 与 LIA 联合使用时取得最高 PSNR (36.1750 dB),而额外 FLOPs 仅增加 0.01G,显示 RoPE-2D 可作为几乎无参数的轻量增强,为 LIA 的局部聚合提供更稳定的二维位置参照,从而在不显著增加复杂度的前提下进一步提升重建质量与空间一致性。

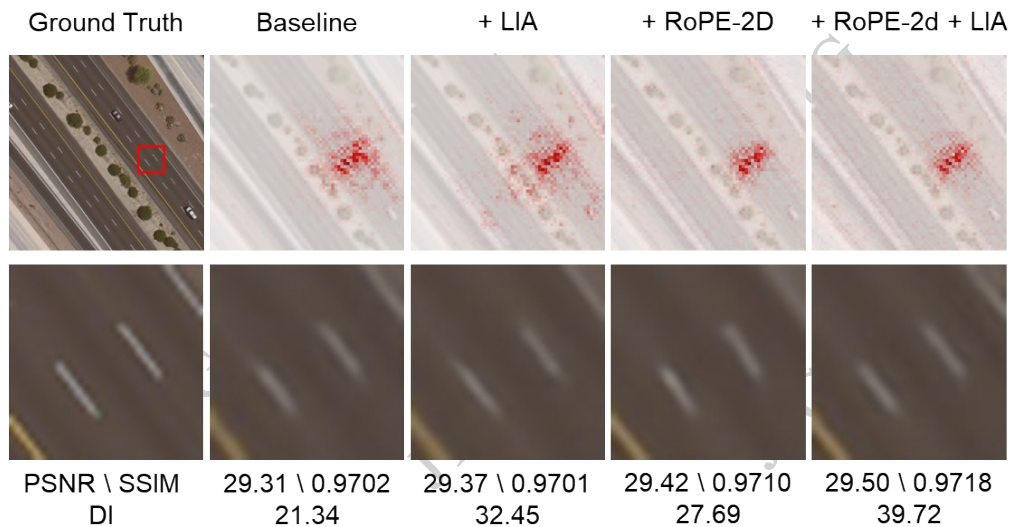


图6 局部归因图可视化

Fig. 6 Visualization of LAMs

总体而言,实验结果强调了在局部聚合机制中引入位置编码的重要性。尽管 RoPE-2D 本身能够增强空间感知能力,但与 LIA 的结合进一步提升了模型在高质量图像重建中的表现,这一点由显著提高的 PSNR 结果所验证。

进一步地,为了从可解释性角度验证 RoPE-2D 与 LIA 模块在特征表达与空间建模方面的作用,本文在 UCMerced $\times 4$ 数据集上进行了基于局部归因图(local attribution map, LAM)(Gu 等, 2021)的可视化分析,如图 6 所示。

LAM 结果直观地揭示了模型在重建特定区域时所关注的输入像素分布范围及强度。可以观察

到,基线模型的归因响应范围较窄,主要集中在目标边缘附近,表明其对上下文信息的利用有限;仅引入 LIA 模块时,归因响应明显扩散,模型能够感知到更宽的局部邻域,但由于缺乏明确的位置关系约束,关注区域呈现一定的漂移;仅引入 RoPE-2D 时,模型在空间结构上表现出更稳定的响应分布,能够保

持合理的边界定位;而在 RoPE-2D 与 LIA 联合使用的情况下,显著图既呈现更广的扩散范围,又能保持精确的空间聚焦,说明两者在位置建模与局部信息聚合方面形成了互补。

此外,基于 LAM 中提出的扩散指数(Diffusion Index, DI)对不同模型的特征聚合范围进行了量化。

表5 不同 RoPE-2D 与 LIA 组合的性能比较

Table 5 Performance comparison of different combinations of RoPE-2d and LIA

	Baseline			
RoPE-2D	X	√	X	√
LIA	X	X	√	√
参数量	7.85 M	7.85 M	8.10 M	8.10 M
FLOPs	17.39 G	17.40 G	20.22 G	20.23 G
PSNR	36.0539	36.1233	36.1674	36.1750

结果显示,如图中所示,联合模型的 DI 值达 39.72,显著高于基线模型的 21.34,表明在重建过程中,模型能够综合利用更大范围的输入像素信息。结合 PSNR 与 SSIM 的同步提升(29.50 dB / 0.9718),这进一步验证了 RoPE-2D 与 LIA 在结构保持与细节恢复方面的有效协同。

综上,LAM 可视化与 DI 指标均从可解释性视角印证了模块设计的合理性。RoPE-2D 通过引入旋转不变的位置编码增强了空间感知,LIA 则通过局部信息自适应聚合扩展了特征感受野,两者协同作用使模型在保持边缘锐度的同时,更充分地利用上下文特征,从而实现了更高质量的重建表现。

3 结论

本文提出了一种新型的基于 Mamba 架构的遥感图像超分辨率网络—LAMA-SR。该模型在选择性状态空间主干网络中引入旋转位置编码(RoPE-2D)与局部交互聚合模块(LIA),能够在捕获全局依赖关系的同时,有效保留细粒度的局部细节信息。基于四个典型遥感数据集的实验结果表明,RoPE-2D 能在二维特征空间中引入显式的相对位置信息,使模型在结构保持与纹理细节恢复方面较基线模型均有显著提升;LIA 模块通过自适应局部特征聚合机制,使网络在细粒度区域的表达更充分,边缘与纹理重建更加自然。当两者协同作用时,模型在 PSNR、SSIM 等定量指标上均取得最佳表现,验证了位置建模与局部聚合机制在空间感知与重建一致性方面的互补效应。

基于 LAM(Local Attribution Map)的可解释性分析进一步揭示了模型的内部响应规律。相较于基线模型,集成 RoPE-2D 的模型在归因图上展现出更集

中的结构性响应,而引入 LIA 后,响应区域显著扩展,能够覆盖更广泛的上下文特征。当两模块联合使用时,模型在归因分布上呈现出兼具聚焦性与扩散性的特征,反映出在局部与全局依赖建模之间实现了良好的平衡。结合 DI(Diffusion Index)指标的量化结果可见,LAMA-SR 在特征扩散范围与信息利用效率方面均优于其他配置,表现出更强的空间关联建模能力。

尽管模型在重建性能与可解释性方面均取得显著提升,研究仍存在一定局限。LIA 模块在高分辨率任务中引入了额外的计算开销,对推理效率造成一定影响;RoPE-2D 的位置编码形式仍基于固定函数映射,对复杂几何变换或非规则地物场景的适应性有限。此外,LAM 分析主要聚焦于局部特征层面的响应,尚未充分揭示深层语义特征间的关系。未来的工作将进一步优化位置编码与聚合机制的可学习性,并探索其在多模态融合与实时遥感应应用中的潜力,以实现更高效、更具泛化能力的超分辨率重建模型。

参考文献(References)

- Cai H, Li J, Hu M, Gan C and Han S. 2023. EfficientViT: Lightweight Multi-Scale Attention for High-Resolution Dense Prediction // Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). IEEE: 17256 - 17267 [DOI: 10.1109/ICCV51070.2023.01587]
- Chen W, Luo L, Qu S and Dang C. 2025. MFEM: Multiscale Frequency-Enhanced Mamba for lightweight remote sensing image super-resolution. IEEE Geoscience and Remote Sensing Letters, 22: 1-5 [DOI: 10.1109/LGRS.2025.3587491]
- Chen X, Wang X, Zhou J, Qiao Y and Dong C. 2023. Activating more pixels in image super-resolution transformer // Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition

- tion (CVPR). Vancouver, BC, Canada: IEEE: 22367-22377. [DOI: 10.1109/CVPR52729.2023.02142]
- Dai D and Yang W. 2011. Satellite image classification via two-layer sparse coding with biased image representation. *IEEE Geoscience and Remote Sensing Letters*, 8 (1) : 173-176 [DOI: 10.1109/LGRS.2010.2055033]
- Dong C., Loy C C., He K M and Tang X O. 2016. Image super-resolution using deep convolutional networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38 (2) : 295-307 [DOI: 10.1109/TPAMI.2015.243928]
- Dong Z., Wang M., Wang Y., Zhu Y and Zhang Z. 2019. Object detection in high resolution remote sensing imagery based on convolutional neural networks with suitable object scale features. *IEEE Transactions on Geoscience and Remote Sensing*, 58 (3) : 2104-2114. [DOI: 10.1109/TGRS.2019.2953119]
- Gao X B., Lu W., Tao D C and Li X L. 2009. Image quality assessment based on multiscale geometric analysis. *IEEE Transactions on Image Processing*, 18 (7) : 1409 - 1423 [DOI: 10.1109/TIP.2009.2018014]
- Ghaffarian S., Kerle N., and Filatova T. 2018. Remote sensing-based proxies for urban disaster risk management and resilience: A review. *Remote sensing*, 10 (11) : 1760. [DOI: 10.3390/rs10111760]
- Gu A and Dao T. 2023. Mamba: Linear-time sequence modeling with selective state spaces [EB/OL]. [2025-11-15]. <https://arxiv.org/abs/2312.00752>
- Gu A., Goel K and Ré C. 2021. Efficiently modeling long sequences with structured state spaces [EB/OL]. [2025-11-15]. <https://arxiv.org/abs/2111.00396>
- Gu J and Dong C. 2021. Interpreting Super-Resolution Networks with Local Attribution Maps//Proceedings of the IEEE/CVF Conference on Computer Vision. Nashville, TN, USA: IEEE: 9195-9204 [DOI: 10.1109/CVPR46437.2021.00908]
- Guo H., Li J., Dai T., Ouyang Z., Ren X., and Xia S T. 2024. MambaIR: A simple baseline for image restoration with state-space model//Proceedings of 2024 European conference on computer vision. Cham: Springer Nature Switzerland: 222-241 [DOI: 10.1007/978-3-031-72649-1_13]
- He K., Zhang X., Ren S and Sun J. 2016. Deep residual learning for image recognition// Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Las Vegas, USA: IEEE: 770 - 778 [DOI: 10.1109/CVPR.2016.90]
- Helber P., Bischke B., Dengel A., and Borth D. 2019. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7) : 2217-2226. [DOI: 10.1109/JSTARS.2019.2918242]
- Heo B., Park S., Han D., and Yun S. 2024. Rotary Position Embedding for Vision Transformer//Proceedings of 2024 European Conference on Computer Vision. Cham: Springer Nature Switzerland. 289-305 [DOI: 10.1007/978-3-031-72684-2_17]
- Hu J., Shen L and Sun G. 2018. Squeeze-and-excitation networks // Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Salt Lake City, USA: IEEE: 7132 - 7141 [DOI: 10.1109/CVPR.2018.00745]
- Jiang K., Yang M., Xiao Y., Wu J., Wang G., Feng X and Jiang J. 2025. Rep-Mamba: Re-Parameterization in Vision Mamba for lightweight remote sensing image super-resolution. *IEEE Transactions on Geoscience and Remote Sensing*, 63: 1-12 [DOI: 10.1109/TGRS.2025.3597745]
- Kim J., Lee J K and Lee K M. 2016. Accurate image super-resolution using very deep convolutional networks // Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Las Vegas, USA: IEEE: 1646-1654 [DOI: 10.1109/CVPR.2016.182]
- Kingma D P and Ba J. 2014. Adam: A method for stochastic optimization [EB/OL]. [2025-11-15]. <https://arxiv.org/abs/1412.6980>
- Lei S and Shi Z. 2022. Hybrid-Scale Self-Similarity Exploitation for Remote Sensing Image Super-Resolution. *IEEE Transactions on Geoscience and Remote Sensing*, 60: 1-10 [DOI: 10.1109/TGRS.2021.3069889]
- Lei S., Shi Z and Mo W. 2022. Transformer-based multistage enhancement for remote sensing image super-resolution. *IEEE Transactions on Geoscience and Remote Sensing*, 60: 1 - 11 [DOI: 10.1109/TGRS.2021.3136190]
- Li H., Deng W., Zhu Q., Guan Q and Luo J. 2024. Local-Global Context-Aware Generative Dual-Region Adversarial Networks for remote sensing scene image super-resolution. *IEEE Transactions on Geoscience and Remote Sensing*, 62: 1-14 [DOI: 10.1109/TGRS.2024.3355419]
- Li Q., Gong M., Yuan Y and Wang Q. 2022. Symmetrical feature propagation network for hyperspectral image super-resolution. *IEEE Transactions on Geoscience and Remote Sensing*, 60: 1-12 [DOI: 10.1109/TGRS.2022.3203749]
- Liang J Y., Cao J Z., Sun G L., Zhang K., Van Gool L and Timofte R. 2021. SwinIR: image restoration using swin Transformer//Proceedings of 2021 IEEE/CVF International Conference on Computer Vision Workshops. Montreal, Canada: IEEE: 1833-1844 [DOI: 10.1109/iccvw54120.2021.00210]
- Lim B., Son S., Kim H., Nah S., and Lee K M. 2017. Enhanced Deep Residual Networks for Single Image Super-Resolution//Proceedings of the 30th IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW 2017). pp. 1132 - 1140 [DOI: 10.1109/CVPRW.2017.151]
- Lu Z., Li J., Liu H., Huang C., Zhang L and Zeng T. 2022. Transformer for single image super-resolution//Proceedings of the IEEE/CVF

- Conference on Computer Vision and Pattern Recognition Workshops (CVPRW). IEEE: 456 - 465 [DOI: 10.1109/CVPRW56347.2022.00061]
- Shazeer N, Mirhoseini A, Maziarz K, Davis A, Le Q V, Hinton G E and Dean J. 2017. Outrageously Large Neural Networks: The Sparsely-Gated Mixture-of-Experts Layer//Proceedings of the 5th International Conference on Learning Representations. Toulon, France. [DOI: 10.48550/arXiv.1701.06538]
- Su J, Ahmed M, Lu Y, Pan S, Wen B, and Liu Y. 2024. Roformer: Enhanced transformer with rotary position embedding//Neurocomputing, 568: 127063 [DOI: 10.1016/j.neucom.2023.127063]
- Waswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez A N, Kaiser Ł and Polosukhin I. 2017. Attention is all you need//Proceedings of the 7th Neural Information Processing Systems. Long Beach, USA: Curran Associates Inc.: 6000-6010.
- Wang J, Wang B, Wang X, Zhao Y and Long T. 2023. Hybrid attention-based U-shaped network for remote sensing image super-resolution. IEEE Transactions on Geoscience and Remote Sensing, 61: 5612515 [DOI: 10.1109/TGRS.2023.3283769]
- Wang S, Zhou T, Lu Y and Di H. 2022. Contextual transformation network for lightweight remote-sensing image super-resolution. IEEE Transactions on Geoscience and Remote Sensing, 60: 1-13 [DOI: 10.1109/TGRS.2021.3132093]
- Wang Z, Bovik A C, Sheikh H R and Simoncelli E P. 2004. Image quality assessment: From error visibility to structural similarity. IEEE Transactions on Image Processing, 13(4): 600 - 612 [DOI: 10.1109/TIP.2003.819861]
- Wang Z, Li L, Xue Y, Jiang C, Wang J, Sun K and Ma H. 2022. FeNet: Feature enhancement network for lightweight remote-sensing image super-resolution. IEEE Transactions on Geoscience and Remote Sensing, 60: 1-12 [DOI: 10.1109/TGRS.2022.3168787]
- Wu T, Zhao R, Lv M, Jia Z, Li L, Liu M and Vivone G. 2025. Efficient Mamba-Attention Network for remote sensing image super-resolution. IEEE Transactions on Geoscience and Remote Sensing, 63: 1-14 [DOI: 10.1109/TGRS.2025.3578879]
- Xia G S, Hu J, Hu F, Shi B, Bai X, Zhong Y, Zhang L and Lu X. 2017. AID: A benchmark data set for performance evaluation of aerial scene classification. IEEE Transactions on Geoscience and Remote Sensing, 55(7): 3965 - 3981 [DOI: 10.1109/TGRS.2017.2685945]
- Xiao Y, Yuan Q, Jiang K, He J, Lin C W and Zhang L. 2024. TTST: A top-k token selective transformer for remote sensing image super-resolution. IEEE Transactions on Image Processing, 33: 738 - 752 [DOI: 10.1109/TIP.2023.3349004]
- Yang Y and Newsam S. 2010. Bag-of-visual-words and spatial extensions for land-use classification//Proceedings of the 18th SIGSPATIAL International Conference on Advances in Geographic Information Systems. San Jose, California, USA: ACM: 270 - 279 [DOI: 10.1145/1869790.1869829]
- Zhang Y L, Li K P, Li K, Wang L C, Zhong B N, and Fu Y. 2018. Image super-resolution using very deep residual channel attention networks//Proceedings of the 2018 European Conference on Computer Vision. Munich, Germany: Springer: 286-301 [DOI: 10.1007/978-3-030-01234-2_18]
- Zhao M, Si F, Wang Y, Zhou H, Wang S, Jiang Y, and Liu W. 2020. First year on-orbit calibration of the Chinese environmental trace gas monitoring instrument onboard GaoFen-5. IEEE Transactions on Geoscience and Remote Sensing, 58(12): 8531-8540. [DOI: 10.1109/TGRS.2020.2988573]
- Zhi R, Fan X and Shi J. 2024. MambaFormerSR: A lightweight model for remote-sensing image super-resolution. IEEE Geoscience and Remote Sensing Letters, 21: 1-5 [DOI: 10.1109/LGRS.2024.3453428]
- Zou Q, Ni L, Zhang T and Wang Q. 2015. Deep learning based feature selection for remote sensing scene classification. IEEE Geoscience and Remote Sensing Letters, 12(11): 2321 - 2325 [DOI: 10.1109/LGRS.2015.2475299]

作者简介

赵公博,男,硕士研究生,主要研究方向为遥感图像处理、超分辨率重建。E-mail:gbzhao23@m.fudan.edu.cn

于珂,女,硕士研究生,主要研究方向为多模态地物分类。E-mail:kyy23@m.fudan.edu.cn

王峰,通信作者,男,副教授,主要研究方向为SAR/ISAR成像、目标识别与无人机载遥感。E-mail:fengwang@fudan.edu.cn