

引文格式:

刘茜娜, 顾津锦, 董超. 图像背景在图像超分辨率中的作用研究 [J]. 集成技术, 2023, 12(5): 76-91.

Liu XN, Gu JJ, Dong C, et al. Investigating the function of image background in image super-resolution [J]. Journal of Integration Technology, 2023, 12(5): 76-91.

图像背景在图像超分辨率中的作用研究

刘茜娜^{1,2} 顾津锦³ 董超^{1*}

¹(中国科学院深圳先进技术研究院 深圳 518055)

²(中国科学院大学 北京 100049)

³(悉尼大学 悉尼 2006)

摘要 图像超分辨率是底层视觉领域的一项代表性任务, 相关研究发现, 图像某个像素位置的重建质量与其周围的背景有关。基于这一发现, 该文探索了通过分割输入图像解释网络的新视角, 提出了一种简单组合数据集, 该数据集信息量丰富, 但单张图中仅包含单一的纹理信息。实验证明, 与目标区域纹理相近的背景, 较有利于模型在该区域的超分辨率重建; 对比分析注意力机制与传统卷积神经网络, 结果显示, 注意力结构更能帮助网络关注长程有效信息。

关键词 超分辨率; 网络可解释性; 图像背景; 数据集

中图分类号 TP 39; TP 751.1 文献标志码 A doi: 10.12146/j.issn.2095-3135.20230215001

Investigating the Function of Image Background in Image Super-Resolution

LIU Xina^{1,2} GU Jinjin³ DONG Chao^{1*}

¹(Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, Shenzhen 518055, China)

²(University of Chinese Academy of Sciences, Beijing 100049, China)

³(The University of Sydney, Sydney 2006, Australia)

*Corresponding Author: chao.dong@siat.ac.cn

Abstract As a representative low-level vision problem, image super-resolution (SR) aims to reconstruct the high-resolution image from its low-resolution counterpart. For a long time, the analysis of SR tasks is based on the whole image, while little works observe the input partition. In this paper, we find that the restoration quality of a certain position is inseparable from its surrounding image background. This phenomenon provides us a new perspective to explain the networks by splitting the input image. We construct a new hybrid dataset, of which the foreground and background contain only one kind of texture information. And then, we prove that the similar background could benefit the network restoration. By analyzing similarity and difference between the attention

收稿日期: 2023-02-15 修回日期: 2023-03-20

作者简介: 刘茜娜, 硕士研究生, 研究方向为图像处理; 顾津锦, 博士研究生, 研究方向为图像处理; 董超(通讯作者), 研究员, 研究方向为计算机视觉, E-mail: chao.dong@siat.ac.cn。

mechanism and the traditional CNN network, we show that the attention structure could help the network focus on long-range effective information. Moreover, a data enhancement method to improve the network final performance and potential future works are also proposed.

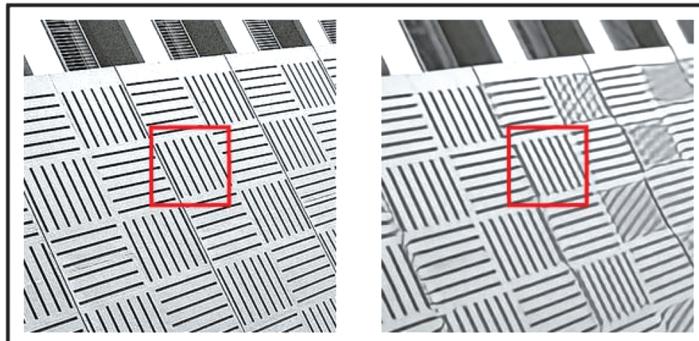
Keywords super-resolution; network interpretability; image background; dataset

1 引言

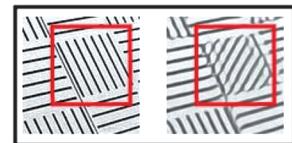
图像超分辨率(image super-resolution, SR),以下简称“超分”,是一项经典的底层视觉任务,旨在从低分辨率输入中恢复高分辨率图像。继单帧超分辨率卷积神经网络^[1]成功将卷积神经网络引入超分任务后,许多工作设计了新的网络结构^[2-6],大大提升了网络的拟合能力。近年来,随着注意力机制^[7-9]和变换神经网络(transformer neural networks, Transformer)结构^[10-12]的加入和改进,超分模型取得了新一轮的性能飞跃。上述工作虽取得了新的进展,但他们成功的原因仍然神秘,研究人员只能看到测试结果,无法解释模型的行为。上述仅依靠性能驱动的网络分析方式,限制更好的模型结构诞生。

在图1中,针对红色方框内的区域,仅改变方框周围的补充像素,增强型深度残差网络(enhanced deep residual networks, EDSR)^[5]给

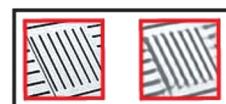
出了大相径庭的恢复结果。这一现象让研究者对方框周围像素的作用产生了疑问:邻域的加入是否必要?有些相邻的像素是否对图像重建有害?什么样的邻域对恢复结果有益?上述问题启发本项研究从一个全新的视角分析超分网络的重建过程:通过分割输入图像来观察和解释网络。本文选择输入图像的中心区域作为分析对象“前景”,将周围的其他像素视为分析“背景”。实验证明,有些邻域的出现是不必要的,甚至对结果有害。此外,为支持该分析方法,本文在自然图像之外制作了一个户外场景简单组合数据集(outdoor scene simple combined dataset, OSSCD)。这些数据的特殊之处在于,作为前景的中心区域与背景来自不同的图像。本项工作通过改变前景的邻域探索超分任务中背景的影响,借助新的分析视角和数据集可发现,当网络重建一个区域时,附近的像素与该区域越相似,该区域的性能越好。经实验证,该



(a) 左图是正确图片,右图是输入为 64×64 时的输出结果



(b) 左图是正确图片,右图是输入为 24×24 时
(输入去掉红框外部分像素) 的输出结果



(c) 左图是正确图片,右图是输入为 16×16 时
(输入去掉红框外全部像素) 的输出结果

图1 EDSR的超分结果

Fig. 1 SR results of EDSR

结论对领域内几个代表性的网络结构均成立。代表性网络指：EDSR^[5]，深度残差通道卷积网络（very deep residual channel attention networks, RCAN）^[9]，和变换神经网络（image restoration using swin transformer, SwinIR）^[10]。本文研究团队将该结论应用于网络行为的解释，成功发现了注意力机制和 Transformer 结构与传统卷积神经网络（convolutional neural networks, CNN）的异同，注意力结构的确能帮助网络关注长程有效信息。

本文的主要贡献有：(1)提出了新的分析视角——将输入图像拆分分析；(2)制作了更适合解释网络的 OSSCD；(3)拓展解释网络结构和网络行为，提出并简单验证了本项工作未来可能的发展方向。

2 国内外的研究现状

2.1 超分辨率

作为底层视觉中的一项代表性任务，SR 的目的是通过学习高分辨率和低分辨率图像对的映射，从低分辨率输入中重建高分辨率图像。超分辨率卷积神经网络^[1]是首个基于 CNN 提出的超分网络。此后，研究人员开发了大量超分深度学习模型^[2-6]，包括残差体系结构^[13]、循环体系结构^[14]等。在注意力机制得到推广后，许多网络增加了对注意力模块的设计和使用^[7-9]。研究人员认为，注意力机制可帮助网络获得关于恢复目标更详细的信息，并抑制网络利用其他的无用信息。事实上，与传统的 CNN 性能相比，基于注意力的网络性能更好。近年来，Transformer 在超分领域也取得了突破性进展^[10-12]，得到更清晰的恢复结果已成为现实。在上述网络中，SwinIR^[10]利用移位窗口对长程依赖进行建模，并因性能优异受到广泛关注。然而，迄今为止，超分领域的研究人员仍不清楚网络性能提升的原因。

2.2 可解释性

自深度学习被应用于计算机视觉以来，许多工作开始关注神经网络的可解释性。在高级视觉任务中，对深度神经网络的解释方法已得到较为广泛的关注和应用^[15-16]。如在卷积神经网络的基础上，引入反卷积神经网络^[17]，即利用卷积层中卷积核的转置进行反卷积，进而可视化出每个网络层学习到的特征，该方法可对隐层特征进行有效的定性分析。此外，还有研究在深度网络中引入注意力机制^[18-20]，即在不影响模型效果的前提下，引入注意力向量，对特征及网络中的隐层特征赋予不同的权重，并在训练过程中对该权重进行学习，可得到各个特征对模型学习的重要性，从而达到解释模型的效果。归因（显著性图）是解释网络较为优异的方法，它可视化了哪部分输入负责模型预测的结果。基于显著性图，许多方法^[21-25]获得了人类可理解的归因表示。超分网络具有深度学习和深度神经网络难以解释的性质。2021 年，Gu 等^[26]提出一种专门针对超分网络的解释方法，称为局部归因图（local attribution map, LAM），以定位影响网络输出的输入特征。此外，在底层视觉领域，Liu 等^[27]发现，在模型测试期间，特征图会根据退化类型聚集在一起。现有的解释工作有助于研究人员从不同的角度解释网络的工作机制，为更好的设计奠定了基础。但类似的研究较少，因此，加深对超分网络可解释性的探究是当前超分领域亟须解决的问题。

2.3 数据集

作为超分领域常用的训练集，DIV2K^[28]的数据质量已得到广泛认可，可帮助研究人员训练出更好的模型。此外，Set5^[29]、Set14^[30]、BSD100^[31]、Manga109^[32]和 Urban100^[33]等数据集也被广泛用于训练和测试。除上述语义复杂的数据集外，Wang 等^[34]提出了户外场景数据集（outdoor scene dataset, OST），该数据集纹理丰

富, 但单张图片仅包含单一的纹理, 如动物毛发、建筑物的砖块、水的波纹等。实验表明, OST 也能使网络学习到有效的信息。

3 图像背景的作用

3.1 超分网络中的特殊现象

超分领域内有许多用于评估图像质量的指标, 其中一种较为重要的是峰值信噪比(peak signal to noise ratio, PSNR)。如图 2 所示, 现阶段有两种常见的测试图像质量的方法: (1)将整张测试图像输入训练好的超分网络, 获得恢复结果, 并计算与原始高分辨率图像的差异, PSNR 可量化此类差异; (2)将整图切割成小块, 分别输入网络, 先将获得的结果进行组合, 再将组合结果作为最终恢复结果进行评估计算。当使用的硬件设备计算能力不足时, 研究人员通常会采用第二种测试方法。但这两种测试方法计算的 PSNR 结果可能会存在较大差距(单图差距超过 1 dB), 且第二种方法的量化数值总是较高, 但缘由并不明确。在上述两种方法中, 将整图输入和切片输入进行对比, 受影响最大的部分是小切片图的边缘, 这些像素的边缘(图 2(b)中的红线附近)变化较大, 周围像素在剪切测试时完全消失。这种“背景”的缺失是否导致两种测试方法之间的差异? 但迄今为止, 没有类似的办法帮助理解当网络重建某部分图像区域时, 相邻像素甚至远处像素带来的影响。因此, 本文将图像的目

标区域与其他区域分开研究, 以便找出未被选择的像素在超分过程中扮演的角色。

3.2 分析方法

目前, 对超分网络的研究主要集中在整张退化图像的恢复上, 很少有方法关注所选区域如何受其附近像素的影响。若该效应可得到解释, 则研究人员可了解更多的网络行为, 这对完整图像的恢复意义重大。启发本文找到新的分析视角的是解释工具 LAM, 这是一项从超分网络输入中定位重要像素的工作。LAM 选定图像中心的 16×16 局部块作为分析区域进行了归因分析, 输入图的其余部分被视为网络恢复中心域的补充输入。在选定图像的部分区域后, LAM 可确定网络为重建该部分对附近像素的利用率。受此启发, 本文建议分离图像, 操作步骤如下: 首先, 选择图像的中心区域作为前景; 然后, 在前景固定的前提下, 不断变换背景, 即为该小区域改变它的邻域, 以观察网络对前景恢复的变化。改变结果如图 3(a)所示。对于相同的训练模型, 在网络重建过程中, 仅改变区域附近的像素是否会对该部分的恢复结果产生很大影响? 如果有, 这一现象主要出现在哪些情况下? 找到、衡量和利用该影响, 更多地了解网络行为特性, 是本文设计这一分析方法的主要目的。

3.3 简单组合数据集

关于邻域像素的作用, 基于第 3.2 节中提到的方法, 在改变背景的过程中, 背景和前景之间的关系发生了变化, 从而造成中心域重建结果的



图 2 两种常见的测试方法

Fig. 2 Two common test methods

不同，这是衡量图像背景产生作用的关键。对于这种关系来说，背景和前景之间存在相似性和差异性，二者的相似性和差异性与中心区域恢复程度之间的关系是现阶段需要关注的问题。但类似图3(a)中的自然图像信息含量往往较高，这给分析图片之间的关系带来了较大困难。为了更容易地将所选区域与其背景分离，并降低分析的干扰和复杂程度，需重新建立一个可控的数据集。便于衡量图像背景对图像中心影响的新数据集必须满足两个特征：(1)单一前景到替换背景的多样性。在确定中心区域的内容后，需要为该前景替换多个背景。对于某个前景而言，良好的性能(高PSNR)表示当前背景对模型恢复这部分纹理相对有利，低PSNR则含义相反。(2)前景和背景纹理的单调性。使用自然图像进行分析的困难在于：若背景纹理复杂多样，则即使发现中心区域因替换某个背景而显示出良好的结果，也很难确定背景中具体是哪个部分起到了关键作用。此时，建立图像之间相关性的方法会受到干扰，结果可能无法控制。综上所述，现阶段常见的超分数据集并不满足此时的分析需求。

Wang等^[34]提出使用先验类别信息解决超分纹理不真实的问题，并设计制作了OST。该数据集包含各类纹理，且每张图像中的信息单调。

基于OST，本项工作制作了OSSCD，图3(b)为部分示例。简单是指组合图片的前景和背景仅包含一种纹理信息(背景中的纹理不一定与前景中的纹理相同)，如草地、天空、建筑物等。在此类组合数据中，背景和前景的纹理丰富，但不复杂。一方面保证了信息含量，可保证网络学习有效的信息；另一方面，与自然图像相比，新数据较易分析背景和前景之间的相关性。以该组合数据集为分析对象，在中心区域的恢复过程中，模型除利用该区域自己的纹理外，只能使用背景中的另一种纹理。

3.4 方法概括

综上所述，基于超分任务中的一些特殊现象，本文以一种新的视角分析超分辨率网络的重建过程，即通过分割输入图像来观察和解释网络。目标是观察图像背景在图像超分辨率中的作用，确定背景是否影响像素的重建。为便于分析，本节提出一个简单组合数据集，该数据集将和自然组合数据集一起支撑下文中的实验和分析。

4 实验结果

4.1 数据集收集

为保证自然组合图片和OSSCD之间的结论



图3 组合数据集

Fig. 3 Combined dataset

一致, 下述实验将遵循先在自然组合图片中找到规律, 再扩展到 OSSCD 的原则。首先, 从 DIV2K 验证集^[24]和 Urban100^[29]中采样 300 张 128×128 的子图像, 选取 150 张作为前景源, 另外 150 张作为背景源。从前景源中, 选择中心 32×32 的区域逐一与背景源中的图像组合, 该操作为本阶段带来了 22 500 张组合图像。其次, 遵循解释复杂案例的原则, 手动删除具有重复和不可识别内容的图像, 使其中心内容具有意义。最终筛选出 100 组数据, 每组 100 张, 共 10 000 张。单个组中的前景内容相同, 仅中心之外的背景被不断替换(其中一个是该中心区域的原始背景)。此外, 本阶段以相同的方式准备了 10 000 个 OSSCD 数据(其前景和背景均来自 OST)。

4.2 背景的不同作用

第 3.1 节中提到, 将整图切割成小块后输入网络测试, 获得的重建结果可能更好。本小节将通过实验证明这一结果。LAM 中收集了 150 张对超分网络具有挑战性的图像作为分析的测试集, 大小为 256×256 。这些图像选自 DIV2K 验证集和 Urban100 中在不同超分网络之间具有低平均 PSNR 和高方差的子图像。本节使用 LAM 提出的 150 张测试集, 尝试恢复完整图像, 或舍弃部分背景, 将中间的 64×64 区域直接输入训练好的网络。另外, 选择 56×56 的中心区域, 作为目标对象, 测量其 PSNR, 并进行比较。缩小测试区域是为了消除边缘损坏带来的干扰。如表 1 所示, 许多区域从其背景中单独拿出后恢复出了更好的性能。在舍弃背景后, 近一半的中心区域(62/150)被 EDSR 恢复得更好。对于 RCAN 和 SwinIR 而言, 超过 20% 的数据有相同现象。这验证了本文之前的想法: 并非所有背景均能给前景恢复带来收益。对于带有注意力机制的 RCAN^[9]和经典的 Transformer 结构 SwinIR 来说, 网络似乎只关注了应该关注的部分, 补充像

素对中心区域的恢复损害不大。但对于传统的 CNN 网络(EDSR)来说, 盲目扩大背景范围, 很有可能损伤网络的表征能力。那么, 什么样的背景对前景恢复有害, 又是什么样的补充输入对中心区域的恢复有益呢?

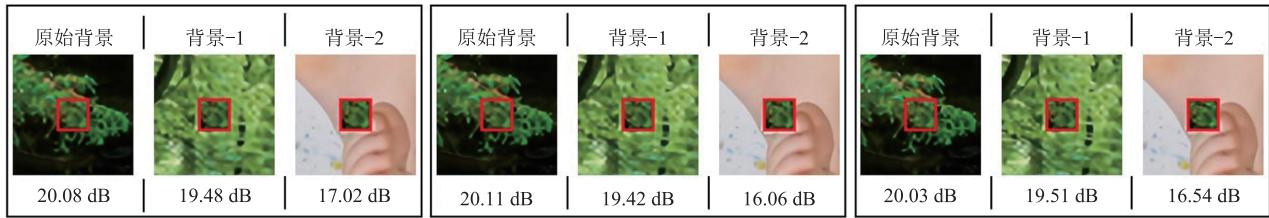
表 1 输入为全图/部分时模型对中心区域的恢复

Table 1 Models' recovery of central area when the input is full / partial

模型类别	恢复结果	
	全图更好	部分更好
EDSR	88	62
RCAN	117	33
SwinIR	111	39

4.3 相似性和恢复效果

在发现背景发挥着不同的作用后, 本小节将通过实验探索各类应用场景。这一阶段从 DIV2K 数据集中采样高分辨率图像进行训练, 最后选择经过训练的 EDSR、RCAN 和 SwinIR 模型进行测试。实验证明, 中心块在改变周边背景后, 恢复结果有了很大变化。此外, 这一阶段的实验还带来了一些新的发现。例如, 尽管数据集中为单个中心区域提供了 99 种不同的背景, 但它仍然更倾向于原始背景, 这块区域在其自身背景下恢复得更准确, PSNR 更高, 如图 4 所示。然而, 值得注意的是, 一些其他背景看上去也有助于前景的恢复。主观分析表明, “友好”的背景与前景本身非常相似。在 100 组数据中, 每个组中重建良好的中心块具有与自身高度相似的新邻居。该结论适用于本节使用的 3 个测试模型。若组合图像的前景和背景相似, 则 EDSR、RCAN 和 SwinIR 均可以很好地恢复中间部分, 图 5 证实了该结论。这说明网络在恢复某一部分区域时, 周边补充像素的输入并非越多越好, 相似的信息才是有用的, 它们能帮助网络学习和提升对类似纹理的恢复。当邻域都是无用和有害信息时, 网络利用这些像素会使中心区域面临恢复较差的困境。



注：图中所有图片均为模型输出

图 4 更换背景后的恢复结果

Fig. 4 Recovery after changing the background

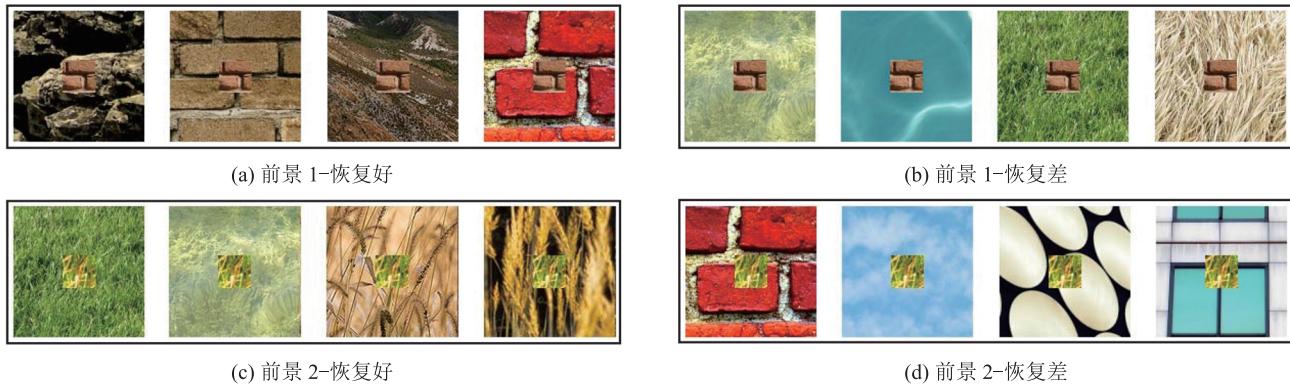


图 5 前景重建好/不好的情况

Fig. 5 Foreground with high/low performance

5 探索发现

图 6 为几个经典的超分模型对同一张输入的重建结果，模型分别为 EDSR、RCAN 和 SwinIR。显然，EDSR 和 RCAN 错误地还原了中央红框中的纹理。为解释网络错误的原因，本节对不同模型的超分结果进行了分析。就 EDSR 而言，全图的恢复效果都不够准确，除周围的像素外，没有任何位置能给中央区域条纹方向的判断带来错误的指导。但对于同一张输入而言，在前景与背景都相同的情况下，SwinIR 的判断较为准确。使 RCAN 和 EDSR 频繁发挥失常的原因是什么？Transformer 和 CNN 的表现在此类情况下有什么区别和联系？针对上述问题，本节将一幅图像拆开，并把前景和背景单独看待，以解释不同类型网络的工作机制。

不同网络之间，结构差异较大：EDSR 是一

个基本的卷积神经网络，RCAN 是带有注意力机制的 CNN，SwinIR 是一个经典的 Transformer 结构。SwinIR 的作者提到，网络中的非局部稀疏注意力结构是基于非局部注意机制设计的。研究人员认为，具有注意力结构的网络能够忽略无关信息，更加关注提升性能需要关注的关键信息^[35-36]。为分析未经验证的网络特性能否反映在上文提出的新方法中，本节使用 OST 制作了一批新数据，图 7 为 6 种方法的示例，以组合相同的前景和背景。完整图像大小为 128×128 ，其中心区域为 16×16 ，周围背景的大小分别为 0（没有外来背景的原始图像）、32、48、64、80 和 96。依照第 4.3 节中的规律，在前景和背景相似的情况下，EDSR、RCAN 和 SwinIR 的恢复效果均较好。本阶段数据集的前景和背景并不相似，随着背景越来越大，对中心恢复有用的信息离中心越来越远。利用这种方法，本实验可观察网络

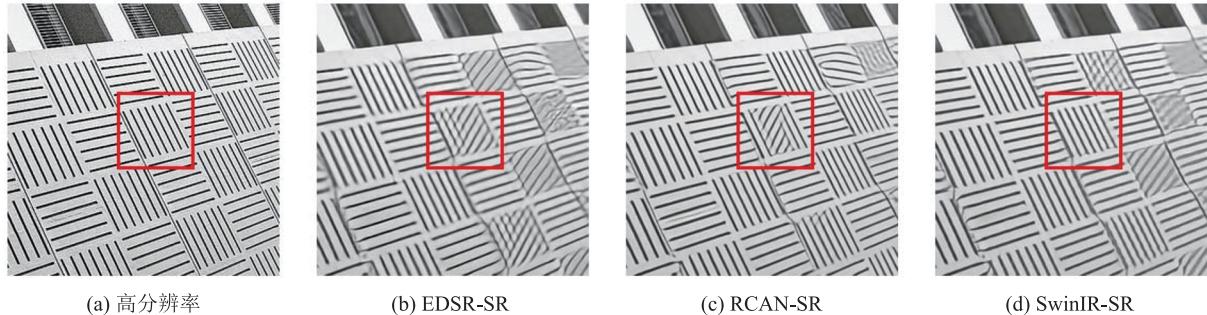


图 6 不同模型的超分结果

Fig. 6 SR results of different models

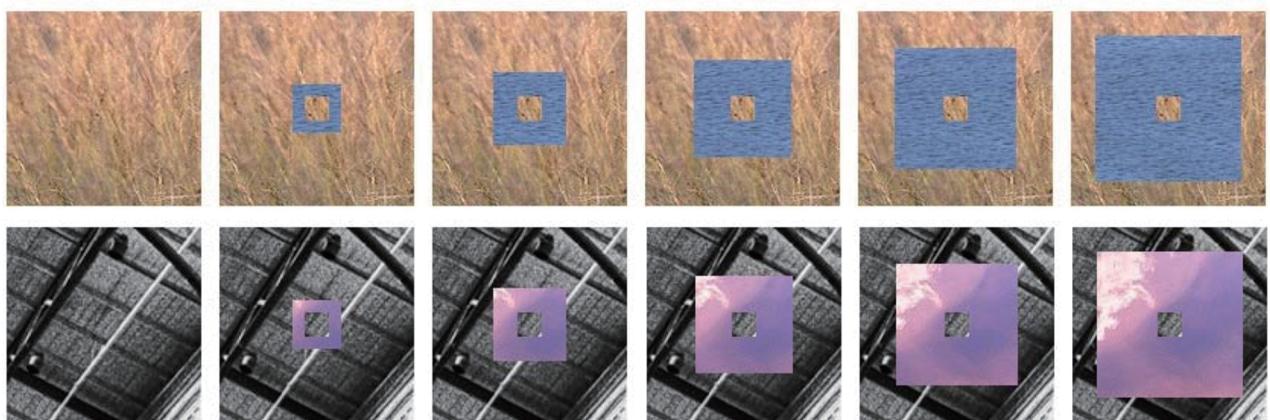


图 7 背景大小不同的新数据集

Fig. 7 Dataset with different size of background

捕获远程信息的能力。依托于新数据，本阶段获得的研究结果总结如下：

(1) PSNR 跃升。对于中心区域的恢复, EDSR 显然受到了目标附近无关背景的干扰, 5 种尺寸的外来背景均导致中心区域的恢复质量下降(21.34 dB 到 19.76 dB 及以下)。反观 RCAN, 当背景大小为 32 时, 就展现出了差异性。虽然此时中心域 PSNR 仍然没有使用原始背景时(26.82 dB)高, 但当背景从 48 缩小至 32(相似的像素更接近中心)时, PSNR 发生跃升(18.88 dB 到 21.05 dB)。此时, 注意力结构似乎帮助 RCAN 跨过无用信息、关注相似信息, 因而获得了更好的重建结果。对于 SwinIR 而言, 当背景为 32 和 48 时, PSNR 均出现了跃升(背景 48 时为 20.61 dB, 背景 32 时为 21.45 dB, 背景 0 时

为 26.33 dB)。具体数值见图 8~10。综上所述，为什么加入无关背景后，EDSR 反而是效果最好的模型呢？

(2) 对有用信息的关注。第 4.2 节介绍了超分领域新提出的解释工具 LAM。当发现中心区域的恢复随有用信息的接近而改变，同时又根据(1)中提出的问题，本文尝试使用 LAM 查看组合图片的归因图，以望有助于解释(1)中的现象。通过分析图 8 中的归因结果可知，对于 EDSR 来说，当背景大小为 32 时，网络试图争取一些有用的信息，但不可避免地利用了很多无关的信息。图中标红的像素点表示网络利用了该位置的信息来重建红色方框内的区域，颜色越深代表利用率越高。归因图显示，EDSR 同时用到背景中的像素和背景外的像素，即在恢复图像中心的过

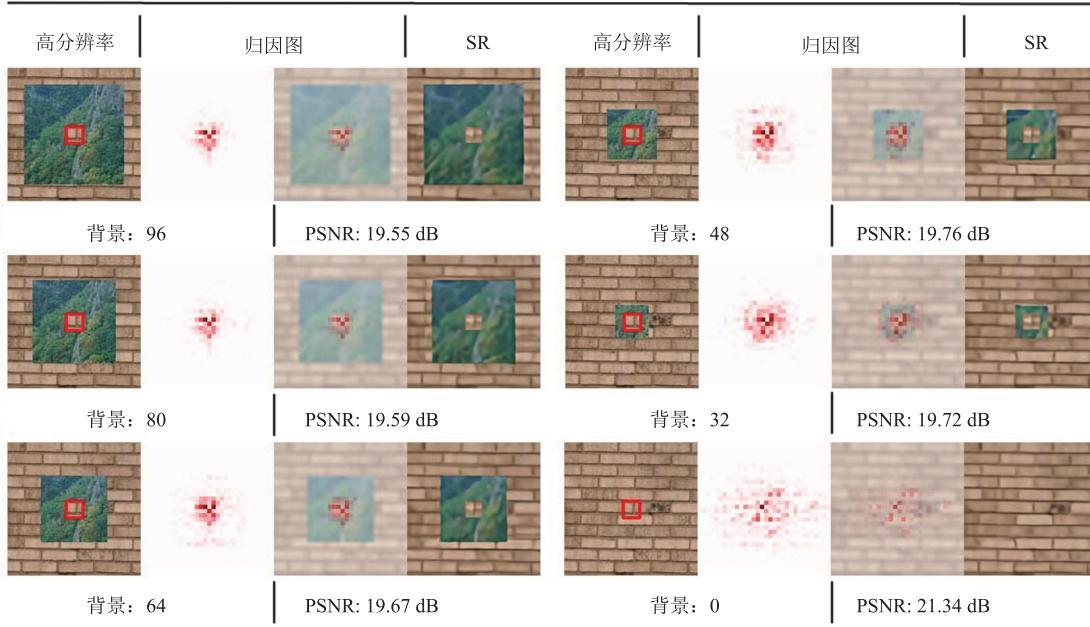


图 8 EDSR 的具体数值和归因图

Fig. 8 Specific values and attribution map of EDSR

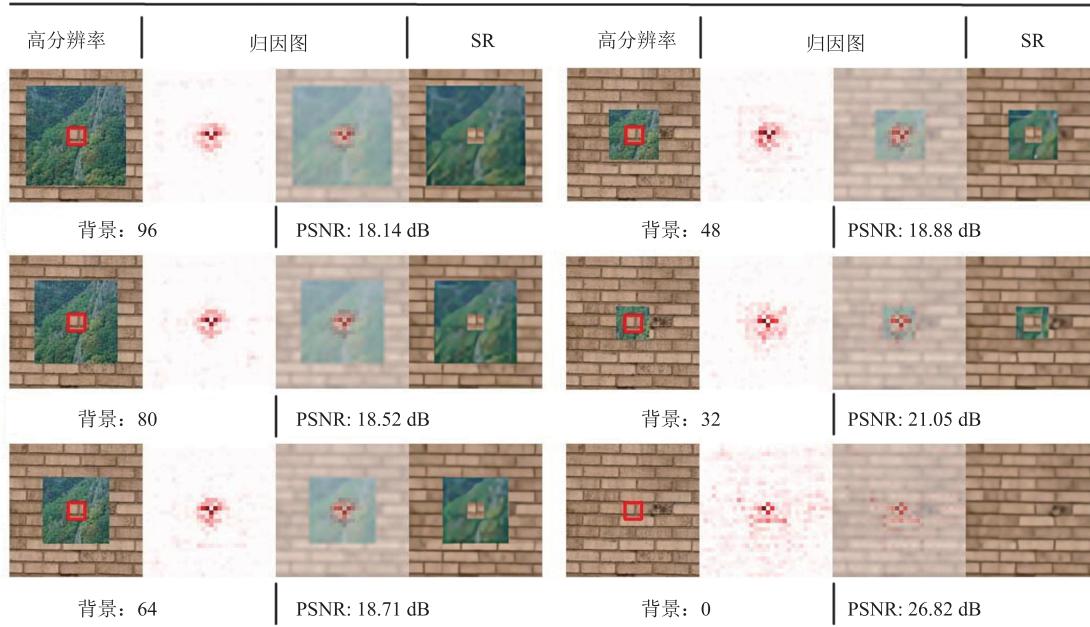


图 9 RCAN 的具体数值和归因图

Fig. 9 Specific values and attribution map of RCAN

程中, EDSR 使用了一些相似的有效信息, 但缺乏规避无效信息的能力, 从而导致中心区域的 PSNR 不高。当背景大小从 32 扩展到 48 或更大时, 网络开始忽视有效信息。显然, EDSR 的远

距离信息捕捉能力还不够强。与 EDSR 一样, RCAN 只能在背景大小为 32×32 时找到有效信息。如图 9 所示, 与 EDSR 不同的是, 注意力结构有助于 RCAN 专心关注关键信息、减少对

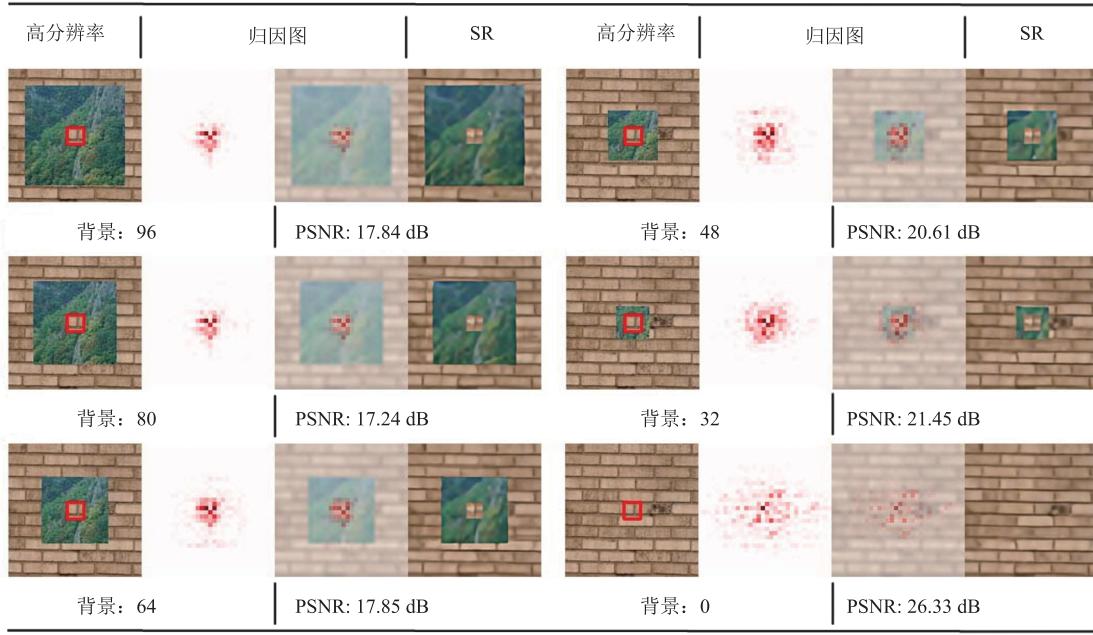


图 10 SwinIR 的具体数值和归因图

Fig. 10 Specific values and attribution map of SwinIR

无用像素的使用, 反映在归因图中为对背景的利用减少, 因此, PSNR 在此处跃升。而对 SwinIR 来说, 它的视野显然更为广阔。如图 10 所示, 当背景大小为 48 时, 网络依然可以关注到障碍之外的有用信息。此外, 当背景大小为 32 和 48 时, 中心区域的恢复结果均有所提高。不过, RCAN 和 SwinIR 对无关信息的忽略能力仍有所欠缺, 它们对长程信息的争取会随着对无关信息的利用, 最终导致模型性能的损伤。这既解释了(1)中的问题, 又说明了现阶段的注意力机制还有待提升。

本节结合解释工具 LAM, 从背景和前景单独分析的角度, 解释了模型结构设计带来的网络行为异同。注意力结构可帮助模型更专注于关键信息, 规避对无用像素的使用, 但规避能力有限。在无关信息较多的情况下, RCAN 和 SwinIR 甚至达不到 EDSR 的恢复效果。

由于 RCAN 和 SwinIR 是具备大感受野的模型, 因此, 本节不仅从局部中心区域扩散非相似背景, 还从非中心区域扩散(整体图像边缘向内

扩散)了无关区域, 如图 11~13 所示, 以更完备地研究不同背景对不同模型的影响。由图 11~13 可知, 当无关背景从整体图像边缘向内扩散时, EDSR 的性能有所下降; RCAN 作为带有注意力结构的模型, 更专注于关键信息, 规避了对无用像素的使用, 性能下降较少; 但对于具备最大感受野的 SwinIR 而言, 强大的长程信息利用能力使它的表征能力下降了近 2 dB。

6 未来方向

本项工作的关键在于发现和揭示现象, 关于这类方法在未来能否有更多的数值结论和更具体的应用, 本文进行了以下尝试和展望。

6.1 量化相似性

许多研究人员通过先提取图像本身的特征, 再计算这些特征的差异性, 以衡量图像之间的关系。已有许多从图像中提取特征的方法被广泛应用。Haralick 等^[37]提出了灰度共生矩阵, 以描述纹理特征, 还有一些工作关注图像中的颜色信

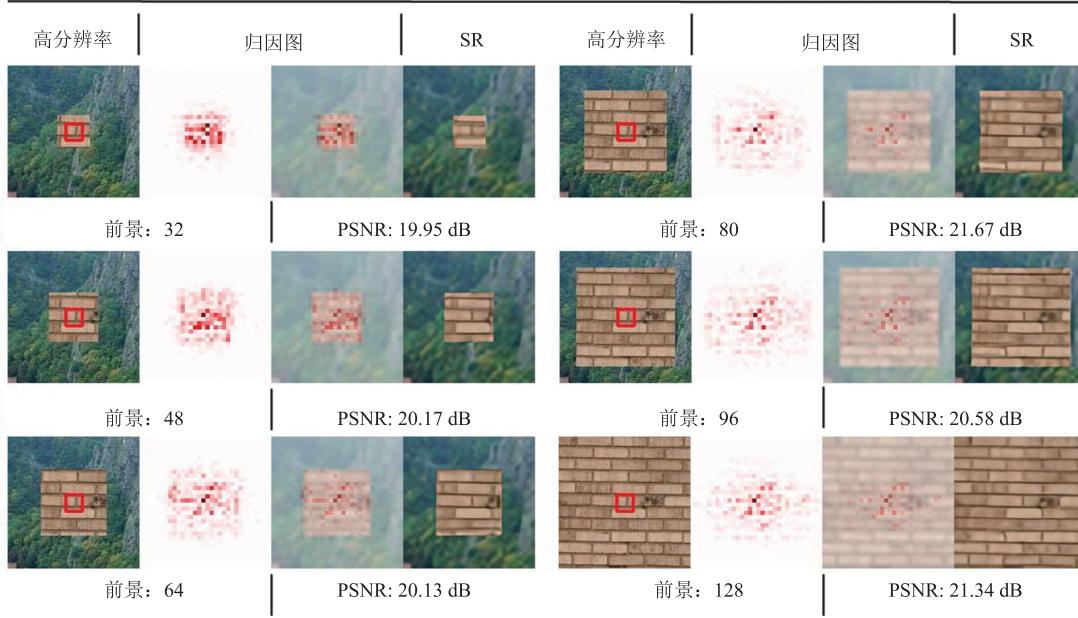


图 11 EDSR 的具体数值和归因图(方法二)

Fig. 11 Specific values and attribution map of EDSR (Method 2)

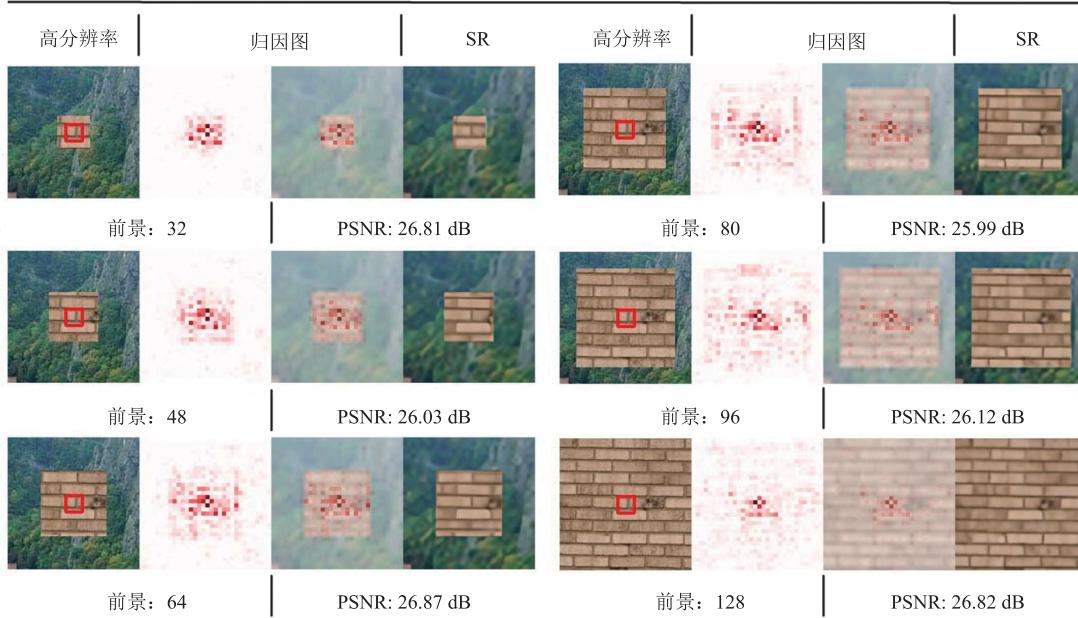


图 12 RCAN 的具体数值和归因图(方法二)

Fig. 12 Specific values and attribution map of RCAN (Method 2)

息^[38]。其中，格拉姆矩阵常被用于风格迁移的深度网络中^[39-43]，具体操作是用神经网络提取图像的浅层特征图，然后计算特征图的格拉姆矩阵值。格拉姆矩阵值结合了图像的纹理、边缘和颜

色特征。风格迁移任务的目标通常优化(降低)目标图像和风格图像之间的格拉姆差异。格拉姆矩阵具有公认的分析图像之间的相似性和差异性的功能。

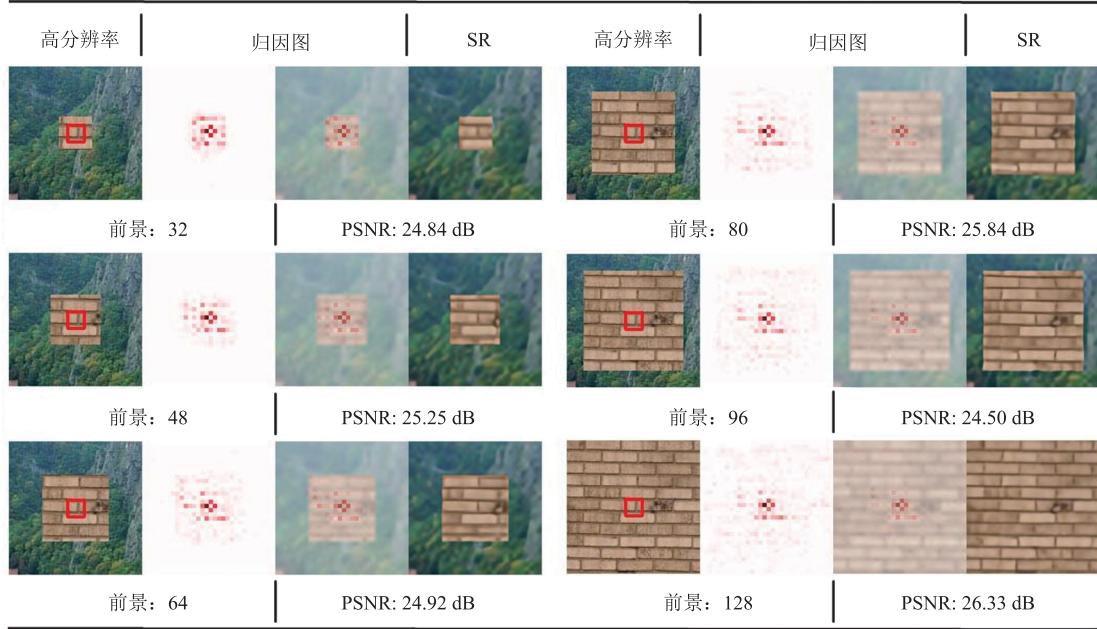


图 13 SwinIR 的具体数值和归因图(方法二)

Fig. 13 Specific values and attribution map of SwinIR (Method 2)

研究发现, 主观上相似的背景有助于前景重建, 因此, 可先使用格拉姆矩阵量化前景与背景之间的相似性, 然后找到二者相似性与前景恢复性能之间的数值关联。具体做法是先将前景和背景图像分别输入 VGG16^[44]网络, 用第一卷积层获得的特征图作为计算图像格拉姆矩阵的基础。关于这一操作, 许多研究^[45-46]中都有类似讨论, 特征图计算格拉姆矩阵是风格迁移任务中的常见操作。然而, 格拉姆矩阵只被用于掌握整个图像的一般风格, 其局限性在于, 当图片之间风格差异较小时, 它不能非常准确地反映相似性。当量化格拉姆差和中心域 PSNR 的相关性时, 结果并不明显, 找到合适的相关性指标是该项研究在下一阶段需要解决的主要问题。

6.2 数据增强方法

第 4.3 节提到, 当中心区域周围有与它类似的像素时, 该区域的性能会更好。依据这一发现, 能否通过重新排布图像, 提高图像中某区域的恢复性能? 图 14~15 展示了可能的重排结果, 即选定某个区域作为重建目标, 然后为该

目标创建高度相似的背景。背景可能由其自身组成, 也可能由其翻转、旋转和缩放的结果组成。翻转可以是水平翻转或竖直翻转, 旋转角度有顺时针 90°、180°、270° 3 种, 缩放倍数为 2 或 4, 缩放方法为最近邻插值。最终生成的图像面积为原始区域面积的 25 倍, 即长和宽均为原始图像的 5 倍。将最终生成的图像输入到训练好的网络中进行测试, 可以发现目标区域恢复得更好(PSNR 更高)。此处的 PSNR 不是针对整个画面, 仅指所选区域的恢复程度。

从该方法很容易联想到测试时数据增强(test-time augmentation, TTA), 这是一种以计算资源为代价消除干扰的做法。TTA 的具体操作为: 测试时, 将原始数据进行不同形式的增强后, 再输入网络, 如水平翻转、放大缩小和中心旋转等, 然后逆回结果, 取平均值, 作为最终输出。与 TTA 不同, 本阶段是为选定区域添加有用的邻居, 是一种不同的数据增强方式。相较于 TTA 的结果会随所选的增强方法而变化, 此处的测试只需进行一次, 结果不会因测试次数而改变。为公

平地显示这两种方法之间的差异，在消耗相同算力的前提下，本节将 TTA 的结果与所提方法的恢复结果进行比较。选择一个要重建的区域，利用 TTA 进行 25 次数据增强测试。结果表明，TTA 最多可带来 0.1 dB 的改善，但在某些情况下，由于增强方法的选取不当，结果可能比不使用 TTA 更差。本文方法不仅可以稳定地提高 PSNR，而且没有损害性能的风险，图 14 为两个显著改进的示例。然而，本文方法也有其局限性。当测试图像(图 15)随机生成时，大多数选定的测试区域(白色方框中)不会有很大的性能提高(但仍然没有降低性能的风险)。本文从 DIV2K 数据集中随机选取 200 张图片(这些数据大多和图 15 类似，并非来自人工挑选)，执行上述随机增强操作。由表 2 的统计结果可知，这种增强方法的增益稳定，没有损害性能的风险。当所选区域在原始图

表 2 输入为原图/增强图时模型对中心区域的恢复

Table 2 Models' recovery of central area when the input is original / enhanced

模型类别	恢复结果	
	原图 (dB)	增强图 (dB)
EDSR	23.02	23.28
RCAN	23.11	23.39
SwinIR	23.24	23.33

像中特殊时，即原始图像中，该区域相似的部分占比较少时，该测试增强方法对性能是友好的，网络需要“看到”更多类似的纹理来训练对此特殊区域的恢复能力。

7 结语

超分任务中某个位置的重建质量与其周围的

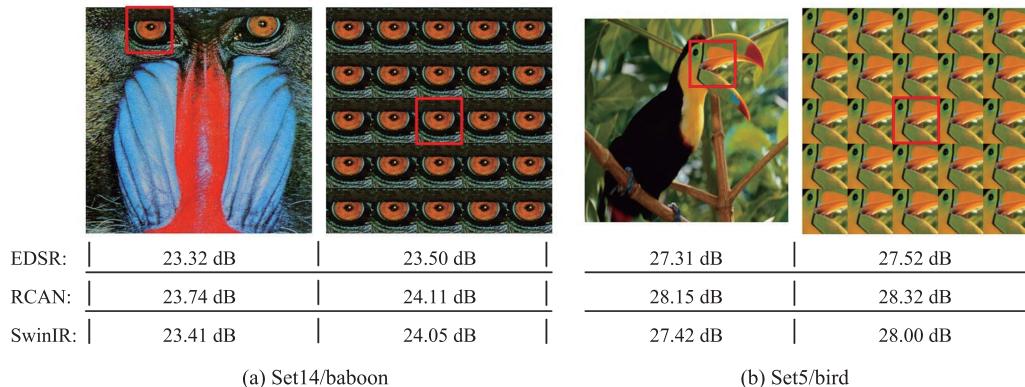


图 14 改进较大的示例

Fig. 14 Examples of great improvement

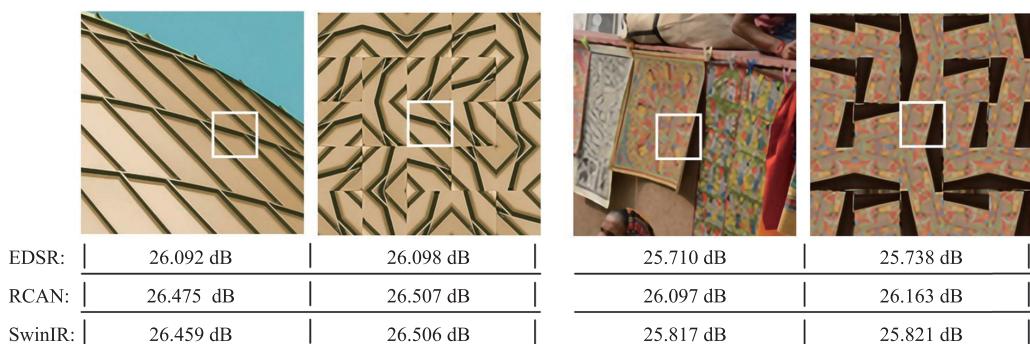


图 15 改进较小的示例

Fig. 15 Examples of less improvement

“背景”密不可分,当需要恢复的区域与其附近的像素相似时,性能更好。为证实这一结论,本研究设计了分析网络的新方法,即目标区域与其他区域的分割分析,并制作了针对该问题的数据集。借助这一发现,本文进一步探索了网络的工作机制,发掘了量化相似性的可能性,并提出一种重做图像背景的数据增强方法。希望这项工作能为超分任务的分析带来新的视角,进而帮助研究人员更好地理解网络行为,指导设计更好的网络和评估算法。

参 考 文 献

- [1] Dong C, Loy CC, He KM, et al. Image super-resolution using deep convolutional networks [J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2015, 38(2): 295-307.
- [2] Kim J, Lee JK, Lee KM. Accurate image super-resolution using very deep convolutional networks [C] // Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016: 1646-1654.
- [3] Ledig C, Theis L, Huszár F, et al. Photo-realistic single image super-resolution using a generative adversarial network [C] // Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017: 4681-4690.
- [4] Kong XT, Zhao HY, Qiao Y, et al. ClassSR: a general framework to accelerate super-resolution networks by data characteristic [C] // Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021: 12016-12025.
- [5] Lim B, Son S, Kim H, et al. Enhanced deep residual networks for single image super-resolution [C] // Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, 2017: 136-144.
- [6] Zhang YL, Tian YP, Kong Y, et al. Residual dense network for image super-resolution [C] // Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018: 2472-2481.
- [7] Chen HY, Gu JJ, Zhang Z. Attention in attention network for image super-resolution [Z/OL]. arXiv Preprint, arXiv: 2104.09497, 2021.
- [8] Dai T, Cai JR, Zhang YB, et al. Second-order attention network for single image super-resolution [C] // Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019: 11065-11074.
- [9] Zhang YL, Li KP, Li K, et al. Image super-resolution using very deep residual channel attention networks [C] // Proceedings of the European Conference on Computer Vision, 2018: 286-301.
- [10] Liang JY, Cao JZ, Sun GL, et al. SwinIR: image restoration using Swin Transformer [C] // Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021: 1833-1844.
- [11] Lu ZS, Li JC, Liu H, et al. Transformer for single image super-resolution [Z/OL]. arXiv Preprint, arXiv: 2108.11084, 2021.
- [12] Yang FZ, Yang H, Fu JL, et al. Learning texture Transformer network for image super-resolution [C] // Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020: 5791-5800.
- [13] Dong C, Loy CC, Tang XO. Accelerating the super-resolution convolutional neural network [C] // Proceedings of the Computer Vision-ECCV 2016, 2016: 391-407.
- [14] Kim J, Lee JK, Lee KM. Deeply-recursive convolutional network for image super-resolution [C] // Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016: 1637-1645.
- [15] Mahendran A, Vedaldi A. Visualizing deep convolutional neural networks using natural pre-images [J]. *International Journal of Computer Vision*, 2018, 125(1): 1-16.

- Vision, 2016, 120: 233-255.
- [16] Zhou BL, Bau D, Oliva A, et al. Interpreting deep visual representations via network dissection [J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2018, 41(9): 2131-2145.
- [17] Yosinski J, Clune J, Nguyen A, et al. Understanding neural networks through deep visualization [Z/OL]. arXiv Preprint, arXiv: 1506.06579, 2015.
- [18] Zhao B, Wu X, Feng JS, et al. Diversified visual attention networks for fine-grained object classification [J]. IEEE Transactions on Multimedia, 2017, 19(6): 1245-1256.
- [19] Xiao TJ, Xu YC, Yang KY, et al. The application of two-level attention models in deep convolutional neural network for fine-grained image classification [C] // Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015: 842-850.
- [20] Jaderberg M, Simonyan K, Zisserman A, et al. Spatial Transformer networks [C] // Proceedings of the 28th International Conference on Neural Information Processing Systems-Volume 2, 2015: 2017-2025.
- [21] Lundberg SM, Lee SI. A unified approach to interpreting model predictions [C] // Proceedings of the 31st International Conference on Neural Information Processing Systems, 2017: 4768-4777.
- [22] Shrikumar A, Greenside P, Kundaje A. Learning important features through propagating activation differences [C] // Proceedings of the International Conference on Machine Learning, 2017: 3145-3153.
- [23] Simonyan K, Vedaldi A, Zisserman A. Deep inside convolutional networks: visualising image classification models and saliency maps [Z/OL]. arXiv Preprint, arXiv:1312.6034, 2013.
- [24] Springenberg JT, Dosovitskiy A, Brox T, et al. Striving for simplicity: the all convolutional net [Z/OL]. arXiv Preprint, arXiv: 1412.6806, 2014.
- [25] Sundararajan M, Taly A, Yan QQ. Axiomatic attribution for deep networks [C] // Proceedings of the 34th International Conference on Machine Learning, 2017: 3319-3328.
- [26] Gu JJ, Dong C. Interpreting super-resolution networks with local attribution maps [C] // Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021: 9199-9208.
- [27] Liu YH, Liu AR, Gu JJ, et al. Discovering “semantics” in super-resolution networks [Z/OL]. arXiv Preprint, arXiv: 2108.00406, 2021.
- [28] Agustsson E, Timofte R. NTIRE 2017 challenge on single image super-resolution: dataset and study [C] // Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, 2017: 126-135.
- [29] Bevilacqua M, Roumy A, Guillemot C, et al. Low-complexity single-image super-resolution based on nonnegative neighbor embedding [C] // Proceedings of the 23rd British Machine Vision Conference, 2012: 135.1-135.10.
- [30] Yang JC, Wright J, Huang TS, et al. Image super-resolution via sparse representation [J]. IEEE Transactions on Image Processing, 2010, 19(11): 2861-2873.
- [31] Martin D, Fowlkes C, Tal D, et al. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics [C] // Proceedings of the Eighth IEEE International Conference on Computer Vision, 2001: 416-423.
- [32] Matsui Y, Ito K, Aramaki Y, et al. Sketch-based manga retrieval using Manga109 dataset [J]. Multimedia Tools and Applications, 2017, 76: 21811-21838.
- [33] Huang JB, Singh A, Ahuja N. Single image super-resolution from transformed self-exemplars [C] // Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015:

- 5197-5206.
- [34] Wang XT, Yu K, Dong C, et al. Recovering realistic texture in image super-resolution by deep spatial feature transform [C] // Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018: 606-615.
- [35] Cohen MR, Maunsell JHR. Attention improves performance primarily by reducing interneuronal correlations [J]. *Nature Neuroscience*, 2009, 12(12): 1594-1600.
- [36] Muqeet A, Iqbal MTB, Bae SH. HRAN: hybrid residual attention network for single image super-resolution [J]. *IEEE Access*, 2019, 7: 137020-137029.
- [37] Haralick RM, Shanmugam K, Dinstein IH. Textural features for image classification [J]. *IEEE Transactions on Systems, Man, and Cybernetics*, 1973, SMC-3(6): 610-621.
- [38] Han J, Ma KK. Fuzzy color histogram and its use in color image retrieval [J]. *IEEE Transactions on Image Processing*, 2002, 11(8): 944-952.
- [39] Gatys LA, Ecker AS, Bethge M, et al. Controlling perceptual factors in neural style transfer [C] // Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017: 3985-3993.
- [40] Gupta A, Johnson J, Alahi A, et al. Characterizing and improving stability in neural style transfer [C] // Proceedings of the IEEE International Conference on Computer Vision, 2017: 4067-4076.
- [41] Jing YC, Yang YZ, Feng ZL, et al. Neural style transfer: a review [J]. *IEEE Transactions on Visualization and Computer Graphics*, 2019, 26(11): 3365-3385.
- [42] Li YH, Wang NY, Liu JY, et al. Demystifying neural style transfer [Z/OL]. arXiv Preprint, arXiv: 1701.01036, 2017.
- [43] Li YJ, Fang C, Yang JM, et al. Universal style transfer via feature transforms [C] // Proceedings of the 31st International Conference on Neural Information Processing Systems, 2017: 385-395.
- [44] Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition [Z/OL]. arXiv Preprint, arXiv: 1409.1556, 2014.
- [45] Shen FL, Yan SC, Zeng G. Neural style transfer via meta networks [C] // Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018: 8061-8069.
- [46] Gatys LA, Ecker AS, Bethge M. Texture synthesis using convolutional neural networks [J]. arXiv Preprint, arXiv: 1505.07376, 2015.