



武汉大学学报(信息科学版)

Geomatics and Information Science of Wuhan University

ISSN 1671-8860, CN 42-1676/TN

《武汉大学学报(信息科学版)》网络首发论文

题目：融合注意力与序列单元的文本超分辨率
作者：韦豪东，易尧华，余长慧，林立宇
DOI：10.13203/j.whugis20220158
收稿日期：2022-07-14
网络首发日期：2022-08-19
引用格式：韦豪东，易尧华，余长慧，林立宇. 融合注意力与序列单元的文本超分辨率[J/OL]. 武汉大学学报(信息科学版). <https://doi.org/10.13203/j.whugis20220158>



网络首发：在编辑部工作流程中，稿件从录用到出版要经历录用定稿、排版定稿、整期汇编定稿等阶段。录用定稿指内容已经确定，且通过同行评议、主编终审同意刊用的稿件。排版定稿指录用定稿按照期刊特定版式（包括网络呈现版式）排版后的稿件，可暂不确定出版年、卷、期和页码。整期汇编定稿指出版年、卷、期、页码均已确定的印刷或数字出版的整期汇编稿件。录用定稿网络首发稿件内容必须符合《出版管理条例》和《期刊出版管理规定》的有关规定；学术研究成果具有创新性、科学性和先进性，符合编辑部对刊文的录用要求，不存在学术不端行为及其他侵权行为；稿件内容应基本符合国家有关书刊编辑、出版的技术标准，正确使用和统一规范语言文字、符号、数字、外文字母、法定计量单位及地图标注等。为确保录用定稿网络首发的严肃性，录用定稿一经发布，不得修改论文题目、作者、机构名称和学术内容，只可基于编辑规范进行少量文字的修改。

出版确认：纸质期刊编辑部通过与《中国学术期刊（光盘版）》电子杂志社有限公司签约，在《中国学术期刊（网络版）》出版传播平台上创办与纸质期刊内容一致的网络版，以单篇或整期出版形式，在印刷出版之前刊发论文的录用定稿、排版定稿、整期汇编定稿。因为《中国学术期刊（网络版）》是国家新闻出版广电总局批准的网络连续型出版物（ISSN 2096-4188，CN 11-6037/Z），所以签约期刊的网络版上网络首发论文视为正式出版。

DOI:10.13203/j.whugis20220158

引用格式：

韦豪东, 易尧华, 余长慧, 等. 融合注意力与序列单元的文本超分辨率[J]. 武汉大学学报·信息科学版, 2022, DOI: 10.13203/j.whugis20220158 (WEI Haodong, YI Yaohua, YU Changhui, et al. Text Super-resolution Method with Attentional Mechanism and Sequential Units[J]. *Geomatics and Information Science of Wuhan University*, 2022, DOI: 10.13203/j.whugis20220158)

融合注意力与序列单元的文本超分辨率

韦豪东¹ 易尧华¹ 余长慧¹ 林立宇²

1 武汉大学遥感信息工程学院, 湖北 武汉, 430079

2 武汉大学测绘遥感信息工程国家重点实验室, 湖北 武汉, 430079

摘要：街景影像中的文本信息是感知与理解场景的关键线索，低分辨率街景影像文本区域细节缺乏导致文本识别准确率降低。文本超分辨率通过增强文本区域边缘及纹理细节提高文本识别准确率，本文提出融合注意力与序列单元的街景影像文本超分辨率方法。首先采用混合残差注意力结构提取影像文本区域空间信息、通道信息并融合特征，序列单元通过双向门控循环结构提取影像中文本间序列先验信息；再利用梯度先验知识作为约束条件，重构街景影像文本区域。本文采用TextZoom真实场景影像及合成文本影像进行对比分析，试验结果表明超分辨率重构的街景影像文本区域边缘清晰、纹理细节丰富，可以提高街景影像文本识别准确率。

关键词：街景影像；超分辨率；注意力机制；序列信息；梯度先验损失

中国分类号：P237

文献标识码：A

街景影像中的文本包含丰富语义信息，是感知与理解场景的关键线索。街景影像中的文本通常形状不规则、背景复杂且分辨率低，这些都直接导致了街景影像文本识别准确率降低^[1]。超分辨率重建(super-resolution, SR)算法对影像文本区域进行超分辨率预处理，通过提高影像文本区域清晰度提升文本识别准确率^[2]。

基于浅层卷积神经网络(convolutional neural networks, CNN)的文本超分辨率算法学习低分辨率(low resolution, LR)文本到高分辨率(high resolution, HR)文本的映射。Dong 等人^[3]在 2015 国际文档分析与识别比赛^[4]文本超分辨率赛道上首次提出基于 CNN 的文本超分辨率重建方法

(super-resolution CNN, SRCNN)^[5]。文献[6]使用转置卷积和亚像素卷积对网络提取的特征图进行上采样，处理黑白文档超分辨率问题。文献[7]分别训练文本 SRCNN 分支与影像 SRCNN 分支，融合双分支提升场景图像超分辨率重建质量。这些方法卷积层数较少，无法充分学习 LR 文本和 HR 文本之间复杂的映射关系，也没有充分利用文本区域先验知识。

近年来深层 CNN 成功应用于影像文本超分辨率任务。文献[8]引入自然场景图像抠图算法对提取的文本、前景和背景层在两路分别重建文本区域边缘与颜色信息，并使用文本内容作为监督标签训练网络。文献[1]提出使用文本内容识别结果构建文本感知损失训练网络。文献[2]构建了真实自然场景影像文本数据集 TextZoom，并提出了融合序列残差块

收稿日期：2022-07-14

项目资助：国家重点研发计划(2021YFB2206200)。

第一作者：韦豪东，硕士生，主要从事图像超分辨率重建及相关应用的理论与方法研究。WUWeihaodong@163.com

通讯作者：易尧华，博士，教授。yyh@whu.edu.cn

与中央对齐模块的文本超分辨率重建网络（text super-resolution network, TSRN）验证其有效性。上述方法从不同角度使网络模型专注于处理场景影像的文本区域，增强 LR 文本细节，提升文本识别准确率。这些网络普遍采用通用图像超分辨率重建网络结构，没有充分利用图像中文本区域与背景差异、文本之间语义序列关系等先验信息。

本文提出融合混合残差注意力机制与序列单元的文本超分辨率重建网络（hybrid sequential residual attentional text SR network, HSRATN）。为了使模型关注场景影像文本区域以充分学习文本区域先验知识，提出了混合残差注意力结构，融合空间注意力机制和拉普拉斯通道注意力机制，利用特征图在空间和通道上各自的依赖关系，学习文本部分的多层次特征表示。此外，为了利用影像中文本字符间的序列先验信息，进一步提出混合序列残差注意力模块，在混合残差注意力结构中融合包含双向门控循环结构的序列单元，提取字符序列关系。采用梯度先验损失函数衡量 HR 文本图像与 SR 文本图像在梯度场上的差距，并对网络的中间特征图进行可视化，验证梯度先验损失对提高文本边缘重建效果的有效性。

1 相关工作

1.1 超分辨率重建

图像超分辨率重建指从退化的 LR 图像重建生成对应的 HR 图像^[9]，广泛应用于文档图像^[6]、街景影像^[2]、人脸图像^[10]、遥感影像^[11]重建等。基于学习的 SR 从样本数据中学习先验知识，模型一方面通过加深网络来扩大网络感受野，提升性能^[5]；另一方面则通过残差连接、密集连接等结构提取并保存各层次特征图的信息，提高模型训练效率^[12]。

1.2 注意力机制

注意力机制^[13]在自然语言处理领域取得成功应用后，也开始应用于计算机视觉语义分割^[14]、图像分类^[15]和图像识别^[16]等任务。文献^[17]首次将注意力机制应用于图像超分辨率重建，利用通道注意力机制学习通道间相互关系来为各通道赋予权重，忽略低频信息而加强高频信息。文献^[18]提出混合残差注意力网络（hybrid residual attentional network, HRAN），融合空间注意力和通道注意力机制，学习特征图空间和通道之间的关系。文献^[19]将多注意力结构应用于文本超分辨率任务。文献^[20]提出密集残差拉普拉斯注意力网络（densely residual laplacian SR network, DRLN），在通道注意力中增加拉普拉斯金字塔结构，增强网络特征学习能力。注意力机制可以根据特征的相对重要性对其赋予权重，为了充分学习影像中文本区域特征，本文算法利用注意力机制加强文本高频信息。

2 本文方法

2.1 超分辨率网络结构

如图 1 所示，本文算法网络结构包括预处理网络、推理网络与重建网络，选用 HRAN 作为基础模型进行优化。

预处理网络由两个核大小为 3×3 的卷积层组成，将输入从图像空间映射到特征空间，进行初级特征提取。通过适应性图像阈值算法生成 RGB 图像的二值化语义掩膜，对两者在通道维度上连结构成 RGBM 四通道特征图，将其作为网络的输入。适应性图像阈值算法首先将图像转为灰度图像，然后计算图像的平均亮度，将大于平均亮度的区域像素值变成 0 而其它区域像素值变成 255^[19]。如图 2 所示，二值化语义掩膜可以看作图像文本区域和背景区域的语义分割图，作为语义先验输入网络有助于增强文本重建效果并提高模型学习效率。

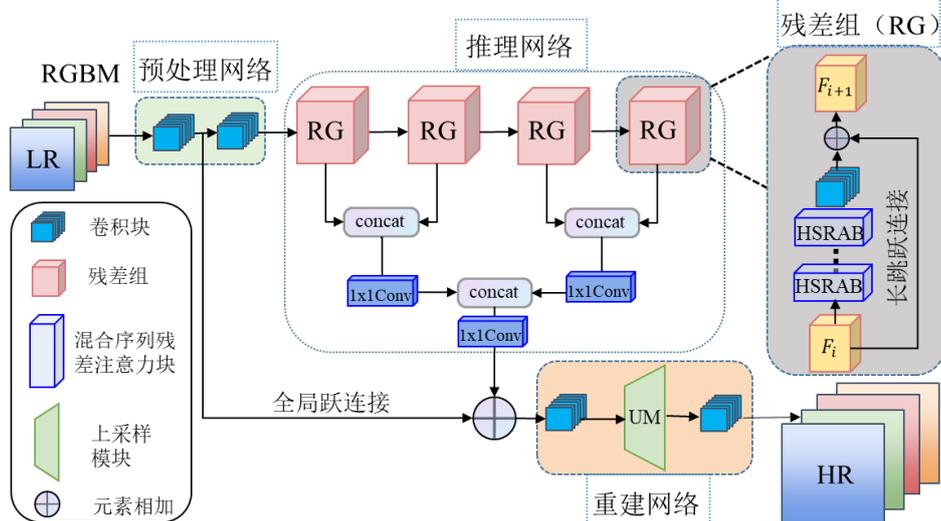


图1 本文算法网络结构

Fig.1 The Architecture of HSRATN

推理网络由多个残差组 (residual group, RG) 和特征融合结构组成。残差组包括多个混合序列残差注意力模块 (hybrid sequential residual attentional block, HSRAB) 和一个3×3卷积层, 提取特征图中的空间信息、通道信息和文本序列信息。特征融合结构使用1×1卷积融合相邻RG输出的特征图, 融合提取的低、中与高层次特征。网络中全局跳跃连接保存输入图像的初级特征, 使其只需学习残差映射关系以恢复丢失的高频细节, 不必学习完整图像之间复杂的映射关系。



图2 二值化语义掩膜

Fig. 2 The Demonstration of Binary Semantic Masks

重建网络将信息从特征空间映射到图像空间, 包括两个3×3卷积层和一个用像素洗牌层^[21]实现的上采样模块。

2.2 混合序列残差注意力模块

如图3所示, HSRAB 由混合残差注意力结构和序列单元 (sequential unit, SU) 组成。HSRAB 模块中加入短跳跃连接, 通过建立恒等映射通道提高输入特征的复用率, 解决训练过程中出现的梯度异常和网络性能退化问题^[22]。第*i*个RG中第*j*个HSRAB可以表示为:

$$F_{i,j+1} = H_{SUv}(H_{SUh}(H_{SA}(F_{ij}) \cdot H_{LCA}(F_{ij})) + F_{ij}) \quad (1)$$

式中, F_{ij} 、 $F_{i,j+1}$ 为输入输出特征, H_{SA} 、 H_{LCA} 表示空间 (spatial attention, SA) 和拉普拉斯通道注意力 (laplacian channel attention, LCA), H_{SUh} 和 H_{SUv} 表示处理图像水平行和竖直列的两个序列单元。

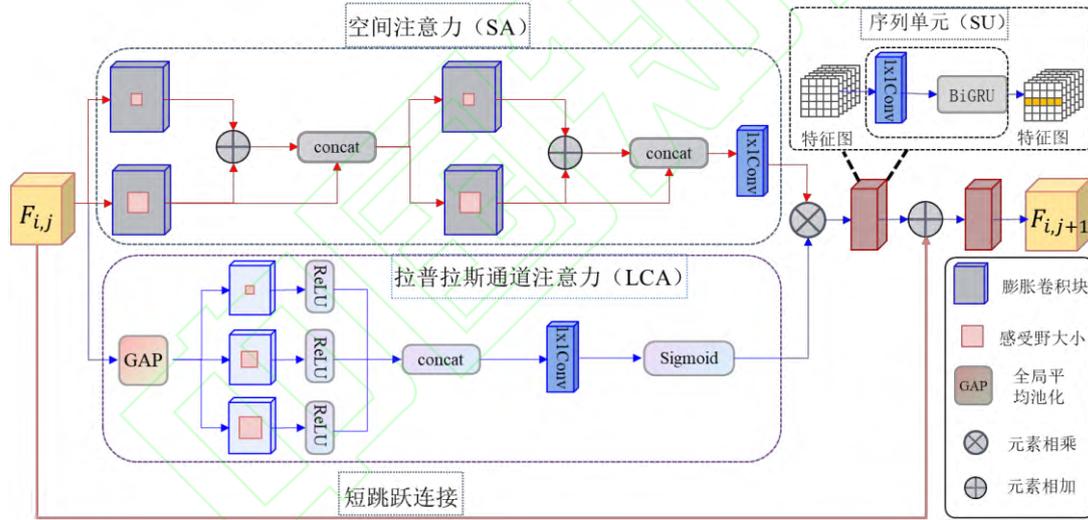


图3 混合序列残差注意力模块结构

Fig.3 The Proposed Hybrid Sequential Residual Attentional Block

2.2.1 混合残差注意力结构

混合残差注意力结构融合了空间注意力和拉普拉斯通道注意力机制。SA 通过提取不同感受野大小的特征图获得多层次空间特征, 引入的 LCA 通过学习多频率子带特征, 适应性地调整特征依赖关系来为特征通道赋权。融合二者使模型更好地关注图像的文本区域, 学习文本部分的多层次特征表示。

多尺度超分辨率重建网络 (multi-scale super-resolution network, MSRN)^[23]利用卷积核大小不同的卷积感受野大小不同的特性, 提取多尺度特

征图。HRAN 使用膨胀因子大小不同的膨胀卷积获得不同的感受野, 同时能减少网络参数数量。同样地, 本文算法使用多个不同的膨胀卷积提取文本多尺度特征, 将多尺度特征提取看作是空间注意力。SA 的整个过程可以表示为:

$$F_{S_1} = \tau(H_{D_1}(F_{ij})) \quad (2)$$

$$F_{S_2} = \tau(H_{D_2}(F_{ij}) + F_{S_1}) \quad (3)$$

$$F_S = [F_{S_1}; F_{S_2}] \quad (4)$$

$$F_{S_1} = \tau(H_{D_1}(F_S)) \quad (5)$$

$$F_{S_2} = \tau(H_{D_2}(F_S) + F_{S_1}) \quad (6)$$

$$H_{SA}(F_{ij}) = H_{1 \times 1}([F_{S_1}; F_{S_2}]) = \tau(W_D[F_{S_1}; F_{S_2}]) \quad (7)$$

式中, H_{D_1} 、 H_{D_2} 表示膨胀因子为 1 和 2 的膨胀卷积, τ 表示 LeakyReLU 激活函数, $H_{1 \times 1}$ 表示 1×1 卷积, W_D 为其权重。最后在通道维度上连结两个膨胀卷积层的输出, 然后使用 1×1 卷积进行通道压缩, 保持输出与输入特征图通道数相同。

特征图各个通道所含信息的重要程度不同, 通道注意力机制利用通道间相互关系来为各通道赋予权重。为了充分学习这种关系, 和 DRLN 一样, 本文在通道注意力结构内引入了拉普拉斯金字塔, 组成 LCA。首先, LCA 使用全局平均池化 (global average pooling, GAP) 操作获得图像的一维统计特征, 数学表达式如下:

$$F_C = H_{GAP}(F_{ij}) = \frac{1}{h \times w} \sum_{a=1}^h \sum_{b=1}^w F_{ij}(a,b) \quad (8)$$

式中, 输入特征图 F_{ij} 高、宽和通道数分别为 h 、 w 、 C , $F_{ij}(a,b)$ 表示 F_{ij} 在 (a,b) 位置的值。处理之后, 特征 F_C 的尺寸为 $1 \times 1 \times C$ 。然后, LCA 使用拉普拉斯金字塔学习多频率子带特征, 其中拉普拉斯金字塔由膨胀因子不同的膨胀卷积组成, 输出的多层次特征在通道维度上连结。然后经过一个 1×1 卷积进行上采样, 得到的一维特征经 sigmoid 激活函数处理后获得 LCA 权重。LCA 可以表示为:

$$F_{D_3} = \text{ReLU}(H_{D_3}(F_C)) \quad (9)$$

$$F_{D_5} = \text{ReLU}(H_{D_5}(F_C)) \quad (10)$$

$$F_{D_7} = \text{ReLU}(H_{D_7}(F_C)) \quad (11)$$

$$F_{Laplace} = [F_{D_3}; F_{D_5}; F_{D_7}] \quad (12)$$

$$H_{LCA}(F_{ij}) = f(W_U F_{Laplace}) \quad (13)$$

式中, H_{D_3} 、 H_{D_5} 、 H_{D_7} 表示膨胀因子为 3、5、7 的膨胀卷积, ReLU 为激活函数, F_{D_3} 、 F_{D_5} 、 F_{D_7} 、 $F_{Laplace}$ 分别表示多层次特征和拉普拉斯金字塔特征, W_U 、 f 表示上采样卷积层权重和 sigmoid 激活函数。最后, 将 SA 和 LCA 分支的输出相乘进行融合, 加权重要的子波段特征, 即提取最有效的文本区域特征。

2.2.2 序列单元

通用图像超分辨率算法一般只考虑重建图像中纹理、边缘等高频细节, 而街景影像文本中字符间有确切语义序列关系。即使某一字符模糊而无法直接辨识, 人也可以通过前后字符语义关系进行“完形填空”, 从而识别模糊字符。在文本检测识别领域, 许多模型使用循环神经网络及其变体提取文本语境信息^[24]。在文本超分辨率领域, TSRN 使用双向长短期记忆结构提取文本序列关系。为了简化模型,

提高计算效率, 本文提出采用结构更简单的双向门控循环单元 (bidirectional gating cycle unit, BiGRU) 提取文本字符间的序列关系。如图 3 右部所示, 提出的序列单元由卷积层和 BiGRU 组成。首先, 使用卷积层提取特征, 然后将特征图水平或者垂直列作为一维序列特征向量, 输入 BiGRU 更新隐藏层的内部状态, 学习序列特征之间的语义先验关系, 该过程可以表示为:

$$S_{t_h} = H_{\text{BiGRU}_h}(X_{t_h}, S_{t_h-1}), t_h=1,2,\dots,H \quad (14)$$

$$S_{t_v} = H_{\text{BiGRU}_v}(X_{t_v}, S_{t_v-1}), t_v=1,2,\dots,W \quad (15)$$

式中, S_t 表示隐藏状态, X_t 表示输入特征, H_{BiGRU_h} 、 H_{BiGRU_v} 分别表示处理水平文本行和垂直文本列的两个 BiGRU, t_h 、 t_v 表示沿着输入特征水平和垂直方向构造特征向量, H 、 W 为特征图宽和高。最后, SU 将 BiGRU 输出的一维序列特征转化成特征图。

2.3 梯度先验损失

现有基于逐像素求差的损失函数无法描述 HR 与 SR 图像在梯度轮廓上的差距, 超分辨率重建的文本边缘不够清晰。Sun 等人^[25]提出在 SR 网络中使用梯度轮廓先验重建出更锐利的边缘, 而 TSRN 将其应用在文本超分中。Tran 等人^[26]提出梯度差异损失, 将其与像素损失相结合, 增强重建图像的边缘。上述方法提出的梯度轮廓损失和梯度差异损失利用图像中的梯度先验知识, 指导网络梯度流动, 锐化文本边缘。因此, 引入梯度先验损失 (gradient prior loss, GPL) 重建文本边缘等细节。本文使用的 GPL 公式如下:

$$L_{GPL} = \frac{1}{N} \sum_{i=1}^N \|\nabla I_{sr}([I_{SR}, M_{SR}]^i) - \nabla I_{hr}([I_{HR}, M_{HR}]^i)\|_1 \quad (16)$$

式中, N 为批大小, ∇I_{sr} 和 ∇I_{hr} 表示图像梯度场。

本文算法提出的联合损失约束如图 4 所示, 通过逐像素计算 SR 与 HR 图像的差距, 在亮度层面对算法进行约束; 同时计算 SR 与 HR 图像梯度场 (像素 RGB 值的空间梯度) 之间的差距, 最后利用联合约束损失进行训练。HSRATN 的损失函数为:

$$L(\Theta) = \frac{1}{N} \sum_{i=1}^N (\| [I_{SR}, M_{SR}]^i - [I_{HR}, M_{HR}]^i \|_2 + \lambda_l \|\nabla I_{sr}([I_{SR}, M_{SR}]^i) - \nabla I_{hr}([I_{HR}, M_{HR}]^i)\|_1) \quad (17)$$

式中, Θ 表示网络参数, λ_l 为损失权重。

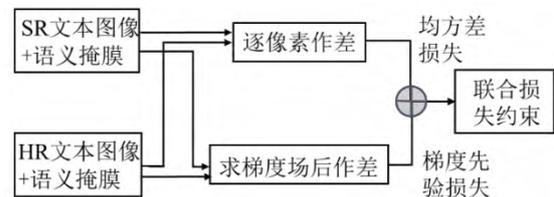


图 4 HSRATN 的联合损失约束

Fig.4 Joint Constraint Loss of HSRATN

2.4 实现细节

本文算法的网络配置与 HRAN 相同,包含 4 个 RG, 每个 RG 由 8 个 HSRAB 和一个 3×3 卷积层组成。模型能处理彩色和灰度图像,最后一层滤波器数目相应为 4 或 2, 其余各层滤波器数目均为 64。大部分卷积层采用 LeakyReLU 激活函数, 表达式为 $y = \max(0, x) + \theta * \min(0, x)$, θ 是很小的常数。

3 实验结果分析

3.1 测试数据集

本文使用真实场景文本超分辨率数据集 TextZoom 进行实验验证。数据集包含放大倍率为 2 的 LR (64×16) 和 HR(128×32)文本图像对, 以及对应的文本标签。训练集约 17000 张图像, 测试集包含 easy、medium、hard 三个部分, LR 图像质量依次降低。因为拍摄时相机的抖动、偏移以及裁剪操作, TextZoom 中文本图像更模糊, LR 和 HR 图像对之间有偏移问题。同 TSRN 一样, LR 图像经空间转换网络 (spatial transformer networks, STN) [27] 对齐处理后再输入本文网络, 以减轻偏移问题对重建结果造成的影响。

3.2 训练设置

选用 Adam 算法作为损失函数优化器, 梯度一阶和二阶矩估计的系数为 0.9 和 0.999。学习率为 0.0001, 训练周期为 200, 批量数据尺寸为 25, 联合约束损失中损失权重 λ 设置为 0.0001。采用 PyTorch 深度学习框架实现网络, 代码地址为 <https://github.com/Slupiter/HSRATN>, 硬件参数为: Intel Xeon(R) CPU E5-1620 3.5GHz, GeForce GTX 1080 GPU。

3.3 消融实验

为了分析 HSRATN 中 GPL、LCA、SU 的作用, 依次修改网络的配置, 比较重建效果的差异, 验证其有效性并构建最佳网络。使用文本识别模型 ASTER (attentional scene text recognizer) [24] 的识别准确率作为评价指标, 实验结果如表 1 所示。为了真实地反映各算法的文本重建结果在视觉感知效果上的差异, 如图 5 所示, 展示了部分结果。每个图像下方字符串为识别结果, 标红表示识别错误。

表1 不同配置HSRATN模型消融实验结果

Table.1 Ablation Study for Different Settings of HSRATN

方法	配置	损失函数	ASTER 准确率
----	----	------	-----------

			easy	medium	hard	平均
1	HRAN	L_2	70.5%	55.9%	38.5%	56.0%
2	HRAN	$L_2 + L_{GP}$	71.9%	55.9%	37.8%	56.3%
3	HRAN+LCA	L_2	71.3%	56.2%	38.9%	56.5%
4	HRAN+SU	L_2	72.8%	57.1%	39.6%	57.5%
5	HRAN+SU	$L_2 + L_{GP}$	72.9%	57.6%	39.8%	57.8%
6*	HRAN+SU+LCA	$L_2 + L_{GP}$	73.5%	58.4%	40.8%	58.6%

注: 加粗字体为每列最优值, *表示最佳网络。

本文分别基于像素损失与联合约束损失训练网络, 表 1 中方法 1 和 2 的结果表明, 加入 GPL 后, 平均识别准确率提高了 0.3%。虽然提升较少, 但从图 5 第四行的结果可以看出, 重建文本边缘清晰, 视觉质量更好。方法 1 和 3 的结果表明, 在通道注意力中加入拉普拉斯通金字塔结构后, 平均识别准确率提高了 0.5%。这说明混合注意力机制可以使模型更好地学习图像文本区域的多层次特征表示, 重建清晰的文本细节。从方法 1 和 4 的结果可以看出, 序列单元使得准确率提高了 1.5%。图 5 第五行的结果显示, 重建图像在视觉上更有辨识度, 如“supervisor”中的“s、e、o”。这说明在网络中加入 BiGRU 可以提取文本字符间的序列先验信息, 有利于提高模糊字符重建效果。

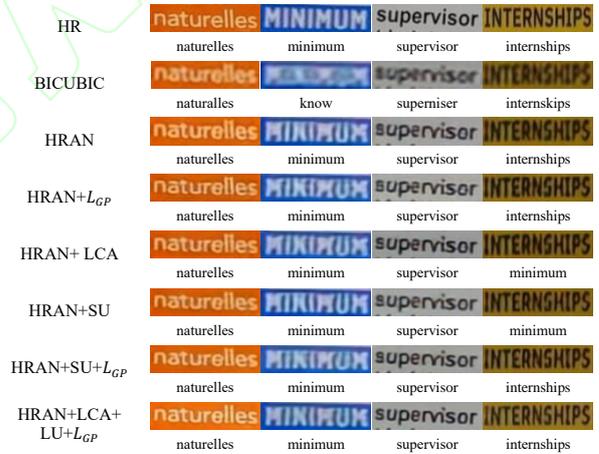


图 5 消融实验的结果

Fig.5 Results of Ablation Study

此外, 本文对六种方法的中间特征图进行可视化, 绘出彩色热力图, 如图 6 所示, 从蓝到红为热度增加方向。所选特征图为推理网络的输出。从 (b)、(c) 和 (e)、(f) 可以看出, GPL 增强了模型对文本边缘的关注度, 有利于获得清晰锐利的文本。此外, 从 (b)、(d)、(e) 和 (c)、(f)、(g) 可以看出, SU 和 LCA 使模型更加关注文本区域, 这有利于提取文本区域先验知识, 重建高质量文本。这进一步验证了本文算法各模块的有效性。



图 6 中间特征热力图

Fig. 6 The Intermediate Feature Heatmap

3.4 实验结果分析

将本文算法与双三次插值 (BICUBIC)、SRCNN、深度超分 (very deep SR, VDSR)、残差超分网络 (SR residual network, SRResNet)、拉普拉斯超分网络 (laplacian SR network, LapSRN)^[28]、增强深度残差超分 (enhanced deep residual SR, EDSR)、残差密集网络 (residual dense network, RDN)、MSRN、HRAN 和 DRLN 等通用超分方法以及文本超分

TSRN 进行对比实验。为了消除实验偏差, 使用 TextZoom 和公布的代码对所有模型重新训练。

表 2 列出了重建影像的文本识别准确率, 这里使用的主流文本识别模型 ASTER、MORAN (multi-object rectified attention network)^[29] 和 CRNN (convolutional recurrent neural network)^[30] 均为原作者公开的代码。可以看出, 分辨率放大 2 倍时, HSRATN 重建结果识别准确率优于其他模型。与基线算法 HRAN 相比, HSRATN 重建结果的 ASTER、MORAN 和 CRNN 平均识别准确率分别提高了 2.6%、3.2%和 4.7%。与该领域领先的通用超分算法 DRLN 相比, HSRATN 结果的识别准确率分别提高了 2%、2.4%和 3.3%。与文本超分算法 TSRN 相比, HSRATN 的 ASTER 和 CRNN 平均识别准确率在×2 尺度下分别获得 0.3%和 4.1%的提升, 达到了先进水平。

表2 TextZoom真实数据集超分辨率模型重建效果

Table.2 Performance of SR Models on Three Subsets in TextZoom

模型	损失函数	ASTER 准确率				MORAN 准确率				CRNN 准确率			
		easy	medium	hard	平均	easy	medium	hard	平均	easy	medium	hard	平均
BICUBIC	-	64.7%	42.4%	31.2%	47.2%	60.6%	37.9%	30.8%	44.1%	36.4%	21.1%	21.1%	26.8%
SRCNN	L_2	69.4%	43.3%	32.2%	49.5%	63.2%	39.0%	30.2%	45.3%	38.7%	21.6%	20.9%	27.7%
VDSR	L_2	71.7%	43.5%	34.0%	51.0%	62.3%	42.5%	30.5%	46.1%	41.2%	25.6%	23.3%	30.7%
SRResNet	$L_2 + L_{tv}$	69.6%	47.6%	34.3%	51.3%	60.7%	42.9%	32.6%	46.3%	39.7%	27.6%	22.7%	30.6%
RRDB	L_1	70.9%	44.4%	32.5%	50.6%	63.9%	41.0%	30.8%	46.3%	40.6%	22.1%	21.9%	28.9%
EDSR	L_1	72.3%	48.6%	34.3%	53.0%	63.6%	45.4%	32.2%	48.1%	42.7%	29.3%	24.1%	32.7%
RDN	L_1	70.0%	47.0%	34.0%	51.5%	61.7%	42.0%	31.6%	46.1%	41.6%	24.4%	23.5%	30.5%
LapSRN	Charbonnier	71.5%	48.6%	35.2%	53.0%	64.6%	44.9%	32.2%	48.3%	46.1%	27.9%	23.6%	33.3%
MSRN	L_2	70.2%	54.6%	37.0%	55.0%	64.2%	47.9%	35.1%	50.0%	49.80%	34.90%	29.9%	38.9%
HRAN	L_2	70.5%	55.9%	38.5%	56.0%	64.70%	48.80%	36.20%	50.8%	52.80%	36.80%	30.40%	40.8%
DRLN	L_2	72.3%	55.1%	39.1%	56.6%	66.58%	49.10%	36.30%	51.6%	53.40%	40.60%	30.50%	42.2%
TSRN	$L_2 + L_{GP}$	75.1%	56.3%	40.1%	58.3%	70.1%	53.3%	37.9%	54.8%	52.5%	38.2%	31.4%	41.4%
本文算法	$L_2 + L_{GP}$	73.50%	58.4%	40.8%	58.6%	67.2%	53.4%	38.8%	54.0%	56.2%	44.4%	33.7%	45.5%

	easy			medium			hard		
HR	statistics	Education	ACCESS	naturelles	MINIMUM	supervisor	NATIONAL	PARKING	MUSICA ALTA
BICUBIC	(stati)which	lack(ation)	more(s)	naturalles	know	supemiser	upon	pasking	mus(i)caalta
SRResNet	statistics	Education	ACCESS	naturelles	MINIMUM	supervisor	NATIONAL	PARKING	MUSICA ALTA
MSRN	statistics	Education	ACCESS	naturelles	MINIMUM	supervisor	NATIONAL	PARKING	MUSICA ALTA
HRAN	statistics	Education	ACCESS	naturelles	MINIMUM	supervisor	NATIONAL	PARKING	MUSICA ALTA
DRLN	statistics	Education	ACCESS	naturelles	MINIMUM	supervisor	NATIONAL	PARKING	MUSICA ALTA
TSRN	statistics	Education	ACCESS	naturelles	MINIMUM	supervisor	NATIONAL	PARKING	MUSICA ALTA
本文算法	statistics	Education	ACCESS	naturelles	MINIMUM	supervisor	NATIONAL	PARKING	MUSICA ALTA

图7 主流超分辨率模型结果

Fig.7 Results of State-of-the-Art SR Models

本文算法主要目的是提升在测试集上的文本识别准确率。表3也展示了各算法的图像重建质量评价指标峰值信噪比（peak signal-to-noise ratio, PSNR）和结构相似性指数（structural similarity index, SSIM）结果。可以看出，本文算法的PSNR在easy和medium测试集上数值较差，经分析有两方面原因：1）PSNR的定义和像素损失函数高度相关，最小化像素损失等同于直接最大化PSNR值。而本文采用联合约束损失训练网络，没有为了追求较高的PSNR值采用单一的像素损失。2）在处理真实数据集中LR/HR图像不对齐问题时，使用的STN可能造成图像像素偏移。SSIM在easy测试集上数值较差，推测本文算法更适用于较低质量图像重建。通常来说，PSNR和SSIM不能准确地反映图像视觉感知质量。因此，在本文中，文本识别准确率评价指标更为重要。

图7展示了BICUBIC、SRResNet、MSRN、HRAN、DRLN、TSRN和本文算法的部分重建图像以及Aster文本识别结果。可以看出，由于文本特征的结构化程度较高，对某些部位的形变起到限制作用，现有通用超分算法得到的结果过于平滑，文本边缘和纹理较为模糊，无法重建不同文本图像的细节特征，有些得到了错误的文本识别结果。HSRATN则受益于混合残差注意力结构提取的文本区域多层次特征表示、序列残差单元提取的文本序列信息和梯度先验损失，获得了丰富的文本先验知识，从而能够重建出更清晰的边缘和纹理细节，如“supervisor”中的“s”、“education”中的“a”和“national”中的“n”。此外，本文模型对倾斜文本也有效，如“access”和“musicalta”。

表3 TextZoom真实数据集SR模型PSNR和SSIM值

Table.3 PSNR and SSIM of SR Models on TextZoom

方法	PSNR			SSIM		
	easy	medium	hard	easy	medium	hard
BICUBIC	22.35	18.98	19.39	0.7884	0.6254	0.6592
SRCNN	23.48	19.06	19.34	0.8379	0.6323	0.6791
VDSR	24.62	18.96	19.79	0.8631	0.6166	0.6989
SRResNet	24.36	18.88	19.29	0.8681	0.6406	0.6911
RRDB	22.12	18.35	19.15	0.8351	0.6194	0.6856
EDSR	24.26	18.63	19.14	0.8633	0.6440	0.7108
RDN	22.27	18.95	19.70	0.8249	0.6427	0.7113
LapSRN	24.58	18.85	19.77	0.8556	0.6480	0.7087
MSRN	23.20	19.01	20.12	0.8538	0.6589	0.7301
HRAN	23.08	18.87	19.86	0.8621	0.6652	0.7355
DRLN	22.91	18.98	20.00	0.8622	0.6676	0.7425
TSRN	25.07	18.86	19.71	0.8897	0.6676	0.7302

本文算法	23.00	18.94	20.21	0.8710	0.6751	0.7542
------	-------	-------	--------------	--------	---------------	---------------

图8展示了本文算法和部分主流模型在性能和参数数量方面的实验结果，均在2倍真实TextZoom数据集上进行。可以看到，HSRATN比RDN、DRLN、EDSR参数数量更少，重建效果却更好。相比于EDSR模型，参数数量少了约72%，Aster平均识别准确率提高了5.6%。相比于基线模型HRAN，参数数量多了19%，准确率却提高了2.3%。这表明本文算法在性能和参数数量之间有良好的权衡。

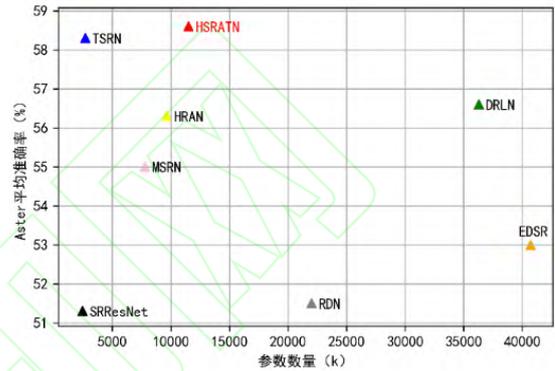


图8 性能和参数数量比较

Fig.8 Comparisons for Performance and Number of Parameters

3.5 合成TextZoom数据集实验

为了进一步验证本文算法的泛化性能，使用双三次插值下采样合成的TextZoom数据对部分算法进行训练和测试。表4展示了HSRATN与BICUBIC、通用超分算法SRResNet、MSRN、HRAN以及文本超分算法TSRN的ASTER文本识别准确率。可以看出，分辨率放大4倍时，HSRATN重建结果识别准确率优于其他模型。相比于基线模型HRAN，平均识别准确率提高了0.9%。这表明本文算法在其他超分倍率及退化模型上有较强的泛化能力。

表4 超分模型在合成TextZoom数据集上的结果

Table.4 Results of SR Models on Synthetic TextZoom

方法	放大倍率	ASTER 准确率			
		测试1	测试2	测试3	平均
BICUBIC	×4	21.7%	31.5%	45.9%	32.3%
SRResNet	×4	43.1%	54.9%	58.4%	51.6%
MSRN	×4	45.6%	57.8%	61.1%	54.3%
HRAN	×4	51.3%	59.8%	62.3%	57.4%
TSRN	×4	51.8%	60.9%	60.4%	57.4%
本文算法	×4	52.8%	60.7%	62.4%	58.3%

3.6 ICDAR 2015 TextSR数据集实验

受文献[19]启发，使用ICDAR 2015 TextSR数据集^[4]进行实验，验证本文方法重建效果，该数据

集源自法语视频字幕。此处采用数据集规定的 Tesseract-OCR v3.02 软件测试文本识别准确率 OCR_{ac} ，可以表示为：

$$OCR_{ac} = 1 - \frac{1}{K} \sum_{i=1}^N (d_i) \quad (18)$$

其中， N 为测试集图像数量 141， K 为图像中字符总数 2929， d_i 表示图像文本标签和识别结果字符串的编辑距离。

表 5 展示了部分通用 SR 方法、文本超分辨率方法^[19, 31-32]和本文方法 2 倍放大实验的客观评价指标结果。可以看到，本文方法在该数据集上文本识别率高于其它方法， OCR_{ac} 为 79.04，比 HR 原图识别准确率 78.8% 高 0.24%，比 HRAN、文献[19]和文献[32]分别提高 0.41%、0.24%、0.21%。图 9 展示了本文方法超分重建的部分结果，可以看出重建图像文本边缘清晰、纹理细节丰富，接近 HR 图像。这表明本文算法在其它场景文本图像数据上有较强的泛化能力。

表5 ICDAR 2015 TextSR数据集上的结果

Table.5 Results on ICDAR 2015 TextSR Dataset

方法	放大倍率	PSNR	SSIM	OCR 准确率
SRCNN	×2	31.75	0.980	77.19%
SRResNet	×2	29.04	0.950	76.55%
LapSRN	×2	31.40	0.972	76.48%
EDSR	×2	32.09	0.976	77.53%
RDN	×2	32.61	0.978	77.67%
HRAN	×2	32.92	0.987	78.63%
文献[19]	×2	32.86	0.979	78.80%
文献[31]	×2	33.21	0.978	78.78%
文献[32]	×2	33.94	0.982	78.83%
本文算法	×2	32.97	0.987	79.04%



图 9 ICDAR 2015 TextSR 数据集本文方法结果 (×2)

Fig.9 Results of HSRATN on ICDAR 2015 TextSR Dataset (×2)

4 结 语

针对街景影像分辨率低导致文本识别准确率降低的问题，提出了一种融合混合残差注意力机制与序列单元的影像文本超分辨率重建算法。融合空间注意力和拉普拉斯通道注意力，充分利用特征图在

空间和通道上各自的依赖关系，使网络专注于提取文本区域的多层次特征，恢复高频信息。利用序列单元提取图像中文字符间的序列先验信息，对算法性能进行优化，提高文本模糊字符重建效果。采用结合梯度先验损失与内容损失的联合损失，重建清晰文本边缘。实验结果表明本文算法可以更好地利用街景影像中文字区域的先验知识，超分辨率重建影像文本边缘清晰、纹理细节丰富，提高了文本识别准确率。

参 考 文 献

- [1] Wang W J, Xie E Z, Sun P Z, et al. TextSR: Content-Aware Text Super-Resolution Guided by Recognition, 2019[OL]. <https://arxiv.org/pdf/1909.07113.pdf>, 2022
- [2] Wang W J, Xie E Z, Liu X B, et al. Scene Text Image Super-Resolution in the Wild[M]//Computer Vision - ECCV 2020. Cham: Springer International Publishing, 2020: 650-666
- [3] Dong C, Zhu X M, Deng Y B, et al. Boosting Optical Character Recognition: A Super-Resolution Approach, 2015[OL]. <https://arxiv.org/pdf/1506.02211.pdf>, 2022
- [4] Peyrard C, Baccouche M, Mamalet F, et al. ICDAR2015 competition on text image super-resolution[C]//2015 13th International Conference on Document Analysis and Recognition (ICDAR). Tunis, Tunisia.: 1201-1205
- [5] Dong C, Loy C C, He K M, et al. Image Super-Resolution Using Deep Convolutional Networks[J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2016, 38(2): 295-307
- [6] Pandey R K, Vignesh K, Ramakrishnan A G, et al. Binary Document Image Super Resolution for Improved Readability and OCR Performance, 2018[OL]. <https://arxiv.org/pdf/1812.02475.pdf>, 2022
- [7] Nakao R, Iwana B K, Uchida S. Selective super-resolution for scene text images[C]//2019 International Conference on Document Analysis and Recognition (ICDAR). Sydney, NSW, Australia.: 401-406
- [8] Lin K, Liu Y B, Li T H, et al. Text image super-resolution by image matting and text label supervision[C]//2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW). Long Beach, CA, USA.: 1722-1727
- [9] Wang Z H, Chen J, Hoi S C H. Deep Learning for Image Super-Resolution: A Survey[J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021, 43(10): 3365-3387
- [10] Liao Haibin, Chen Youbin, Chen Qinghu. Non-Local Similarity Dictionary Learning Based Super-Resolution for Improved Face Recognition[J]. *Geomatics and Information Science of Wuhan University*, 2016, 41(10): 1414-1420 [10] (廖海斌, 陈友斌, 陈庆虎. 基于非局部相似字典学习的人脸超分辨率与识别[J]. 武汉大学学报·信息科学版, 2016, 41(10): 1414-1420)
- [11] Chen Hang, Luo Bin. Multi-Angle Remote Sensing Images Super-Resolution Reconstruction Using Dynamic Upsampling Filter Deep Network[J]. *Geomatics and Information Science of Wuhan University*, 2021, 46(11): 1716-1726 (陈行, 罗斌. 利用

- 动态上采样滤波深度网络进行多角度遥感影像超分辨率重建[J]. 武汉大学学报·信息科学版, 2021, 46(11): 1716-1726
- [12] Lim B, Son S, Kim H, et al. Enhanced deep residual networks for single image super-resolution[C]//2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops. Honolulu, HI, USA.: 1132-1140
- [13] Vaswani A, Shazeer N, Parmar N, et al. Attention Is All You Need, 2017[OL]. <https://arxiv.org/pdf/1706.03762.pdf>, 2022
- [14] Fu J, Liu J, Tian H J, et al. Dual attention network for scene segmentation[C]//2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Long Beach, CA, USA.: 3141-3149
- [15] Dosovitskiy A, Beyer L, Kolesnikov A, et al. An Image Is Worth 16x16 Words: Transformers for Image Recognition at Scale, 2021[OL]. <https://arxiv.org/pdf/2010.11929.pdf>, 2022
- [16] Zhao H S, Jia J Y, Koltun V. Exploring self-attention for image recognition[C]//2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Seattle, WA, USA.: 10073-10082
- [17] Zhang Y L, Li K P, Li K, et al. Image Super-Resolution Using Very Deep Residual Channel Attention Networks[C]//Proceedings of the 2018 European Conference on Computer Vision, Munich, Germany, 2018
- [18] Muqet A, Iqbal M T B, Bae S H. HRAN: Hybrid Residual Attention Network for Single Image Super-Resolution[J]. *IEEE Access*, 7: 137020-137029
- [19] Wang Y Y, Su F, Qian Y. Text-attentional conditional generative adversarial network for super-resolution of text images[C]//2019 IEEE International Conference on Multimedia and Expo. Shanghai, China.: 1024-1029
- [20] Anwar S, Barnes N. Densely Residual Laplacian Super-Resolution[J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022, 44(3): 1192-1204
- [21] Shi W Z, Caballero J, Huszár F, et al. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network[C]//2016 IEEE Conference on Computer Vision and Pattern Recognition. Las Vegas, NV, USA.: 1874-1883
- [22] He K M, Zhang X Y, Ren S Q, et al. Deep Residual Learning for Image Recognition, 2015[OL]. <https://arxiv.org/pdf/1512.03385.pdf>, 2022
- [23] Li J C, Xie E Z, Fang F M. Multi-Scale Residual Network for Image Super-Resolution[C]//Proceedings of the 2018 European Conference on Computer Vision, Munich, Germany, 2018
- [24] Shi B G, Yang M K, Wang X G, et al. ASTER: An Attentional Scene Text Recognizer with Flexible Rectification[J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019, 41(9): 2035-2048
- [25] Sun J, Sun J, Xu Z B, et al. Gradient Profile Prior and Its Applications in Image Super-Resolution and Enhancement[J]. *IEEE Transactions on Image Processing*, 2011, 20(6): 1529-1542
- [26] Tran H T M, Phuoc T H. Deep Laplacian Pyramid Network for Text Images Super-Resolution[J]//Proceedings of the 2019 IEEE-RIVF International Conference on Computing and Communication Technologies, Danang, Vietnam, 2019
- [27] Jaderberg M, Simonyan K, Zisserman A, et al. Spatial Transformer Networks[C]//Proceedings of the 29th Annual Conference on Neural Information Processing Systems, Montreal, Canada, 2015
- [28] Lai W S, Huang J B, Ahuja N, et al. Deep Laplacian Pyramid Networks for Fast and Accurate Super-Resolution[C]//Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition, Honolulu, HI, USA, 2017
- [29] Luo C J, Jin L W, Sun Z H. MORAN: A Multi-Object Rectified Attention Network for Scene Text Recognition[J]. *Pattern Recognition*, 2019, 90: 109-118
- [30] Shi B G, Bai X, Yao C. An End-to-End Trainable Neural Network for Image-Based Sequence Recognition and Its Application to Scene Text Recognition[J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017, 39(11): 2298-2304
- [31] Geng C, Chen L, Zhang X, et al. Adversarial Text Image Super-Resolution using Sinkhorn Distance[C]//ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2020.
- [32] Xue M, Huang Z, Liu R, et al. A Novel Attention Enhanced Residual-In-Residual Dense Network for Text Image Super-Resolution[C]//2021 IEEE International Conference on Multimedia and Expo (ICME), 2021.

Text Super-resolution Method with Attentional Mechanism and Sequential Units

WEI Haodong¹ YI Yaohua¹ YU Changhui¹ LIN Liyu²

¹ School of Remote Sensing and Information Engineering, Wuhan University, Wuhan 430079, China

² State Key Laboratory of Information Engineering in Surveying, Mapping and Remote Sensing, Wuhan University, Wuhan 430079, China

Abstract: Objectives: The text in street view images is the clue to perceive and understand scene information. Low-resolution street view images lack details in the text region, leading to poor recognition accuracy. Super-resolution can be introduced as pre-processing to reconstruct edge and texture details of the text region. To improve text recognition accuracy, we propose a text super-resolution

network combining attentional mechanism and sequential units. **Methods:** A hybrid residual attention structure is proposed to extract spatial information and channel information of the image text region, learning multi-level feature representation. A sequential unit is proposed to extract sequential prior information between texts in the image through bidirectional gated recurrent units. Using gradient prior knowledge as the constraint, a gradient prior loss is designed to sharpen character boundaries. **Results and Conclusions:** In order to verify the effectiveness of the proposed method, we use real scene text images in TextZoom and synthetic text images to carry out comparative analysis experiments. Experimental results show that the proposed method can reconstruct clear text edges and rich text texture details, and improve text recognition accuracy of street view images.

Key words: street view images; super-resolution; attentional mechanism; sequential information; gradient prior loss

First author: WEI Haodong, master, specializes in the theories and methods of super-resolution. E - mail: WHUweihaodong@163.com

Corresponding author: YI Yaohua, PhD, professor. E - mail: yyh@whu.edu.cn

Foundation support: The National Key Research and Development Program of China (2021YFB2206200).

网络首发:

标题: 融合注意力与序列单元的文本超分辨率

作者: 韦豪东, 易尧华, 余长慧, 林立宇

DOI: 10.13203/j.whugis20220158

收稿日期: 2022-07-14

引用格式:

韦豪东, 易尧华, 余长慧, 等. 融合注意力与序列单元的文本超分辨率[J]. 武汉大学学报·信息科学版, 2022, DOI: 10.13203/j.whugis20220158 (WEI Haodong, YI Yaohua, YU Changhui, et al. Text Super-resolution Method with Attentional Mechanism and Sequential Units[J]. *Geomatics and Information Science of Wuhan University*, 2022, DOI: 10.13203/j.whugis20220158)

网络首发文章内容和格式与正式出版会有细微差别, 请以正式出版文件为准!

您感兴趣的其他相关论文:

基于非局部相似字典学习的人脸超分辨率与识别

廖海斌, 陈友斌, 陈庆虎

武汉大学学报·信息科学版, 2016, 41(10): 1414-1420

<http://ch.whu.edu.cn/cn/article/doi/10.13203/j.whugis20140498>