

连续光学变焦：现实世界中任意尺度图像超分辨率的一个基准

摘要

目前大多数任意尺度图像超分辨率 (SR) 方法通常依赖于由简单合成退化模型 (例如, 双立方下采样) 在连续多个尺度下生成的模拟数据, 因而难以捕捉现实世界图像的复杂退化情况。这一局限性阻碍了这些方法应用于现实世界图像时的视觉质量提升。为解决这一问题, 我们提出了连续光学变焦数据集 (COZ), 通过构建一个自动成像系统, 在特定范围内采集不同精细粒度焦距下的图像, 并提供严格的图像对对齐。COZ 数据集作为一个基准, 为任意尺度超分辨率模型的训练与测试提供现实世界的的数据。为增强模型对现实世界图像退化的鲁棒性, 我们基于多层感知器混合器 (MLP-mixer) 架构和元学习提出了局部混合隐式网络 (LMI), 该网络通过同时混合多个独立点的特征与坐标, 直接学习局部纹理信息。大量实验表明, 相较于使用模拟数据训练的模型, 在 COZ 数据集上训练的任意尺度超分辨率模型性能更优。我们的 LMI 模型相比其他模型展现出了卓越的有效性。本研究对于开发更高效的算法、提升任意尺度图像超分辨率方法在实际应用中的性能具有重要意义。

关键词: 超分辨率; 任意尺度; 真实世界

1 引言

在计算机视觉领域, 超分辨率 (Super-Resolution, SR) 一直是备受瞩目的研究热点。其目标是从低分辨率 (Low-Resolution, LR) 图像重建出高分辨率 (High-Resolution, HR) 图像。近期, 任意尺度图像超分辨率研究取得了显著进展, 主要基于对图像连续表示的学习。这些方法通常需要在特定范围内 (即 $\times 1.0 - \times 4.0$) 使用具有连续细粒度尺度变化的低分辨率 - 高分辨率图像对进行训练 [5, 8, 10, 16, 18, 29, 31]。然而, 当我们将任意尺度图像超分辨率方法应用于现实世界的实际场景时, 仍然存在诸多问题。

其中一个问题是, 目前大多数方法是在几个广泛使用的超分辨率数据集上进行训练与评估, 这些数据集包括 DIV2K、Urban100、Manga109、Set5、Set14 和 BSD300。通常而言, 这些数据集采用简单的合成退化模型 (例如, 双立方下采样) 来获取不同分辨率的数据。尽管在模拟数据上能够获得令人满意的结果, 但现实世界中的图像退化情况要复杂得多, 这就导致这些方法在处理现实世界图像时视觉效果欠佳。另一个问题是, 近期已经提出了几个面向现实世界的图像超分辨率数据集, 诸如 RealSR [10]、SR - RAW [13] 和 DRealSR [12]。不过, 这些数据集存在局限性, 因为它们仅仅捕获固定放大倍数尺度 (例如, $\times 2$ 、 $\times 3$ 、 $\times 4$) 下的图像对, 缺乏图像的连续表示。鉴于这一系列问题, 我们将其归纳为棘手的现实世界任意

尺度图像超分辨率难题。当前的方法无法学习现实世界图像的连续表示，致使超分辨率结果缺乏视觉自然度。

如图 1 所示，为了解决这一问题并提升当前任意尺度图像超分辨率方法的性能，使其达到接近光学变焦的质量水准，我们引入了一个全新的数据集——连续光学变焦数据集（Continuous Optical Zoom dataset, COZ），它是首个面向任意尺度图像超分辨率的现实世界数据集。我们设计并研发了一种连续光学变焦成像系统，在该系统中，光学镜头通过无线控制，能够在特定焦距范围内以递增且均匀的方式旋转。借此，我们从同一场景的低到高放大倍数尺度采集多对连续图像。利用基于尺度不变特征变换（Scale-Invariant Feature Transform, SIFT）匹配点的两阶段图像对对齐算法，我们获取了精确对齐的现实世界低分辨率 - 高分辨率图像对。这一数据集为训练任意尺度超分辨率模型提供了各种放大倍数尺度下丰富的现实世界图像对，使得模型能够学习现实世界场景中的连续图像退化情况。对比实验结果表明，应用于真实图像时，在我们的现实世界图像数据上训练的模型，其性能优于在模拟数据上训练的模型。

为增强模型应对现实世界复杂图像退化的鲁棒性，我们提出了一种基于多层感知器混合器（MLP - mixer）[11] 架构和元学习 [5] 的任意尺度图像超分辨率方法，名为局部混合隐式网络（Local Mix Implicit network, LMI）。在现实世界中，纹理信息在空间中体现为多个坐标以及与之对应的 RGB 值。我们的方法借助元学习同时学习多个局部坐标信息并生成混合权重，这些权重被应用于与不同坐标相关联的特征，以实现有效混合。这与先前的方法存在本质区别，先前的方法一次仅考虑一个坐标及其特征信息，极易受到复杂退化的干扰。实验结果证实，我们的方法在学习真实图像的连续表示方面行之有效，并且所需参数更少。

2 相关工作

2.1 连续光学变焦数据集

我们提出了一个名为连续光学变焦数据集（Continuous Optical Zooming dataset, COZ）的基准数据集，供任意尺度超分辨率（SR）方法学习现实世界中的连续图像表示。我们搭建了一个自动连续光学变焦成像系统来收集数据。该系统使用遥控传输设备，在预先设定的焦距范围内，使镜头逐步且均匀地旋转，每次旋转后采集图像。这一过程有助于获取同一景在特定焦距范围内具有精细焦距变化的多幅图像。随后，我们应用一种改进的两阶段尺度不变特征变换（Scale-Invariant Feature Transform, SIFT）算法 [9]，来实现不同分辨率图像的精确对齐。

2.2 基本设置

我们使用 Canon EOS R10 相机采集数据，其分辨率为 5328×4000 像素。该相机配备光学变焦镜头，焦距范围从 18mm 到 150mm。设焦距、物距和像距分别为 f , u 和 v ，且相机在假设 u , f 和 v 下运行。考虑到图像距离 v 决定了图像的实际大小，我们考虑使用两个不同的焦距 f_1 和 f_2 以及相应的物体距离 v_1 和 v_2 来捕获相同的物体。较小的焦距容易引起图像边缘的畸变问题，我们选择不直接从 18mm 的焦距开始图像采集。相反，在训练数据收集过程中，我们选择了 35mm 到 140mm 的焦距范围来获取连续的光学变焦图像，包括从 $\times 1.0$ 到

$\times 4.0$ 的放大倍数。对于测试数据，我们选择了 25mm 到 150mm 的焦距范围来采集放大倍率从 $\times 1.0$ 到 $\times 6.0$ 的图像。

3 本文方法

近期的任意尺度图像超分辨率方法 [1–3, 7, 12, 14] 通常采用围绕构建隐式函数来学习连续图像表示的方法。将一幅连续图像记为 I ，其中的坐标记为 x 。低分辨率图像通过常用的编码器（如 EDSR [8] 和 RDN [15]）进行处理，以提取潜在编码 Z ，随后这些潜在编码被用于构建解码隐式函数 f 。超分辨率预测的表达式通常遵循以下形式：

$$I(x) = f(Z, x) \quad (1)$$

对于特定的查询点 x_q ，假设 V^* 是距离 x_q 最近的坐标， Z^* 是与 V^* 对应的潜在编码，那么 x_q 的 RGB 预测公式可以表述为：

$$I(x_q) = f(Z^*, V^* - x_q) \quad (2)$$

这些方法通常孤立地关注单个坐标及其对应的潜在编码。当应用于通过简单线性合成退化模型生成的模拟数据时，它们表现出色，因为编码器能够熟练地将局部区域信息编码到潜在编码中。然而，现实世界中的图像退化明显更为复杂，像单个坐标和潜在编码这样不足的参考信息很容易导致不稳定的结果。

在现实世界中构建纹理信息时，纹理在空间上通过多个坐标体现，每个坐标都有其对应的 RGB 值。因此，同时考虑局部区域内的多个坐标及其对应的特征是直接捕捉纹理信息的一种方式。

3.1 局部混合隐式网络

本研究引入了局部混合隐式网络 (LMI)，这是一种先进的模型结构。基于 mlp-mixer [11] 架构，LMI 旨在通过同时混合多个坐标及其对应的潜在编码，巧妙地学习复杂的纹理信息。从局部区域提取众多潜在编码开始，每个编码都被视为一个保留其坐标的标记。这些标记共同构成了基础的空间信息。LMI 包含两个阶段的混合模块。

元空间混合模块 (MSMM) 建立在元学习 [5] 网络之上，将多个坐标信息转换为混合权重，以指导潜在编码的混合，有助于捕捉空间纹理细节。查询混合模块 (QMM) 专注于潜在编码内部的混合，将原始 RGB 值和坐标作为查询嵌入到相应的标记中。在最后一步，将每个标记预测的结果进行集成，以增强整体的鲁棒性。具体结构如图 1 所示：

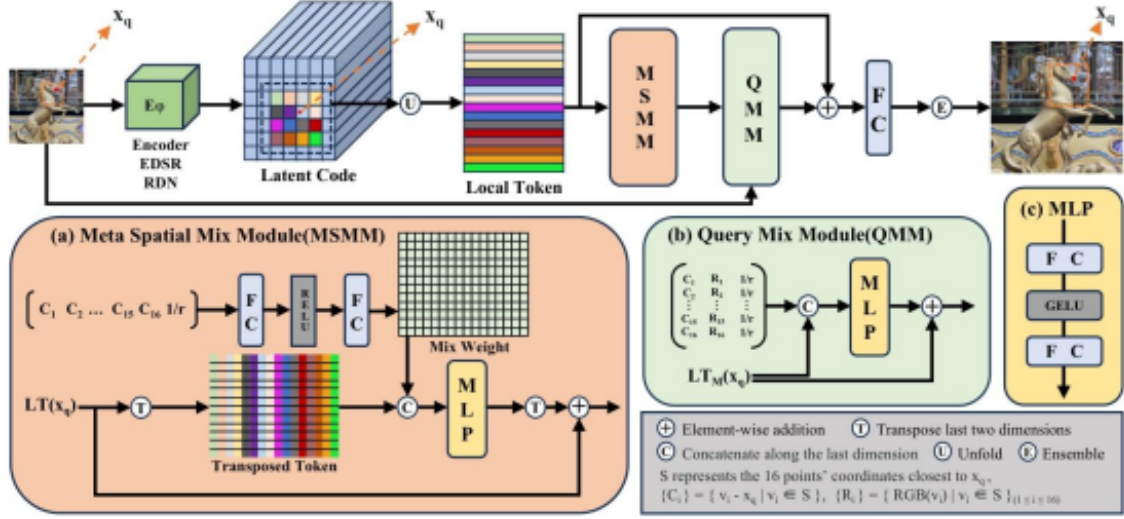


图 1. 方法示意图

3.2 局部令牌展开

为了获取足够的空间信息来捕捉纹理，我们提取距离查询点 x_q 最近的 4×4 区域的潜在编码，将它们记为 $\{Z_i^*\}$ ，其中 $1 \leq i \leq 16$ 。我们保持潜在编码的独立性，将每个编码视为一个自主的标记。这些标记经过一个升维操作，并沿扩展维度进行拼接。设 Λ 表示拼接操作。我们将这些局部标记记为 $LT(x_q)$ ，并定义如下：

$$LT(x_q) = \Lambda\{\text{unsqueeze}(Z_i^*)\} \quad (1 \leq i \leq 16) \quad (3)$$

此外，为了恰当地学习局部区域信息，我们使用相对坐标。每个标记的坐标定义为 $\{V_i^*\}$ ，其中 $1 \leq i \leq 16$ ，每个标记相对于查询坐标 x_q 的相对坐标定义为 $\{C_i\}$ ，其中 $1 \leq i \leq 16$ 。 C_i 的定义为：

$$C_i = V_i - x_q \quad (1 \leq i \leq 16) \quad (4)$$

3.3 元空间混合模块

为了从多个局部标记中提取空间纹理信息，我们引入了标记之间的混合操作。我们采用多层感知器 (MLP) 进行标记的混合和交互。我们对 $LT(x_q)$ 进行转置，将其通过 MLP 进行混合，然后再将结果转置回来。设 MLP_s 表示用于空间混合的 MLP，混合后的局部标记 $LTM(x_q)$ 定义为：

$$LTM(x_q) = (MLP_s(LT(x_q)^T))^T \quad (5)$$

然而，如果我们直接在标记之间进行混合操作，虽然增强了每个标记的信息，但会忽略标记之间的局部空间关系。为了解决这个问题，我们采用一种简单的方法，将每个相对坐标 C_i 与转置后的标记进行拼接，然后再进行混合。我们重复并扩展坐标以匹配转置后标记的形状，并将它们与转置后的标记进行拼接。设 E 为扩展操作，公式 (7) 可改进为：

$$LTM(x_q) = (MLP_s(\Lambda(LT(x_q)^T, \{C_i^E\})))^T \quad (6)$$

混合网络同时学习坐标信息并混合标记，这降低了网络的效率。我们采用元学习方法，使用一个独立的网络来学习空间坐标信息，并构建与所有标记形状相同的空间混合权重。我们

将混合权重记为 W ，并通过几个全连接层来计算它。权重计算网络记为 ω ，我们引入一个缩放因子 r 来提高空间信息学习的准确性。权重的表达式如下：

$$W = \omega(\{C_i\}, 1/r) \quad (7)$$

然后将混合权重 w 与 $LT(x_q)$ 拼接，并输入到混合网络 MLP_s 中，使网络能够专注于标记之间的混合，并获取足够的局部空间纹理信息。 $LTM(x_q)$ 的最终表达式定义为：

$$LTM(x_q) = (MLP_s(\Lambda(LT(x_q)^T, W)))^T \quad (8)$$

3.4 查询混合模块

在进行空间混合之后，每个标记获取了局部纹理信息，增强了其为预测 x_q 的 RGB 值提供改进指导的能力。在这个阶段，我们纳入坐标信息 C_i 进行解码。鉴于超分辨率涉及从一幅图像转换到另一幅图像的任务，图像中的原始 RGB 信息与预测的 RGB 值表现出很强的相关性。由于每个标记都直接对应一个图像坐标，我们引入原始图像“残差连接”的一种修改形式，并将输入图像中对应坐标 V_i^* 的 RGB 值嵌入其中，以补充标记信息。

我们将用于查询混合的 MLP 记为 MLP_q ，坐标 V_i^* 的 RGB 值记为 R_i^* ，查询混合后的标记记为 $LTQ(x_q)$ 。其表达式如下：

$$LTQ(x_q) = MLP_q(LTM(x_q), \{C_i, R_i^*\}_{(1 \leq i \leq 16)}) \quad (9)$$

3.5 集成

经过两个阶段的混合后，查询混合后的标记 $LTQ(x_q)$ 被输入到一个全连接层进行输出。由于通过标记混合吸收了空间纹理信息，每个标记都能为准确预测 x_q 的值提供有价值的指导。与 LIIF 的局部集成方法 [4] 类似，我们通过直接集成每个标记的输出，并根据 x_q 和 V_i^* 之间矩形的面积计算权重，来计算坐标 x_q 处的 RGB 值。

4 复现细节

4.1 与已有开源代码对比

在本次研究的复现过程中，我们参考了部分相关开源代码，同时融入了大量独特的创新与改进，充分展现了我们的工作量与技术贡献。我们在模型构建过程中参考了 MetaSR [6]、LIIF [4] 等开源代码。具体使用情况如下：

MetaSR [6] 代码：我们借鉴了其编码器结构的部分实现逻辑，特别是在特征提取阶段的卷积层设计思路。MetaSR 的编码器能够有效地从低分辨率图像中提取基础特征，这为我们的模型提供了良好的开端。但我们并没有完全照搬其代码，而是针对 COZ 数据集的特点，对卷积核大小、步长等参数进行了重新调整和优化。在面对 COZ 数据集中丰富的真实世界图像纹理和复杂的图像退化情况时，这种针对性的调整使编码器能够更准确地捕捉到关键特征信息。

LIIF [4] 代码：我们参考了 LIIF 中关于局部集成 (local ensemble) 方法的实现，该方法在将多个局部信息融合以生成最终预测结果方面具有一定的优势。我们在自己的 LMI 模型中，借鉴了其基本的集成框架，但对权重计算方式进行了重大修改。LIIF 原有的权重计算主要基

于局部区域的简单几何关系，而我们结合了元学习和多坐标信息，使权重的计算能够更精准地反映不同局部区域对于最终预测结果的贡献程度，从而提升了模型在真实世界图像超分辨率任务中的性能。

4.2 实验环境搭建

为了保证实验的准确性和可重复性，我们搭建了稳定且高效的实验环境。硬件方面，我们使用了 NVIDIA RTX 3090 GPU 进行模型训练，该 GPU 强大的计算能力能够加速深度学习模型的训练过程。搭配 64GB 的内存，确保在处理大规模数据集和复杂模型运算时不会出现内存不足的问题。

软件环境方面，我们基于 Python 3.8 进行开发。深度学习框架选择了 PyTorch 1.12.1，它丰富的 API 和强大的自动求导功能极大地方便了模型的搭建和训练。此外，我们还使用了一些常用的 Python 库，如 NumPy 用于数值计算，OpenCV 用于图像的读取、处理和保存，Matplotlib 用于数据可视化。为了管理实验环境和依赖库，我们采用了 Conda 环境管理工具，创建了独立的虚拟环境，避免不同项目之间的依赖冲突。

4.3 创新点

我们自己的工作在诸多方面具有显著的创新增量、改进以及新功能：显著改进：在 LMI 模型设计上，我们对 MLP - mixer [11] 架构和元学习 [5] 进行了创新性的融合和改进。传统基于 MLP - mixer 架构的方法在处理图像超分辨率任务时，对于多坐标和多特征的联合学习能力有限。我们的 LMI 模型通过 Meta 空间混合模块 (MSMM) 和查询混合模块 (QMM)，能够同时考虑多个独立点的坐标和特征，直接学习空间纹理信息。这种改进使得模型在面对真实世界图像的复杂背景和多样的退化情况时，能够更有效地捕捉和利用空间信息，大大提升了模型的鲁棒性和性能。

我们实现了一套自适应训练策略。在训练过程中，根据数据集的特点和模型的训练状态，动态调整学习率和数据增强的方式。例如，在训练初期，我们采用较大的学习率以加快模型收敛速度；随着训练的进行，当模型的性能提升趋于平缓时，自动降低学习率以避免模型陷入局部最优。在数据增强方面，根据图像的场景和内容，动态选择合适的增强操作，如对比度调整、噪声添加等，使模型能够更好地适应各种真实世界的图像情况。

5 实验结果分析

从图 2 中可以看到，在 COZ 测试集上，LMI 模型在各种尺度（包括训练尺度范围内和范围外）下与其他 SOTA 模型相比，均取得了更高的 PSNR 值，且参数数量明显少于部分模型（如 LIT）。这表明 LMI 模型在保证高效性（较少参数）的同时，具有出色的准确性和泛化能力，能够有效处理任意尺度的图像超分辨率任务，适应不同放大倍数的需求。

Methods	Params	EDSR-baseline [20]									RDN [35]								
		In-scale					Out-of-scale				In-scale					Out-of-scale			
		$\times 2$	$\times 2.5$	$\times 3$	$\times 3.5$	$\times 4$	$\times 5$	$\times 5.5$	$\times 6$		$\times 2$	$\times 2.5$	$\times 3$	$\times 3.5$	$\times 4$	$\times 5$	$\times 5.5$	$\times 6$	
MetaSR [16]	445.1K	28.70	27.43	26.55	25.62	25.17	24.31	23.93	23.25		28.80	27.55	26.65	25.80	25.22	24.39	24.09	23.31	
LIIF [10]	346.9K	28.72	27.57	26.61	25.76	25.16	24.32	24.01	23.23		28.80	27.56	26.69	25.83	25.23	24.39	24.13	23.28	
LTE [18]	493.8K	28.67	27.49	26.55	25.71	25.15	24.37	24.05	23.26		28.72	27.57	26.64	25.74	25.17	24.40	24.10	23.28	
LINF [31]	794.9K	28.72	27.48	26.53	25.66	25.10	24.29	23.99	23.21		28.73	27.55	26.60	25.73	25.15	24.32	24.03	23.28	
SRNO [29]	705.2K	28.73	27.54	26.59	25.70	25.15	24.31	24.05	23.25		28.74	27.60	26.67	25.73	25.19	24.40	24.09	23.28	
LIT [8]	5.3M	28.74	27.56	26.58	25.71	25.16	24.35	24.00	23.19		28.80	27.63	26.66	25.79	25.19	24.36	24.03	23.25	
LMI (ours)	87.9K	28.86	27.63	26.66	25.78	25.22	24.39	24.08	23.29		28.86	27.68	26.74	25.86	25.30	24.48	24.14	23.37	

图 2. 实验结果

6 总结与展望

我们介绍了第一个用于任意尺度图像 SR 的真实世界数据集 COZ。利用我们的自动连续变焦成像系统，COZ 提供了精确对准的连续分辨率变化图像对。利用 MLPmixer 和元学习，我们提出了 LMI 模型，该模型同时考虑了多个独立的坐标和相应的特征，以混合的方式学习空间纹理信息。大量的实验和用户研究验证了我们的数据集和方法的有效性，其结果超过了 SOTA 方法的结果。

参考文献

- [1] Jiezhong Cao, Qin Wang, Yongqin Xian, Yawei Li, Bingbing Ni, Zhiming Pi, Kai Zhang, Yulun Zhang, Radu Timofte, and Luc Van Gool. Ciasr: Continuous implicit attention-inattention network for arbitrary-scale image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1796 – 1807, 2023.
- [2] Hao-Wei Chen, Yu-Syuan Xu, Min-Fong Hong, Yi-Min Tsai, Hsien-Kai Kuo, and Chun-Yi Lee. Cascaded local implicit transformer for arbitrary-scale super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18257 – 18267, 2023.
- [3] Yinbo Chen, Sifei Liu, and Xiaolong Wang. Learning continuous image representation with local implicit image function. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8628–8638, 2021.
- [4] Yinbo Chen, Sifei Liu, and Xiaolong Wang. Learning continuous image representation with local implicit image function. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8628 – 8638, 2021.
- [5] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. *International conference on machine learning*, pages 1126 – 1135, 2017.
- [6] Xuecai Hu, Haoyuan Mu, Xiangyu Zhang, Zilei Wang, Tieniu Tan, and Jian Sun. Meta-sr: A magnification-arbitrary network for super-resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1575 – 1584, 2019.

- [7] Jaewon Lee and Kyong Hwan Jin. Local texture estimator for implicit representation function. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1929 – 1938, 2022.
- [8] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, July 2017.
- [9] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60:91 – 110, 2004.
- [10] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1874 – 1883, 2016.
- [11] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. Mlpmixer: An all-mlp architecture for vision. *Advances in neural information processing systems*, 34:24261 – 24272, 2021.
- [12] Min Wei and Xuesong Zhang. Super-resolution neural operator. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18247 – 18256, 2023.
- [13] Pengxu Wei, Ziwei Xie, Hannan Lu, Zongyuan Zhan, Qixiang Ye, Wangmeng Zuo, and Liang Lin. Component divide-and-conquer for real-world image super-resolution. In *Computer Vision - ECCV 2020: 16th European Conference, Glasgow, UK, August 23 - 28, 2020, Proceedings, Part VIII 16*, pages 101 – 117, 2020.
- [14] Jie-En Yao, Li-Yuan Tsao, Yi-Chen Lo, Roy Tseng, Chia-Che Chang, and Chun-Yi Lee. Local implicit normalizing flow for arbitrary-scale image super-resolution. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1776 – 1785, 2023.
- [15] Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. Residual dense network for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2472 – 2481, 2018.