

Pixel to Gaussian: Ultra-Fast Continuous Super-Resolution with 2D Gaussian Modeling

Long Peng^{1,3†} Anran Wu^{1,2†} Wenbo Li^{3*} Peizhe Xia¹ Xueyuan Dai⁴ Xinjie Zhang⁵

Xin Di¹ Haoze Sun⁶ Renjing Pei³ Yang Wang^{1,4*} Yang Cao¹ Zheng-Jun Zha¹

¹USTC ²AHU ³Huawei Noah's Ark Lab ⁴Chang'an University ⁵HKUST ⁶THU

{longp2001@mail., ywang120@ustc.edu.cn, liwenbo50@huawei.com}

<https://github.com/peylnog/ContinuousSR>

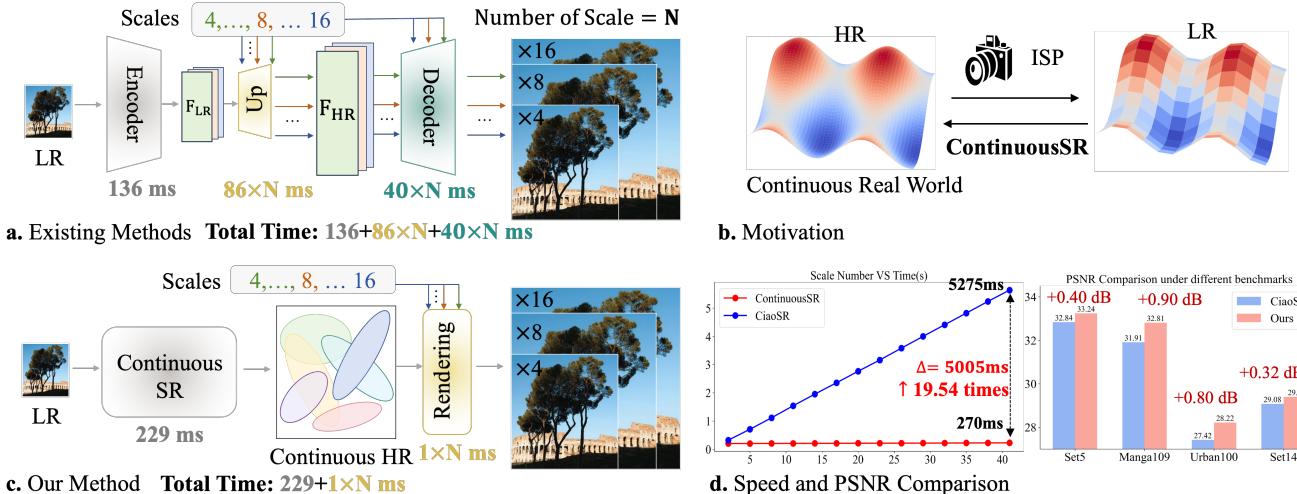


Figure 1. (a) Leveraging implicit modeling, existing ASSR methods rely on multiple upsampling and decoding steps to reconstruct HR images at different scales, which leads to low efficiency and performance. (b-d) Our method explicitly reconstructs 2D continuous HR signals from LR images in a single pass. Then, fast rendering replaces the time-consuming upsampling and decoding process to reconstruct HR images at different scales, significantly improving both performance (0.90 dB in Manga109) and efficiency (19.5× speedup).

Abstract

Arbitrary-scale super-resolution (ASSR) aims to reconstruct high-resolution (HR) images from low-resolution (LR) inputs with arbitrary upsampling factors using a single model, addressing the limitations of traditional SR methods constrained to fixed-scale factors (e.g., $\times 2$). Recent advances leveraging implicit neural representation (INR) have achieved great progress by modeling coordinate-to-pixel mappings. However, the efficiency of these methods may suffer from repeated upsampling and decoding, while

their reconstruction fidelity and quality are constrained by the intrinsic representational limitations of coordinate-based functions. To address these challenges, we propose a novel ContinuousSR framework with a Pixel-to-Gaussian paradigm, which explicitly reconstructs 2D continuous HR signals from LR images using Gaussian Splatting. This approach eliminates the need for time-consuming upsampling and decoding, enabling extremely fast arbitrary-scale super-resolution. Once the Gaussian field is built in a single pass, ContinuousSR can perform arbitrary-scale rendering in just 1ms per scale. Our method introduces several key innovations. Through statistical analysis, we uncover the Deep Gaussian Prior (DGP) and propose DGP-Driven Covariance Weighting, which dynamically optimizes

*Corresponding Authors: Wenbo Li liwenbo50@huawei.com; Yang Wang, ywang120@ustc.edu.cn. †These authors contributed equally to this work.

covariance via adaptive weighting. Additionally, we present Adaptive Position Drifting, which refines the positional distribution of the Gaussian space based on image content, further enhancing reconstruction quality. Extensive experiments on seven benchmarks demonstrate that our ContinuousSR delivers significant improvements in SR quality across all scales, with an impressive 19.5 \times speedup when continuously upsampling an image across forty scales.

1. Introduction

Cameras and smartphones discretize continuous real-world scenes into discrete 2D digital images [4, 27, 56], as illustrated in Figure 1(b). However, limitations in sensor resolution, among other factors, often lead to low-resolution (LR) images that fail to meet user requirements. Image super-resolution (SR) has been proposed to enhance image resolution and finer details [40, 54]. Unlike traditional fixed-scale super-resolution [10, 37, 38, 44, 65], which uses multiple models to learn mappings for fixed scales (*e.g.*, $\times 2$, $\times 3$, $\times 4$), arbitrary-scale super-resolution (ASSR) employs a single model to handle super-resolution with arbitrary scales, which has attracted significant attention [8, 24, 33, 36, 41, 50, 55].

Among these approaches, implicit neural representation (INR) has emerged as a leading technique, delivering visually compelling results [3, 11, 34, 56, 57]. INR aims to learn a continuous mapping from pixel coordinates to pixel values, enabling arbitrary-scale super-resolution through multiple upsampling and decoding steps, as illustrated in Figure 1(a). For example, LIIF [11] is the first to introduce INR into ASSR, employing multi-layer perceptrons to learn this mapping. Later, CiaoSR [3] and CLIT [8] leverage Transformers to enhance the modeling of long-range dependencies in feature upsampling and decoding, achieving state-of-the-art performance. However, the reconstruction fidelity and quality of these methods are inherently constrained by the representational limitations of coordinate-based implicit functions, making it challenging to effectively model continuous high-resolution signals, ultimately leading to sub-optimal performance. Additionally, their reliance on repeated upsampling and decoding significantly reduces efficiency, making real-world deployment impractical.

Given that LR images are discretized from continuous 2D signals, we pose the fundamental question: “*Can we directly reconstruct continuous HR signals from LR images and flexibly choose the desired scale?*” As illustrated in Figure 1(b), this approach not only enhances signal continuity through continuous modeling—leading to improved reconstruction quality—but also significantly boosts efficiency by eliminating the need for time-consuming upsampling and decoding. This idea enables fast and flexible ASSR, making real-world applications more practical.

In this paper, we introduce the novel ContinuousSR framework, built upon the Pixel-to-Gaussian paradigm, which reconstructs 2D continuous HR signals through Gaussian modeling. By first reconstructing a continuous HR Gaussian field, our method enables rapid sampling directly from the continuous representation, effectively replacing traditional time-consuming upsampling and decoding steps. This innovative approach achieves high-quality ASSR in just 1 ms, significantly improving the efficiency.

Directly applying Gaussian modeling to simulate real-world images is highly challenging due to the intricate interweaving of pixel distributions and parameters. To address this, we first identify the Deep Gaussian Prior (DGP) from 40,000 natural images, revealing that the distribution of Gaussian field parameters follows a Gaussian pattern with regularities in their range, as illustrated in Figure 2(a-b). Leveraging this insight, we sample pre-defined Gaussian kernels from the DGP distribution and introduce a novel DGP-Driven Covariance Weighting module, which efficiently optimizes covariance parameters through adaptive weighting. This helps guide the model toward the global optimum. Furthermore, we propose a Adaptive Position Drifting module, which dynamically adjusts the spatial positions of Gaussian kernels based on image content, enhancing structural accuracy. With these innovations, our method not only surpasses state-of-the-art approaches by up to 0.9 dB in reconstruction performance but also achieves a 19.5 \times speedup when continuously upsampling across forty scales. Our main contributions are as follows:

- A novel ContinuousSR is proposed to reconstruct continuous HR signals from LR images by 2D Gaussian modeling, thereby enabling fast and high-quality super-resolution with arbitrary scale.
- The Deep Gaussian Prior (DGP) is discovered, based on which DGP-Driven Covariance Weighting is proposed to facilitate the optimization of covariance. Furthermore, Adaptive Position Drifting is introduced to dynamically learn spatial positions in Gaussian space.
- Extensive experiments demonstrate that our method achieves state-of-the-art performance on seven benchmarks and ultra-fast speed.

2. Related work

2.1. Arbitrary-Scale Super-Resolution

Although traditional fixed-scale super-resolution (FSSR) methods, which use separate models to learn different super-resolution scales, have achieved significant progress [5, 13, 15, 30, 32, 45, 51, 63, 65], they struggle to meet the demand for arbitrary-scale super-resolution in real-world scenarios. Additionally, maintaining multiple models incurs high computational costs, making them less practical. To address these limitations, Arbitrary-Scale

Super-Resolution (ASSR) has been proposed to achieve it with a single model, gaining increasing attention in recent years [3, 17, 19, 19–21, 24, 28, 47, 49, 62, 66, 67]. For example, MetaSR [24] was the first to introduce the meta-upscale module to achieve arbitrary-scale super-resolution, demonstrating promising results. Inspired by the success of implicit neural representation (INR) in 3D reconstruction, LIIF [11] was the first to adapt INR to super-resolution by using a multilayer perceptron to learn the mapping from image coordinates and features to RGB values. To capture more high-frequency details, LTE [34] encodes textures in the Fourier space, while SRNO [55] leverages neural operators to model global relationships. CLIT [8] introduces a cross-scale interaction mechanism to enhance feature learning by integrating information across different resolutions. CiaoSR [3] further improves long-range modeling capability by introducing transformers to INR, achieving state-of-the-art performance. LMF [21] enhances local texture details by combining multi-frequency information in a computationally efficient manner, significantly reducing computational costs while maintaining the reconstruction of fine-grained features. However, these implicit modeling methods struggle to explicitly reconstruct continuous HR signals and require time-consuming upsampling and decoding, leading to low performance and efficiency.

2.2. Gaussian Splatting

Gaussian Splatting (GS) is introduced into 3D as a faster, more efficient alternative to NeRF, using anisotropic 3D Gaussians for real-time rendering and direct scene manipulation [29]. Building on 3DGS, 2D Gaussian Splatting improves the geometric accuracy of radiance fields by combining 2D Gaussians with precise scene projections [25]. Recently, 2D GS finds applications in image processing [6, 16]. For instance, Zhang *et al.* propose leveraging Gaussian Splatting (GS) for image compression and reconstruction [64] through long-time optimization of GS parameters, while Hu *et al.* employ Gaussian Splatting in the feature space to enhance visual quality and speed [23]. However, these methods still struggle to reconstruct continuous HR signals and suffer from long optimization times or multiple upsampling and decoding process.

3. Motivation

To capture the real world, advanced imaging sensors (*e.g.*, CMOS) are used to project the 3D continuous world into 2D and then discretized 2D continuous signals into 2D discrete signals [4, 27, 56], as formulated:

$$I[m, n] = f_c(m\Delta x, n\Delta y). \quad (1)$$

where $f_c(x, y)$ represents the continuous intensity function in the spatial domain (x, y) . The Δx and Δy denote the

sampling step along the spatial dimensions, while $m, n \in \mathbb{Z}$ are the corresponding discrete pixel grids. $I[m, n]$ represent the discrete images. After that, the Image Signal Processor is used to quantize, process, and encode it into a digital low-resolution image \mathbf{I}_{LR} .

Although many methods leveraging implicit modeling have been proposed [11, 33] to achieve ASSR by constructing coordinate-to-pixel mappings, two major challenges remain. On the one hand, the aim of ASSR is to reconstruct $f_c(x, y)$. However, implicit modeling makes it difficult to explicitly model high-quality continuous functions, resulting in limited performance. On the other hand, the pipeline of INR-based ASSR methods suffers from low efficiency, as follows:

$$\begin{aligned} \mathcal{F}_{LR} &= \mathbb{E}(\mathbf{I}_{LR}), & \mathcal{F}_{HR}^s &= \mathbb{U}(\mathcal{F}_{LR}, s), \\ \mathbf{I}_{HR}^s &= \mathbb{D}(\mathcal{F}_{HR}^s). \end{aligned} \quad (2)$$

where \mathbb{E} , \mathbb{U} and \mathbb{D} represent the Encoder, Upsampling, Decoder, respectively, and \mathcal{F}_{HR}^s denotes the high-resolution feature map at scale s . It can be observed that for different scales s , this method requires multiple time-consuming upsampling \mathbb{U} and decoding \mathbb{D} processes to reconstruct HR images F_{HR}^s , as shown in Figure 1(b), resulting in inefficiency. Therefore, we propose the fundamental question: “*Can we directly reconstruct continuous HR signals from LR images?*” This serves as the inverse function of imaging process Eq. 1, as illustrated in Figure 1(c). This approach would not only perform simple sampling to replace multiple upsampling and decoding but also enhance continuity, improving efficiency and performance.

4. Proposed Method

4.1. Continuous Basis Function

Considering that the target function is continuous, it is crucial to select an appropriate continuous basis function. In this work, we choose the Gaussian function for two main reasons: a) Leveraging the Gaussian Mixture Model (GMM) [46], any complex continuous function can be represented as a combination of several Gaussian functions, ensuring broad applicability and theoretical soundness. b) With the recent advancements in the Gaussian splatting community [7, 18], the engineering efficiency and compatibility of Gaussian functions have significantly improved, making them highly suitable for practical implementation. Therefore, we use Gaussian functions $G_i(x, y)$ as fundamental continuous functions to reconstruct real 2D continuous signals $f_c(x, y)$, as shown in the following equation:

$$f_c(x, y) = \sum_{i=1}^N G_i(x, y) \quad (3)$$

where N denotes the number of Gaussian kernels, x and y represents the location in the 2D space. Each Gaussian

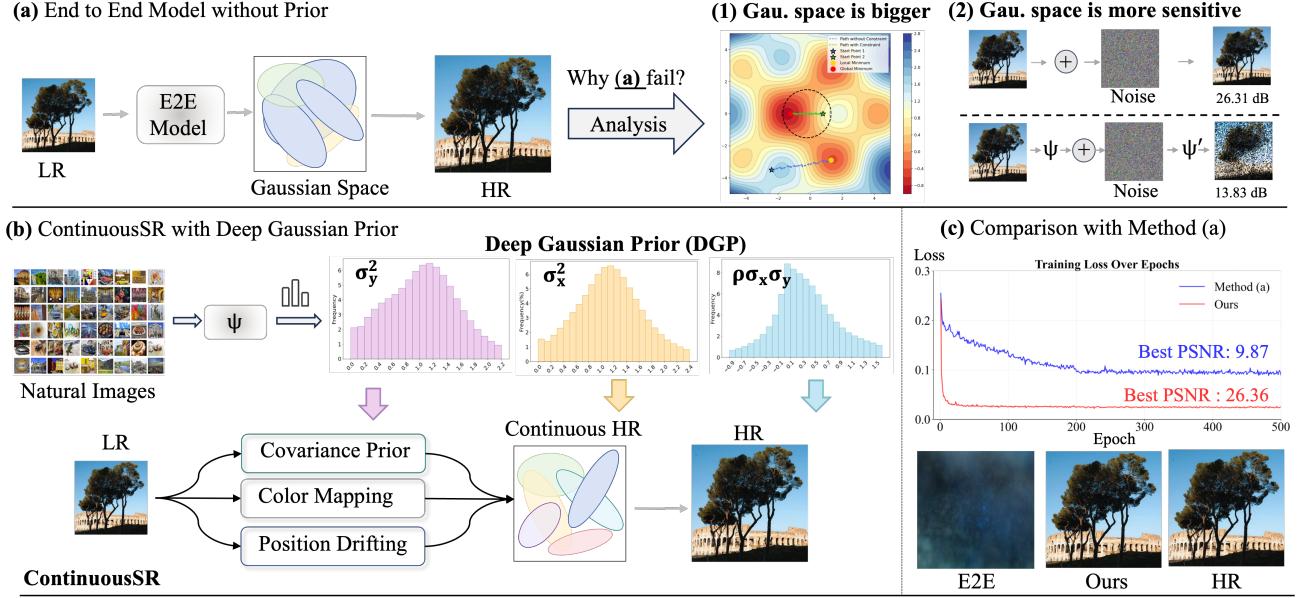


Figure 2. (a) Directly learning the end-to-end model from LR to the Gaussian field is challenging due to the vastness and sensitivity of the Gaussian space. (b-c) Through statistical analysis of 40,000 natural images, we uncover the Deep Gaussian Prior and propose Position Drifting, Covariance Prior, and Color Mapping to propose a novel ContinuousSR, enhancing the quality of the Gaussian field.

kernel has eight parameters needed to optimized, which include:

$$\Sigma = \begin{bmatrix} \sigma_x^2 & \rho\sigma_x\sigma_y \\ \rho\sigma_x\sigma_y & \sigma_y^2 \end{bmatrix}, \mu = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, c_{rgb} = \begin{bmatrix} c_r \\ c_g \\ c_b \end{bmatrix}, \quad (4)$$

where c_{rgb} denotes the RGB parameters of each Gaussian, μ represents the position parameters, and Σ represents the covariance matrix, resulting in a total of eight parameters to be optimized. The value of the Gaussian kernel G_i at the position (x, y) can be expressed as:

$$G_i(x, y, c_{rgb}, \Sigma) = c_{rgb} \frac{1}{2\pi|\Sigma|} \exp\left(-\frac{1}{2}d^\top \Sigma_i^{-1} d\right). \quad (5)$$

where the distance vector d represents the deviations of x and y from their positions μ_x and μ_y .

4.2. Direct End-to-End and Deep Gaussian Prior

A straightforward approach is to learn the parameters of Gaussian kernels directly from low-resolution (LR) images through an end-to-end model. However, this approach is extremely difficult to optimize, as shown in Figure 2(a). As shown in Figure 2(c), the blue loss curve indicates that the optimization process falls into a local optimum, with the PSNR remaining as low as 10 dB. To rule out the possibility of coincidence, we conduct multiple experiments and consistently observe the same conclusion.

Why does direct end-to-end fail? We attribute this to two main challenges: **a) High Complexity:** Each kernel in the

Gaussian space contains numerous difficult-to-learn parameters that need to be optimized, such as position, covariance, and RGB values. Many of these parameters have solution spaces ranging from 0 to positive infinity, resulting in an exceptionally large solution space. For instance, the covariance matrix theoretically only needs to satisfy the condition of being a positive definite matrix. This makes the Gaussian Space significantly larger than traditional image space, while introducing more local traps. Consequently, the complexity of optimization in Gaussian Space increases, making it more prone to local optima, as illustrated in Figure 2 (a). **b) High Sensitivity:** In Gaussian Space, even a slight change in any parameter of a single Gaussian, such as position or covariance, can significantly affect the entire image. This is fundamentally different from the image space, where a single pixel only impacts itself. To further verify this, we add noise with the same distribution to both image space and Gaussian Space to evaluate sensitivity. Note that the Gaussian Space is derived through the optimization method [64], denoted ψ , which requires approximately 1 minute of GPU time per scene to ensure high-quality mapping. The comparison results, as shown in Figure 2(c), indicate that the PSNR in image space is 26.31 dB, whereas it is only 13.83 dB in Gaussian Space. This demonstrates that Gaussian Space is much more sensitive, making the optimization more challenging.

Observation and Deep Gaussian Prior. To uncover the secrets of the Gaussian Space, we conduct statistical experiments to analyze the distribution of Gaussian param-

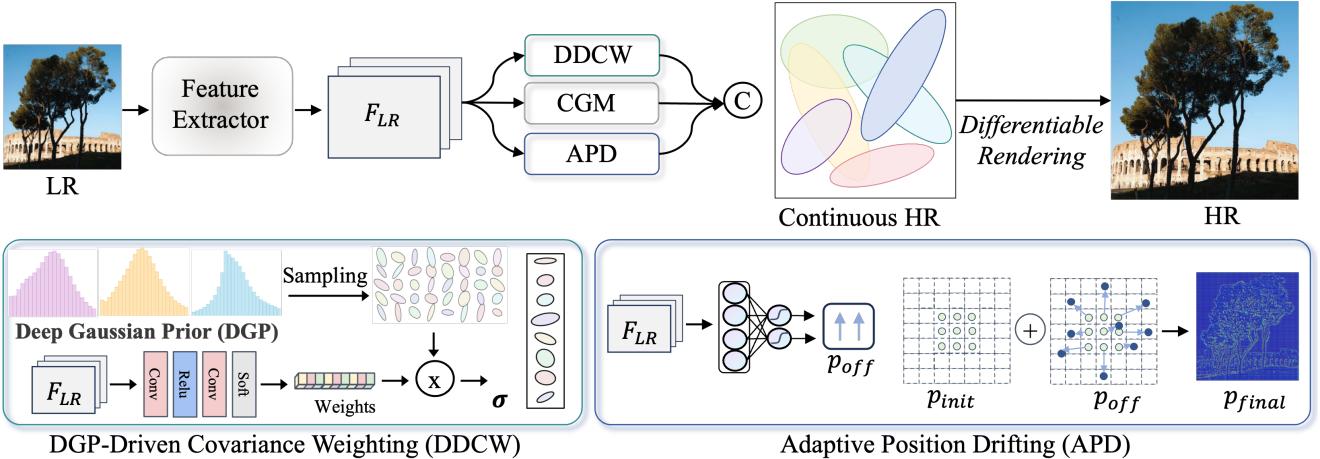


Figure 3. An overview of the proposed ContinuousSR framework, which consists of three key innovations: DGP-Driven Covariance Weighting (DDCW), Adaptive Position Drifting (APD), and Color Gaussian Mapping (CGM).

ters. Specifically, we collect and crop approximately 40,000 high-resolution images [39, 48], and transform them into the Gaussian space using ψ , with optimized over 700 GPU-hours. Subsequently, we statistically analyze the key parameters of Gaussian kernels, including σ_x^2 , σ_y^2 , and $\rho\sigma_x\sigma_y$. The results, as shown in Figure 2(b), indicate that the distribution of most covariances is traceable: a) Approximately 99% of σ_x^2 , σ_y^2 , and $\rho\sigma_x\sigma_y$ fall within the ranges of $0 \sim 2.4$, $0 \sim 2.2$, and $-0.9 \sim 1.5$, respectively. b) The distributions of the three covariances generally follow a Gaussian distribution. We define this finding as the **Deep Gaussian Prior (DGP)**, which provides valuable information to reduce the difficulty of optimization. Based on these observations, we propose an innovative method, ContinuousSR, which for the first time achieves representation learning from low-resolution (LR) images to continuous HR signals. Specifically, ContinuousSR introduces DGP-Driven Covariance Weighting, which simplifies the optimization difficulty in Gaussian Space by constructing pre-defined Gaussian kernels, employing an adaptive weighting mechanism, and incorporating Adaptive Position Drifting based on offset drifting. This approach enables superior performance and achieves fast super-resolution results, as shown in Figure 1 and 3.

4.3. DGP-Driven Covariance Weighting

Directly learning Gaussian covariance parameters remains challenging due to their unknown range and sensitive space. To address this issue, we propose a novel DGP-Driven Covariance Weighting, which leverages the deep Gaussian prior (DGP) to construct a set of pre-defined Gaussian kernels. This approach simplifies the complex task of directly learning covariance parameters into learning a set of weighting coefficients to combine the pre-defined kernels and represent the target kernel, as shown in Figure 3.

Specifically, using the DGP, we sample the three covari-

ance parameters σ_x^2 , σ_y^2 , and ρ from the corresponding distributions of the DGP. These parameters are then used to construct a dictionary of N pre-defined Gaussian kernels. The sampling process is expressed as:

$$\sigma_{i,x}^2, \sigma_{i,y}^2 \sim \mathcal{P}(\sigma_x^2), \mathcal{P}(\sigma_y^2); \rho_i \sigma_{i,x} \sigma_{i,y} \sim \mathcal{P}(\rho \sigma_x \sigma_y), \quad (6)$$

to construct pre-defined Gaussian kernels \mathcal{K} :

$$\mathcal{K} = \{G_i(\begin{bmatrix} \sigma_{i,x}^2 & \rho_i \sigma_{i,x} \sigma_{i,y} \\ \rho_i \sigma_{i,x} \sigma_{i,y} & \sigma_{i,y}^2 \end{bmatrix})\}_{i=1}^N. \quad (7)$$

These candidate covariance kernels cover the majority of types and ranges commonly observed in natural images, providing valuable prior information to facilitate network convergence. We then extract features \mathcal{F}_{LR} from the input LR image \mathbf{I}_{LR} using the backbone encoder \mathbb{E} :

$$\mathcal{F}_{LR} = \mathbb{E}(\mathbf{I}_{LR}). \quad (8)$$

To adaptively generate the target covariance kernel, we introduce an adaptive weighting mechanism that learns a set of weights $\mathbf{W} = \{w_i\}_{i=1}^N$ based on the extracted features. These weights are computed by the adaptive weighting module $\mathcal{M}_{\text{weight}}$, which operates as follows:

$$\mathbf{W} = \text{Softmax}(\mathcal{M}_{\text{weight}}(\mathcal{F}_{LR})). \quad (9)$$

The adaptive weighting module $\mathcal{M}_{\text{weight}}$ is implemented using several layers of convolutional neural networks (CNNs). Finally, each target kernel is generated by performing a weighted combination of the pre-defined kernels in the dictionary:

$$G_{\text{target}} = \sum_{i=1}^N w_i \cdot G_i. \quad (10)$$

Through the proposed method, we achieve effective optimization of Gaussian covariance, providing stronger prior

Table 1. PSNR performance comparison with state-of-the-art methods under different benchmarks. Average Time (AT) is reported in milliseconds (ms). The best and the second-best results are in **bold** and bold. More comparisons are in Appendix Section 10.

PSNR↑	Methods	×4	×6	×8	×10	×12	×16	×18	×20	×32	×48	AT
Urban100 [26]	MetaSR [24]	26.76	24.31	22.92	22.02	21.31	20.35	19.96	19.65	18.38	17.48	41.4
	LIIF [11]	26.68	24.20	22.79	21.84	21.15	20.19	19.80	19.51	18.30	17.45	110.0
	LTE [34]	27.24	24.62	23.17	22.23	21.50	20.47	20.06	19.77	<u>18.47</u>	<u>17.52</u>	151.8
	SRNO [55]	26.98	24.43	23.02	22.06	21.36	20.35	19.95	19.67	18.39	17.51	65.7
	CiaoSR [3]	<u>27.42</u>	<u>24.84</u>	<u>23.34</u>	<u>22.34</u>	<u>21.60</u>	<u>20.54</u>	<u>20.11</u>	<u>19.77</u>	18.45	17.51	341.5
	MambaSR [58]	27.02	24.44	23.01	22.06	21.36	20.34	19.95	19.65	18.29	17.48	90.5
	GaussianSR [23]	26.20	23.76	22.35	21.38	20.66	19.68	19.31	19.03	17.86	17.07	321.4
	Ours	28.22	25.43	23.87	22.86	22.08	20.95	20.54	20.21	18.77	17.70	4.6
DIV2K [1]	MetaSR [24]	29.33	27.03	25.66	24.69	23.94	22.82	22.39	22.01	20.42	19.25	<u>123.5</u>
	LIIF [11]	29.27	26.99	25.60	24.63	23.89	22.77	22.34	21.94	20.36	19.19	480.6
	LTE [34]	29.50	27.20	25.81	24.84	24.09	22.94	22.50	22.12	20.50	<u>19.31</u>	1407.5
	SRNO [55]	29.42	27.12	25.74	24.77	24.03	22.90	22.46	22.06	20.47	19.27	390.9
	CiaoSR [3]	<u>29.59</u>	<u>27.28</u>	<u>25.89</u>	<u>24.91</u>	<u>24.15</u>	<u>22.99</u>	<u>22.54</u>	<u>22.16</u>	<u>20.50</u>	19.30	1857.8
	MambaSR [58]	29.36	27.08	25.70	24.74	23.99	22.87	22.44	22.05	20.46	19.27	398.3
	GaussianSR [23]	29.03	26.73	25.29	24.23	23.44	22.26	21.81	21.42	19.90	18.76	4962.8
	Ours	29.80	27.47	26.07	25.08	24.33	23.18	22.74	22.35	20.68	19.45	4.7
LSDIR [35]	MetaSR [24]	26.54	24.64	23.54	22.79	22.24	21.42	21.09	20.80	19.62	18.68	50.4
	LIIF [11]	26.49	24.59	23.49	22.75	22.21	21.40	21.09	20.75	19.59	18.65	226.4
	LTE [34]	26.73	24.78	23.65	22.88	22.33	<u>21.48</u>	<u>21.15</u>	<u>20.85</u>	<u>19.66</u>	<u>18.71</u>	451.5
	SRNO [55]	26.65	24.72	23.61	22.85	22.30	21.45	21.12	20.83	19.64	18.69	163.6
	CiaoSR [3]	<u>26.80</u>	<u>24.84</u>	<u>23.69</u>	<u>22.92</u>	<u>22.35</u>	21.48	21.14	20.84	19.63	18.67	1289.3
	MambaSR [58]	26.62	24.69	23.59	22.83	22.28	21.44	21.11	20.82	19.64	18.70	197.7
	GaussianSR [23]	26.25	24.39	23.28	22.49	21.92	21.06	20.74	20.45	19.29	18.38	1284.3
	Ours	27.14	25.07	23.91	23.13	22.54	21.89	21.35	21.06	19.79	18.82	4.6

knowledge and avoiding the local optima observed in method (a), as demonstrated in the ablation study in Section 4.2.

4.4. Adaptive Position Drifting

The position parameters are also critical for Gaussian kernels, as they determine their locations in the 2D space. Directly learning these positions is highly challenging, as demonstrated in Section 4.2. Since each LR pixel typically corresponds to multiple pixels in the HR image, a straightforward solution is to fix the positions at the centers of the LR pixels. While this strategy simplifies the optimization process, it significantly limits the model’s representational capacity, making it difficult to adaptively learn the position distribution based on image content.

To address the above issues, we propose a novel method, Adaptive Position Drifting (APD), which not only ensures efficient optimization but also improves representational capacity by allowing the model to adaptively learn positions, as shown in Figure 3. Specifically, we use the center positions of LR pixels as the initialized positions P_{init} and further introduce a dynamic offset mechanism, which learns a dynamic offset from LR features \mathcal{F}_{LR} by \mathcal{M}_{pos} model to adjust the spatial positions. Here, we set the offset range from $-1 \sim 1$ by the Tanh activate function and add the offset P_{off} to the initialized LR center positions to obtain the final

positions P_{final} , as expressed by the following equation:

$$P_{\text{off}} = \text{Tanh}(\mathcal{M}_{\text{pos}}(\mathcal{F}_{\text{LR}})), \quad (11)$$

$$P_{\text{final}} = P_{\text{init}} + P_{\text{off}}. \quad (12)$$

where \mathcal{M}_{pos} is implemented using five multilayer perceptron layers. This P_{off} enables the network to adaptively learn kernel positions based on image content, resulting in denser kernel placement in regions with richer textures and enhancing the network’s performance, as demonstrated in Figure 3, Section 5.3, Appendix Section 13.

In addition, since the RGB is range from 0 to 1 and is relatively easy to optimize, we introduce a simple Color Gaussian Mapping (CGM) to learn the RGB parameters. Specifically, this mapping is implemented using 5 multilayer perceptron (MLP) layers applied to \mathcal{F}_{LR} . In summary, the above three components construct our proposed ContinuousSR framework, as shown in Figure 3.

5. Experiment and Analysis

5.1. Experiment Setting

Datasets. We use the commonly employed DF2K high-quality dataset [52] as HR images, which are degraded using bicubic to generate LR for training. For evaluation, we adopt Set5 [2], Set14 [61], B100 [42], Urban100 [26], Manga109 [43], DIV2K validation [1] and LSDIR [35].

Table 2. Performance comparison of the Urban100 benchmark on SSIM, FID, and DISTs metrics.

Metrics	Methods	$\times 4$	$\times 6$	$\times 8$	$\times 10$	$\times 12$	$\times 16$	$\times 18$	$\times 20$	$\times 32$	$\times 48$
SSIM↑	LIIF [11]	0.7911	0.6861	0.6148	0.5642	0.5270	0.4790	0.4617	0.4503	0.4106	0.3918
	LTE [34]	0.8069	0.7045	0.6321	0.5810	0.5422	0.4900	0.4710	0.4588	0.4145	0.3931
	GaussianSR [23]	0.7751	0.6633	0.5867	0.5334	0.4967	0.4521	0.4369	0.4277	0.3969	0.3835
	CiaoSR [3]	0.8110	0.7126	0.6415	0.5887	0.5503	0.4974	0.4777	0.4637	0.4168	0.3921
	Ours	0.8292	0.7343	0.6624	0.6089	0.5683	0.5097	0.4893	0.4746	0.4216	0.3958
DISTs↓	LIIF [11]	0.1611	0.2178	0.2589	0.2926	0.3209	0.3659	0.3835	0.3990	0.4678	0.5322
	LTE [34]	0.1570	0.2126	0.2541	0.2872	0.3157	0.3611	0.3799	0.3960	0.4695	0.5362
	GaussianSR [23]	0.1740	0.2374	0.2890	0.3296	0.3631	0.4109	0.4302	0.4466	0.5157	0.5713
	CiaoSR [3]	0.1533	0.2074	0.2453	0.2771	0.3049	0.3510	0.3701	0.3863	0.4513	0.4998
	Ours	0.1356	0.1901	0.2299	0.2601	0.2860	0.3324	0.3504	0.3670	0.4439	0.5144
FID↓	LIIF [11]	4.76	24.87	50.05	77.54	102.47	145.25	164.41	179.44	256.95	311.09
	LTE [34]	3.84	21.01	45.15	70.56	92.85	136.54	156.24	170.92	253.52	296.27
	GaussianSR [23]	5.64	29.30	57.00	89.25	120.02	166.20	181.59	202.18	264.27	315.97
	CiaoSR [3]	3.74	20.48	43.25	58.60	92.58	133.77	151.70	168.89	247.84	294.49
	Ours	2.91	16.50	37.09	58.83	78.72	116.30	130.32	143.52	216.21	281.27

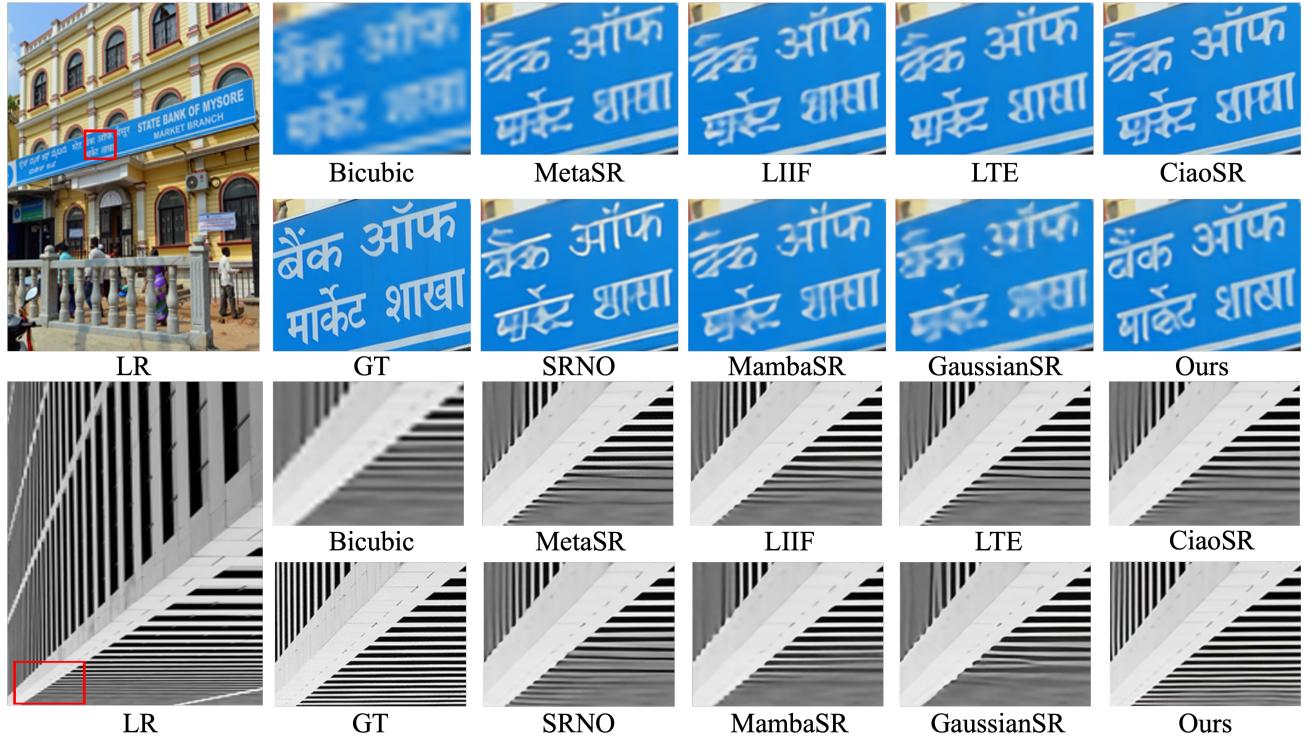


Figure 4. Qualitative comparison. The visual quality of our method outperforms existing methods. Please zoom in for a better view.

Evaluation metrics. Following previous work [11, 34], we use PSNR, SSIM [53], FID [22], and DISTs [14] for evaluation. Note that the PSNR/SSIM value is calculated on the RGB channels for the DIV2K validation set and on the Y channel (*i.e.*, luminance) of the transformed YCbCr space for the other benchmark test sets.

Implementation details. Following previous works [11, 34], we adopt the same way to generate paired images for training. Specifically, initially, we crop image patches of

size 256×256 as ground truth. Then, we use bicubic downsampling to generate corresponding LR images, and the downsampling scaling factor is sampled from a uniform distribution $U(4, 8)$. We employ SwinIR [37] and HAT [10] as backbones. Adam [31] is used as the optimizer, with the initial learning rate setting to $1e-4$ and decaying by a factor of 0.5 every 100 epochs. We utilize the L1 loss [11] and frequency loss [12] for training, with a total batch size of 64 and 1000 training epochs on 8 V100 GPUs.

Table 3. Memory usage (G) comparison with different methods.

Memory Usage	$\times 4$	$\times 6$	$\times 8$	$\times 12$	$\times 16$
LIIF	4.12	6.49	9.79	19.27	OOM
CiaoSR	12.17	22.83	OOM	OOM	OOM
Ours	2.48	2.49	2.50	2.52	2.54

Compared methods. We compare with nine state-of-the-art and popular models: MetaSR [24], LIIF [11], LTE [34], ITSRN [59], SRNO [55], CiaoSR [3], MambaSR [58], GaussianSR [23] and GSASR [6]. The best version from their official code is used for comparison. Details are provided in Appendix Section 8.

5.2. Quantitative and Qualitative Results

Quantitative comparisons. As shown in Table 1 and Table 2, our method achieves the best performance compared to existing approaches across all evaluation metrics and benchmarks. For example, on the Urban100 dataset, our method surpasses the current state-of-the-art (CiaoSR) by 0.80 dB in the $\times 4$, representing a substantial improvement. Similarly, in terms of SSIM and FID, our method achieves further gains in the $\times 4$ scenario, surpassing the current state-of-the-art by 0.0172 and 0.83, respectively. These results demonstrate the effectiveness and superiority of the proposed method.

Complexity comparisons. We present comparisons of runtime and memory usage. Specifically, we evaluate the average runtime across 45 different scales, ranging from $\times 4$ to $\times 48$. As shown in Table 1, our method significantly outperforms existing methods in terms of speed. For instance, it surpasses the current state-of-the-art method (CiaoSR) by nearly 280 times on the LSDIR dataset. Moreover, we also provide comparisons of memory usage. Specifically, we set the input size to 48×48 , disable the tiling strategy, and test the memory usage under different scales. As shown in Table 3, thanks to our efficient pipeline design, our method maintains minimal computational overhead across different scales. In contrast, existing methods, such as LIIF and CiaoSR, fail to handle larger scales and encounter OOM (out of memory) on V100 GPUs.

Qualitative comparisons. We present qualitative comparisons, as shown in Figure 4. Compared to existing methods, our approach reconstructs sharper and more visually pleasing details that are consistent with the ground truth (GT). For instance, in the bottom part of Figure 4, our method effectively reconstructs the texture details inside the building. This highlights the superiority of our method in generating realistic and perceptually satisfying results.

More visual comparisons, user studies, benchmark results, FLOPs comparisons, and details are provided in Appendix Sections 10, 9, and 12.

Table 4. Ablation studies on proposed APD, DDCW, and DGP.

DDCW	APD	PSNR	P_{init}	P_{off}	PSNR	K	PSNR
✓		10.5	✓		27.8	\mathcal{K}_1	27.7
	✓	12.3		✓	10.5	\mathcal{K}_2	27.1
✓	✓	28.2	✓	✓	28.2	\mathcal{K}_{DCP}	28.2

Table 5. Performance comparison under more challenging low-resolution and rainy conditions.

Methods	$\times 4$	$\times 5$	$\times 6$	$\times 7$	$\times 8$
LIIF	24.04	23.53	23.09	22.70	22.39
GaussianSR	24.04	23.51	23.08	22.69	22.38
CiaoSR	23.84	23.45	23.01	22.66	22.33
Ours	24.51	23.95	23.48	23.07	22.76

5.3. Ablation Study

In this section, we present ablation studies on our proposed APD, DDCW, and DGP using the Urban100 $\times 4$ dataset. Specifically, we independently remove APD and DDCW, as illustrated in the right part of Table 4. The exclusion of these modules significantly exacerbates the optimization difficulty, resulting in a considerable decline in PSNR. Then, we validate the effectiveness of P_{init} and P_{off} in APD. As shown in the middle of Table 4, the results demonstrate that the model achieves the best representational capacity and performance when both components are employed. Finally, we evaluate the effectiveness of DGP in DDCW \mathcal{K}_{DCP} . Specifically, we remove the DGP and separately modify the covariance range to $[0,1]$ and $[0,10]$, using uniform sampling to construct \mathcal{K}_1 and \mathcal{K}_2 . As shown in the right in Table 4, incorporating DGP provides a better basis function, thereby enhancing performance. More ablation studies are provided in Appendix Sections 11.

6. Future Work

It is well known that, in real-world scenarios, image degradation is not limited to low resolution but often includes other types of degradation, such as rain and noise. The goal of low-level vision is to remove these degradations while enhancing image resolution and quality. To this end, we evaluate our method on Rain200H [60], simulating low-resolution rainy images with bicubic downsampling. We compare our approach with three existing state-of-the-art methods to validate its effectiveness. As shown in Table 5, our method removes rain degradations more effectively while enhancing resolution and details, outperforming existing methods. This demonstrates the potential of our approach for other low-level vision tasks. In future work, we aim to extend it to more tasks to further enhance its applicability.

To provide a comprehensive understanding of our meth-

ods, we include detailed explanations, additional comparisons, analyses, limitations, future work, and extensive visual examples in the **Appendix**, showcasing superiority.

7. Conclusion

We introduce ContinuousSR, a novel Pixel-to-Gaussian paradigm designed for fast and high-quality arbitrary-scale super-resolution. By explicitly reconstructing 2D continuous HR signals from LR images using Gaussian Splatting, ContinuousSR significantly improves both efficiency and performance. Through statistical analysis, we uncover the Deep Gaussian Prior (DGP) and propose a DGP-driven Covariance Weighting mechanism along with an Adaptive Position Drifting strategy. These innovations improve the quality and fidelity of the reconstructed Gaussian fields. Experiments on seven popular benchmarks demonstrate that our method outperforms state-of-the-art methods in both quality and speed, achieving a 19.5 \times speed improvement and 0.90dB PSNR improvement, making it a promising solution for ASSR tasks.

References

- [1] Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 126–135, 2017. 6
- [2] Marco Bevilacqua, Aline Roumy, Christine Guillemot, and Marie Line Alberi-Morel. Low-complexity single-image super-resolution based on nonnegative neighbor embedding. 2012. 6, 1, 2
- [3] Jiezhang Cao, Qin Wang, Yongqin Xian, Yawei Li, Bingbing Ni, Zhiming Pi, Kai Zhang, Yulun Zhang, Radu Timofte, and Luc Van Gool. Ciaosr: Continuous implicit attention-in-attention 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, 3, 6, 7, 8, 1
- [4] Kenneth R Castleman. *Digital image processing*. Prentice Hall Professional Technical Reference, 1979. 2, 3
- [5] Lukas Cavigelli, Pascal Hager, and Luca Benini. Cas-cnn: A deep convolutional neural network for image compression artifact suppression. In *2017 International Joint Conference on Neural Networks*, pages 752–759, 2017. 2
- [6] Du Chen, Liyi Chen, Zhengqiang Zhang, and Lei Zhang. Generalized and efficient 2d gaussian splatting for arbitrary-scale super-resolution. *arXiv preprint arXiv:2501.06838*, 2025. 3, 8, 2
- [7] Guikun Chen and Wenguan Wang. A survey on 3d gaussian splatting. *arXiv preprint arXiv:2401.03890*, 2024. 3
- [8] 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. 2, 3
- [9] Xiang Chen, Hao Li, Mingqiang Li, and Jinshan Pan. Learning a sparse transformer network for effective image deraining. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5896–5905, 2023. 4
- [10] Xiangyu Chen, Xintao Wang, Wenlong Zhang, Xiangtao Kong, Yu Qiao, Jiantao Zhou, and Chao Dong. Hat: Hybrid attention transformer for image restoration. 2023. 2, 7
- [11] Yinbo Chen, Sifei Liu, and Xiaolong Wang. Learning continuous image representation with local implicit image function. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 8628–8638, 2021. 2, 3, 6, 7, 8, 1
- [12] Yuning Cui, Wenqi Ren, Xiaochun Cao, and Alois Knoll. Focal network for image restoration. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 13001–13011, 2023. 7
- [13] Xin Di, Long Peng, Peizhe Xia, Wenbo Li, Renjing Pei, Yang Cao, Yang Wang, and Zheng-Jun Zha. Qmambabsr: Burst image super-resolution with query state space model. *arXiv preprint arXiv:2408.08665*, 2024. 2
- [14] Keyan Ding, Kede Ma, Shiqi Wang, and Eero P Simoncelli. Image quality assessment: Unifying structure and texture similarity. *IEEE transactions on pattern analysis and machine intelligence*, 44(5):2567–2581, 2020. 7
- [15] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image super-resolution. In *European Conference on Computer Vision*, pages 184–199, 2014. 2
- [16] Jiajun Dong, Chengkun Wang, Wenzhao Zheng, Lei Chen, Jiwen Lu, and Yansong Tang. Gaussiantoken: An effective image tokenizer with 2d gaussian splatting. *arXiv preprint arXiv:2501.15619*, 2025. 3
- [17] Minghong Duan, Linhao Qu, Shaolei Liu, and Manning Wang. Local implicit wavelet transformer for arbitrary-scale super-resolution. *arXiv preprint arXiv:2411.06442*, 2024. 3
- [18] Ben Fei, Jingyi Xu, Rui Zhang, Qingyuan Zhou, Weidong Yang, and Ying He. 3d gaussian splatting as new era: A survey. *IEEE Transactions on Visualization and Computer Graphics*, 2024. 3
- [19] Huiyuan Fu, Fei Peng, Xianwei Li, Yejun Li, Xin Wang, and Huadong Ma. Continuous optical zooming: A benchmark for arbitrary-scale image super-resolution in real world. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3035–3044, 2024. 3, 1
- [20] Zongyao He and Zhi Jin. Dynamic implicit image function for efficient arbitrary-scale super-resolution. In *2024 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6. IEEE, 2024.
- [21] Zongyao He and Zhi Jin. Latent modulated function for computational optimal continuous image representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26026–26035, 2024. 3
- [22] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilib-

- rium. *Advances in neural information processing systems*, 30, 2017. 7
- [23] Jintong Hu, Bin Xia, Bin Chen, Wenming Yang, and Lei Zhang. Gaussiansr: High fidelity 2d gaussian splatting for arbitrary-scale image super-resolution. *arXiv preprint arXiv:2407.18046*, 2024. 3, 6, 7, 8, 1, 2
- [24] Xuecai Hu, Haoyuan Mu, Xiangyu Zhang, Zilei Wang, Tieniu Tan, and Jian Sun. Meta-sr: A magnification-arbitrary network for super-resolution. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1575–1584, 2019. 2, 3, 6, 8, 1
- [25] Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2d gaussian splatting for geometrically accurate radiance fields. In *ACM SIGGRAPH 2024 conference papers*, pages 1–11, 2024. 3
- [26] Jia-Bin Huang, Abhishek Singh, and Narendra Ahuja. Single image super-resolution from transformed self-exemplars. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5197–5206, 2015. 6, 1, 2
- [27] Anil K Jain. *Fundamentals of digital image processing*. Prentice-Hall, Inc., 1989. 2, 3
- [28] Shuguo Jiang, Nanying Li, Meng Xu, Shuyu Zhang, and Sen Jia. Sqformer: Spectral-query transformer for hyperspectral image arbitrary-scale super-resolution. *IEEE Transactions on Geoscience and Remote Sensing*, 2024. 3
- [29] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 3
- [30] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1646–1654, 2016. 2
- [31] Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 7
- [32] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4681–4690, 2017. 2
- [33] 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. 2, 3
- [34] Jaewon Lee and Kyong Hwan Jin. Local texture estimator for implicit representation function. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1929–1938, 2022. 2, 3, 6, 7, 8, 1
- [35] Yawei Li, Kai Zhang, Jingyun Liang, Jiezhang Cao, Ce Liu, Rui Gong, Yulun Zhang, Hao Tang, Yun Liu, Denis Demadola, et al. Lsdir: A large scale dataset for image restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1775–1787, 2023. 6
- [36] Zekun Li, Hongying Liu, Fanhua Shang, Yuanyuan Liu, Liang Wan, and Wei Feng. Savsr: arbitrary-scale video super-resolution via a learned scale-adaptive network. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 3288–3296, 2024. 2
- [37] Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. In *IEEE International Conference on Computer Vision Workshops*, pages 1833–1844, 2021. 2, 7, 1
- [38] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 136–144, 2017. 2
- [39] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 136–144, 2017. 5
- [40] Anran Liu, Yihao Liu, Jinjin Gu, Yu Qiao, and Chao Dong. Blind image super-resolution: A survey and beyond. *IEEE transactions on pattern analysis and machine intelligence*, 45(5):5461–5480, 2022. 2
- [41] Hongying Liu, Zekun Li, Fanhua Shang, Yuanyuan Liu, Liang Wan, Wei Feng, and Radu Timofte. Arbitrary-scale super-resolution via deep learning: A comprehensive survey. *Information Fusion*, 102:102015, 2024. 2
- [42] David Martin, Charless Fowlkes, Doron Tal, and Jitendra Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*, pages 416–423. IEEE, 2001. 6, 1, 2
- [43] Yusuke Matsui, Kota Ito, Yuji Aramaki, Azuma Fujimoto, Toru Ogawa, Toshihiko Yamasaki, and Kiyoharu Aizawa. Sketch-based manga retrieval using manga109 dataset. *Multimedia Tools and Applications*, 76:21811–21838, 2017. 6, 1, 2
- [44] Long Peng, Yang Cao, Renjing Pei, Wenbo Li, Jiaming Guo, Xueyang Fu, Yang Wang, and Zheng-Jun Zha. Efficient real-world image super-resolution via adaptive directional gradient convolution. *arXiv preprint arXiv:2405.07023*, 2024. 2
- [45] Long Peng, Wenbo Li, Renjing Pei, Jingjing Ren, Yang Wang, Yang Cao, and Zheng-Jun Zha. Towards realistic data generation for real-world super-resolution. *arXiv preprint arXiv:2406.07255*, 2024. 2
- [46] Douglas A Reynolds et al. Gaussian mixture models. *Encyclopedia of biometrics*, 741(659-663):3, 2009. 3
- [47] Wei Shang, Dongwei Ren, Wanying Zhang, Yuming Fang, Wangmeng Zuo, and Kede Ma. Arbitrary-scale video super-resolution with structural and textural priors. In *European Conference on Computer Vision*, pages 73–90. Springer, 2024. 3
- [48] Radu Timofte, Eirikur Agustsson, Luc Van Gool, Ming-Hsuan Yang, and Lei Zhang. Ntire 2017 challenge on single image super-resolution: Methods and results. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 114–125, 2017. 5
- [49] Yi Ting Tsai, Yu Wei Chen, Hong-Han Shuai, and Ching-Chun Huang. Arbitrary-resolution and arbitrary-scale face

- super-resolution with implicit representation networks. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 4270–4279, 2024. 3
- [50] Wenbo Wan, Zezhu Wang, Zhiyan Wang, Lingchen Gu, Jiande Sun, and Qiang Wang. Arbitrary-scale image super-resolution via degradation perception. *IEEE Transactions on Computational Imaging*, 2024. 2
- [51] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *European Conference on Computer Vision Workshops*, pages 701–710, 2018. 2
- [52] Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world blind super-resolution with pure synthetic data. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1905–1914, 2021. 6
- [53] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 7
- [54] Zhihao Wang, Jian Chen, and Steven CH Hoi. Deep learning for image super-resolution: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 43(10):3365–3387, 2020. 2
- [55] 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. 2, 3, 6, 8, 1
- [56] Peizhe Xia, Long Peng, Xin Di, Renjing Pei, Yang Wang, Yang Cao, and Zheng-Jun Zha. S3mamba: Arbitrary-scale super-resolution via scaleable state space model. *arXiv preprint arXiv:2411.11906*, 2024. 2, 3
- [57] Xingqian Xu, Zhangyang Wang, and Humphrey Shi. Ultrasr: Spatial encoding is a missing key for implicit image function-based arbitrary-scale super-resolution. *arXiv preprint arXiv:2103.12716*, 2021. 2
- [58] Jin Yan, Zongren Chen, Zhiyuan Pei, Xiaoping Lu, and Hua Zheng. Mambasr: Arbitrary-scale super-resolution integrating mamba with fast fourier convolution blocks. *Mathematics*, 12(15):2370, 2024. 6, 8, 1
- [59] Jingyu Yang, Sheng Shen, Huanjing Yue, and Kun Li. Implicit transformer network for screen content image continuous super-resolution. *Advances in Neural Information Processing Systems*, 34:13304–13315, 2021. 8, 2
- [60] Wenhao Yang, Robby T Tan, Jiashi Feng, Jiaying Liu, Zongming Guo, and Shuicheng Yan. Deep joint rain detection and removal from a single image. In *CVPR*, pages 1357–1366, 2017. 8
- [61] Roman Zeyde, Michael Elad, and Matan Protter. On single image scale-up using sparse-representations. In *Curves and Surfaces: 7th International Conference, Avignon, France, June 24–30, 2010, Revised Selected Papers 7*, pages 711–730. Springer, 2012. 6, 1, 2
- [62] Haonan Zhang, Jie Guo, Jiawei Zhang, Haoyu Qin, Zesen Feng, Ming Yang, and Yanwen Guo. Deep fourier-based arbitrary-scale super-resolution for real-time rendering. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–11, 2024. 3
- [63] Kai Zhang, Yawei Li, Wangmeng Zuo, Lei Zhang, Luc Van Gool, and Radu Timofte. Plug-and-play image restoration with deep denoiser prior. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 2
- [64] Xinjie Zhang, Xingtong Ge, Tongda Xu, Dailan He, Yan Wang, Hongwei Qin, Guo Lu, Jing Geng, and Jun Zhang. Gaussianimage: 1000 fps image representation and compression by 2d gaussian splatting. In *European Conference on Computer Vision*, pages 327–345. Springer, 2024. 3, 4, 1
- [65] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *European Conference on Computer Vision*, pages 286–301, 2018. 2
- [66] Yaoqian Zhao, Qizhi Teng, Honggang Chen, Shujiang Zhang, Xiaohai He, Yi Li, and Ray E Sheriff. Activating more information in arbitrary-scale image super-resolution. *IEEE Transactions on Multimedia*, 2024. 3
- [67] Jinchen Zhu, Mingjian Zhang, Ling Zheng, and Shizhuang Weng. Multi-scale implicit transformer with re-parameterization for arbitrary-scale super-resolution. *Pattern Recognition*, 162:111327, 2025. 3

Pixel to Gaussian: Ultra-Fast Continuous Super-Resolution with 2D Gaussian Modeling

Supplementary Material

8. Details of Compared Methods

To validate the effectiveness of our proposed model, we compare it against seven state-of-the-art (SOTA) and widely adopted models: MetaSR [24], LIIF [11], LTE [34], SRNO [55], CiaoSR [3], MambaSR [58] and GaussianSR [23]. For a fair comparison, we select the best-performing networks for each method based on their official GitHub repositories. Specifically, we use the MetaSR model based on SwinIR, the LIIF model based on RDN, the LTE model based on SwinIR, the SRNO model based on RDN, the CiaoSR model based on SwinIR, the MambaSR model based on RDN, and the GaussianSR model on EDSR-baseline.

9. User Study

To further assess visual quality, we conduct a user study. Ten images are randomly selected from the test datasets, ensuring diversity in image content and complexity. Fifteen participants rate the visual quality of each processed image on a scale from 0 (poor) to 10 (excellent). Each participant evaluates the images independently to ensure unbiased results. As shown in Figure 5, the results demonstrate that existing methods frequently fail to restore high-quality images, particularly in challenging regions with fine details and textures. This leads to lower user satisfaction, with average scores ranging between 4.2 and 7.0 for most competing methods. In contrast, our method achieves the highest average score of 7.7, significantly outperforming all other approaches. The superior performance of our method demonstrates its ability to produce sharper details, better texture preservation, and visually consistent results. Participants consistently note that our method outperforms others, particularly in challenging regions, further validating its effectiveness and generalization in restoring high-quality images.

10. Additional Comparison

More Benchmarks. To further demonstrate the superiority of our proposed method, we conduct experiments to compare its performance against existing methods using the same SwinIR [37] backbone on the Set5 [2], Set14 [61], B100 [42], Urban100 [26] and Manga109 [43] datasets. As shown in Table 6, our method still achieves state-of-the-art performance across all benchmarks and scales.

More Comparisons on Real Datasets. To further demonstrate the superiority of the proposed method in real-world

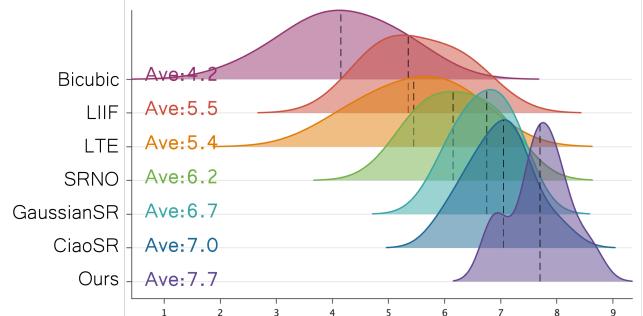


Figure 5. User study.

scenarios, we compare it with existing methods on the real dataset COZ [19] on $\times 5$. The results, as shown in Table 8, indicate that our method consistently outperforms existing approaches in real-world scenarios, validating the superior generalization ability of our method to real-world data.

More Complexity Comparisons. Furthermore, we also provide comparisons of FLOPs and inference time at a single scale on the Manga109 dataset. Specifically, considering that the image shapes in the original dataset may cause other methods to run out of memory, we fix the GT shape to 288 and evaluate the FLOPs and inference time at different scales. As shown in Table 7, our method not only achieves the best performance in terms of PSNR but also maintains the lowest FLOPs and inference time at a single scale, significantly outperforming the current SOTA method, CiaoSR. Moreover, Table 1 in the main text further demonstrates our superiority in total runtime across multiple scales. These results fully validate the efficiency and superiority of the proposed method.

More Performance Metrics. In the main text, we have provided PSNR, SSIM, FID, and DISTS metrics to demonstrate the superiority of the proposed method. Here, we further present a comparison of LPIPS performance on the Urban100 $\times 4$ dataset. As shown in Table 9, our method achieves the best performance in terms of LPIPS. This further validates the superiority of the proposed method in perceptual quality.

More Compared Method. In the main text, we have compared our proposed method with 9 existing methods to demonstrate its superiority. Additionally, we include a comparison with GaussianImage [64]. Specifically, we conduct experiments on the Set5 and Set14 datasets under the $\times 4$ scenario. Since GaussianImage is an optimization-based end-to-end algorithm, we allow this method to optimize on

Table 6. Performance comparison with existing methods using the same SwinIR [37] backbone on the Set5 [2], Set14 [61], B100 [42], Urban100 [26] and Manga109 [43] datasets. Table performance is referred to in [3].

Dataset	Scale	SwinIR [37]	MetaSR [24]	LIIF [11]	ITSRN [59]	LTE [34]	CiaoSR [3]	Ours
Set5 [2]	$\times 4$	32.72	32.47	32.73	32.63	32.81	<u>32.84</u>	32.93
	$\times 6$	-	29.09	29.46	29.31	29.50	<u>29.62</u>	29.67
	$\times 8$	-	27.02	27.36	27.24	27.35	<u>27.45</u>	27.55
	$\times 12$	-	24.82	-	24.79	-	<u>24.96</u>	25.18
Set14 [61]	$\times 4$	28.94	28.85	28.98	28.97	29.06	<u>29.08</u>	29.18
	$\times 6$	-	26.58	26.82	26.71	26.86	<u>26.88</u>	26.96
	$\times 8$	-	25.09	25.34	25.32	25.42	<u>25.42</u>	25.54
	$\times 12$	-	23.33	-	23.30	-	<u>23.38</u>	23.55
B100 [42]	$\times 4$	27.83	27.75	27.84	27.85	27.86	<u>27.90</u>	27.94
	$\times 6$	-	25.94	26.07	26.05	26.09	<u>26.13</u>	26.18
	$\times 8$	-	24.87	25.01	24.96	25.03	<u>25.07</u>	25.13
	$\times 12$	-	23.59	-	23.57	-	<u>23.68</u>	23.76
Urban100 [26]	$\times 4$	27.07	26.76	27.15	27.12	27.24	<u>27.42</u>	27.51
	$\times 6$	-	24.16	24.59	24.50	24.62	<u>24.84</u>	24.93
	$\times 8$	-	22.75	23.14	23.06	23.17	<u>23.34</u>	23.45
	$\times 12$	-	21.31	-	21.34	-	<u>21.60</u>	21.72
Manga109 [43]	$\times 4$	31.67	31.37	31.71	31.74	31.79	<u>31.91</u>	32.19
	$\times 6$	-	27.29	27.69	27.72	27.83	<u>28.01</u>	28.25
	$\times 8$	-	24.96	25.28	25.23	25.42	<u>25.61</u>	25.80
	$\times 12$	-	22.35	-	22.47	-	<u>22.79</u>	22.96

Table 7. Comparison of PSNR (dB), FLOPs (G), and running time (ms) on the Manga109 dataset.

	$\times 4$			$\times 6$			$\times 8$		
	PSNR	FLOPs	Times	PSNR	FLOPs	Times	PSNR	FLOPs	Times
SwinIR-LIIF	31.71	365.11	150	27.69	289.77	141	25.28	271.35	137
SwinIR-CiaoSR	31.91	1949.71	319	28.01	1275.95	256	25.61	1048.08	235
SwinIR-ContinuousSR	32.19	280.99	136	28.25	125.63	112	25.80	79.21	99

LR inputs and adjust the Gaussian mapping scale to perform super-resolution for comparison. As shown in Table 15, this method fails to learn the mapping from LR to HR, resulting in poor performance. Furthermore, it is worth noting that GaussianImage requires nearly 1 minute of optimization per scene on a V100 GPU, which is impractical for real-world applications.

More Compared Methods with GS. Several recent GS-based ASSR methods have been proposed, such as GaussianSR [23] and GSASR [6]. GaussianSR has already been thoroughly analyzed and compared in the main text. Here, we focus on analyzing and comparing GSASR. Although GSASR has made notable progress, it is still constrained by inefficiencies caused by multiple upsampling and decoding processes across different scales. Furthermore, GSASR performs GS in the feature and image space, which makes it struggle to ensure the continuity of reconstructed images across different scales, leading to low performance. In con-

trast, our method leverages 2D GS modeling to reconstruct continuous HR images, enabling both fast and high-quality ASSR. Although the GSASR method has not been open-sourced, we still compare our method against the performance reported in its paper. For example, on the LSDIR benchmark, our method achieves a performance of 27.14 dB at $\times 4$, significantly surpassing GSASR’s best reported performance of 26.73 dB. This demonstrates the superiority of our method in terms of performance. Moreover, in terms of speed, our method requires only 1 ms to generate high-quality HR images across different scales, whereas GSASR takes approximately 91–1573 ms. This further highlights the ultra-fast speed of our proposed method.

More Details in Section 5.2. In Table 1, the Average Time (AT) is calculated by performing super-resolution on LR images across 45 different scales, ranging from $\times 4$ to $\times 48$, and then averaging the total runtime. For each dataset, we select a representative LR shape and downsample it by a

Table 8. Comparison of PSNR on the COZ dataset.

COZ	MetaSR	LIIF	LTE	LINF	SRNO	LIT	CiaoSR	LMI	Ours
PSNR	24.39	24.39	24.4	24.32	24.4	24.36	24.38	24.48	24.68

Table 9. LPIPS \downarrow comparison for Urban100 dataset across different methods.

LPIPS \downarrow	MetaSR	LIIF	MambaSR	LTE	SRNO	CiaoSR	GaussianSR	Ours
Urban100	0.1989	0.2080	0.2073	0.1934	0.1991	0.1872	0.2285	0.1803

factor of 48 to construct the input size for each dataset, ensuring that existing ASSR methods do not encounter out-of-memory (OOM) issues. Specifically, the LR size is 21×13 for Urban100, 42×28 for DIV2K, and 29×19 for LSDIR. As shown in Table 1, our method consistently achieves significant speed advantages over existing methods across different datasets and LR shapes.

11. Additional Ablation Study

Due to space limitations in the main text, we provide additional ablation experiments to demonstrate the effectiveness and rationality of the proposed method. Below, we present detailed descriptions of additional ablation studies and implementation details.

Ablation Study on \mathcal{K} . In the DGP-Driven Covariance Weighting, considering the difficulty for deep learning networks to directly interpret the specific meaning of covariance, we map \mathcal{K} from its original three-dimensional representation (*i.e.*, $(\sigma_x^2, \sigma_y^2, \text{and } \rho)$) to a latent representation space through a convolutional neural network. This approach facilitates better convergence and achieves improved performance. Specifically, we explore the performance when the dimension of the latent space is set to 3, 256, and 512, as shown in Table 10. It can be observed that the best performance is achieved when the dimension is set to 512, showing significant improvements compared to the original three-dimensional setting.

Ablation Study on number of \mathcal{K} . Furthermore, we investigate the impact of the number of \mathcal{K} on the network’s performance. We define 100, 500, and 730 Gaussian covariances, and the results are presented in Table 11. It can be seen that as the number of covariances increases, the network’s performance improves. However, beyond 730, no further performance gains are observed. Therefore, in this work, we set the number of covariances to 730.

Ablation Study on N . Additionally, we study the effect of the number of Gaussians N on the network’s performance. As described in the main text, we initialize one Gaussian kernel at the center of each LR pixel. We further explore the impact of introducing more Gaussian kernels per unit pixel, and the results are shown in Table 12. It can be observed

that the best performance is achieved when the number of kernels is set to 4. Introducing too many kernels increases the optimization complexity, which does not lead to further performance gains. Therefore, we set the number of kernels per pixel to 4 in this work.

Ablation Study on P_{off} . Finally, we examine the impact of the range of P_{off} , denoted as $[-A, A]$, on the network’s performance. Specifically, we set the range to 0.5, 1, and 2, and the results are shown in Table 13. It can be observed that the best performance is achieved when the range is set to 1. Both overly large and overly small ranges have adverse effects on performance.

12. Additional Visual Comparison Results

In this section, we present additional visual comparison results to further demonstrate the superiority of our proposed method, as shown in Figure 7, 8 and 9. It can be observed that our method achieves the best visual satisfaction in terms of detailed textures, while also preserving the highest level of detail fidelity, making it closest to the GT image.

13. Visualization of the Position Distribution

To demonstrate the superior adaptive perception capability of the proposed offset mechanism, which effectively introduces more Gaussian kernels in complex texture regions based on image content, we visualize the learned Gaussian position distribution. As shown in Figure 6, the proposed Adaptive Position Drifting adjusts the original initialization of the position distribution by adaptively perceiving the structural content of the image. The results reveal that regions with richer textures have higher densities of Gaussian kernels.

14. Algorithm Workflow

To clearly demonstrate the details of the proposed method, we design an algorithm workflow, as illustrated in Algorithm 1. This workflow describes the key steps from input to output, including feature encoding, color prediction, offset prediction, covariance estimation, and the final image reconstruction process.

Table 10. Ablation on Dim.

	PSNR	SSIM
3	28.14	0.8273
256	28.17	0.8281
512	28.22	0.8292

Table 11. Number of K .

	PSNR	SSIM
100	28.12	0.8277
500	28.19	0.8286
730	28.22	0.8292

Table 12. Number of N .

	PSNR	SSIM
1	28.01	0.8254
4	28.22	0.8292
9	28.18	0.8284

Table 13. Ablation on P_{off} .

	PSNR	SSIM
0.5	28.03	0.8259
1	28.22	0.8292
2	28.17	0.8283

Table 14. Performance comparison of SR and deraining methods under different scaling factors.

Types	Methods	$\times 4$	$\times 4.5$	$\times 5$	$\times 5.5$	$\times 6$	$\times 6.5$	$\times 7$	$\times 7.5$	$\times 8$
ASSR+Derain	LIIF+DRSformer	17.92	16.76	15.84	15.03	14.47	14.02	13.71	13.58	13.50
	DRSformer+LIIF	20.14	19.87	19.76	19.53	19.39	19.21	19.12	18.98	18.93
	GaussianSR+DRSformer	17.14	16.09	15.25	14.53	14.07	13.75	13.57	13.53	13.52
	DRSformer+GaussianSR	20.12	19.86	19.75	19.52	19.39	19.21	19.13	18.98	18.94
	CiaoSR+DRSformer	19.45	18.25	17.19	16.29	15.55	14.85	14.22	13.94	13.73
	DRSformer+CiaoSR	19.96	19.73	19.65	19.48	19.27	19.15	19.07	18.89	18.87
All in one	LIIF	24.04	23.79	23.53	23.29	23.09	22.90	22.70	22.52	22.39
	GaussianSR	24.04	23.76	23.51	23.28	23.08	22.89	22.69	22.52	22.38
	CiaoSR	23.84	23.66	23.45	22.23	23.01	22.83	22.66	22.50	22.33
	Ours	24.51	24.22	23.95	23.69	23.48	23.28	23.07	22.90	22.76

Table 15. Comparison with GaussianImage [64].

	LIIF	LTE	CiaoSR	GaussianImage	Ours
Set5	32.73	32.81	32.84	28.01	33.24
Set14	28.98	29.06	29.08	25.62	29.40

15. More Exploration and Results

In Section 6 of the main text, we demonstrate the performance of our method in low-resolution and rainy scenarios. Here, we present comparisons across more scaling factors and with the two-stage ASSR+Derain methods DRSformer [9]. The results are shown in Table 14. As shown in Table 14, our method consistently outperforms other methods across all scaling factors. For instance, at the $\times 4$ scale, our method achieves a PSNR of 24.51, significantly higher than the best two-stage method, DRSformer [9]+LIIF (20.14). At the $\times 8$ scale, our method achieves 22.76, outperforming DRSformer+LIIF (18.93). Compared to other "All in one" methods, our approach also achieves superior results, such as 23.95 at the $\times 5$ scale, outperforming both GaussianSR (23.51) and CiaoSR (23.45). These results highlight the robustness, simplicity, and effectiveness of our method for super-resolution and deraining tasks across various scales.

16. Limitation and Future Work

Position Distribution. In this paper, to address the difficulty of optimizing position parameters, we propose Adaptive Position Drifting, which leverages an offset mechanism

Algorithm 1: ContinuousSR.

```

Input: inp : (B, 3, H0, W0)
Output: image : (B, 3, H, W)
F ← Encoder(inp);
C ← CGM(F) ; // [N × 3]
Δx ← APD(F) ; // [N × 2]
Σ ← DDCW(F) ; // [N × 3]
x̃ ← Grid(H, W) + αΔx
for n ∈ [1, N] do
    | (xysn, depthn, radiin, conicn) ←
    | Project(x̃n, Σn);
end
image ← Composite({xys, depth, radii, conic, C});
```

to alleviate the optimization challenges and enhance the representational capacity of the model. However, assigning one or four Gaussian kernels to each LR pixel introduces some limitations. On the one hand, it leads to an overabundance of Gaussian kernels in low-frequency regions, resulting in resource wastage. On the other hand, it increases the optimization difficulty significantly. To address these issues, we plan to explore the adaptive allocation of Gaussian kernels based on the texture complexity of image content in future work. This approach aims to dynamically assign an appropriate number of kernels to different regions, effectively mitigating the aforementioned problems.

Introduce Generation Knowledge. In addition, considering that arbitrary-scale super-resolution sometimes requires large upscaling factors (*e.g.*, $\times 16$, $\times 32$, *etc.*), it is challeng-

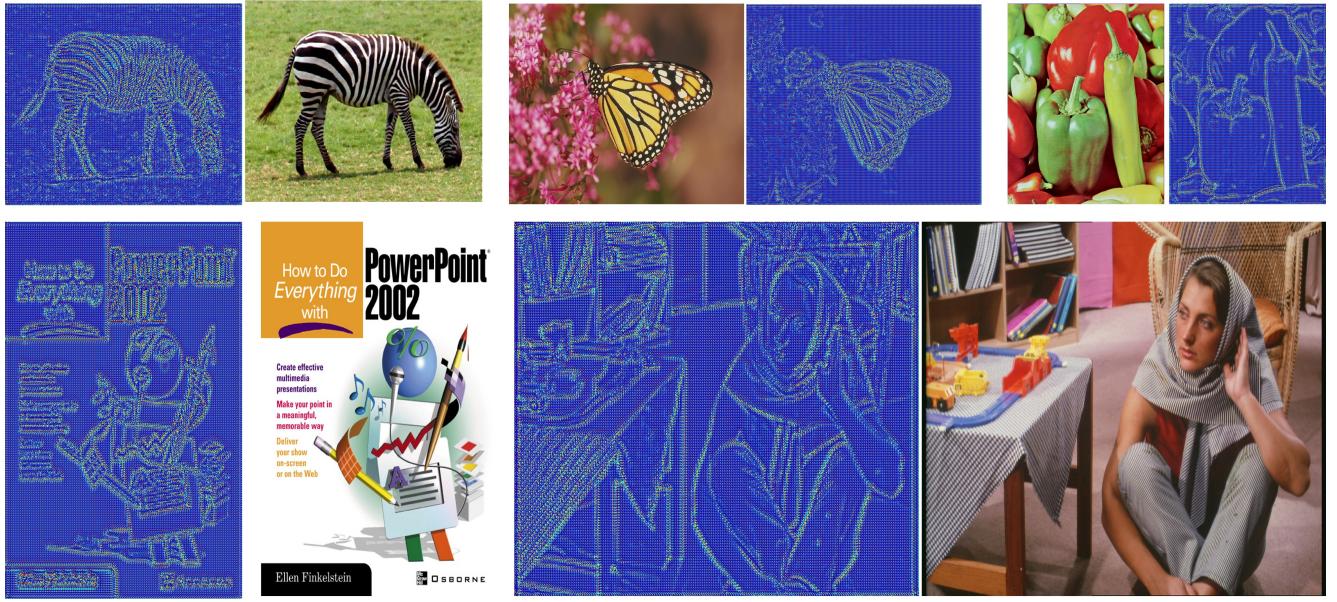


Figure 6. Visualization of the position distribution.

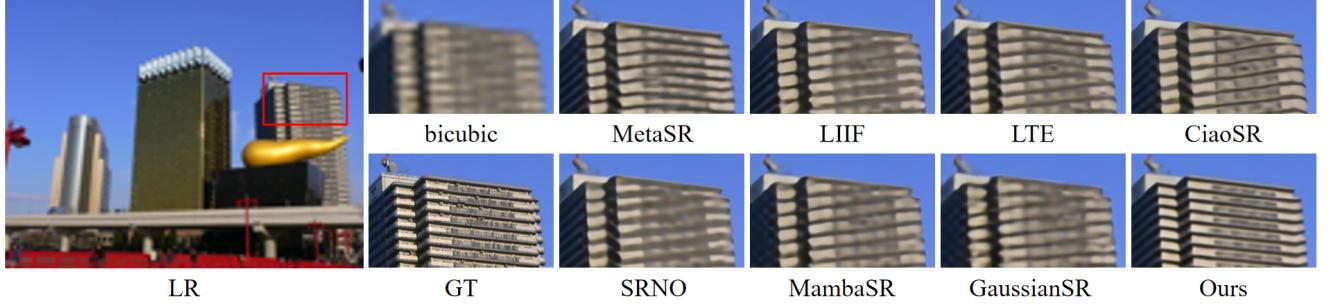


Figure 7. More qualitative comparison. The visual quality of our method outperforms existing methods. Please zoom in for a better view.

ing for the model to generate high-quality details solely relying on the input image and model knowledge. Therefore, in the future, we plan to incorporate more visual knowledge

from diffusion models or semantic knowledge from large vision-language models to help the network generate finer details for high-magnification scenarios.

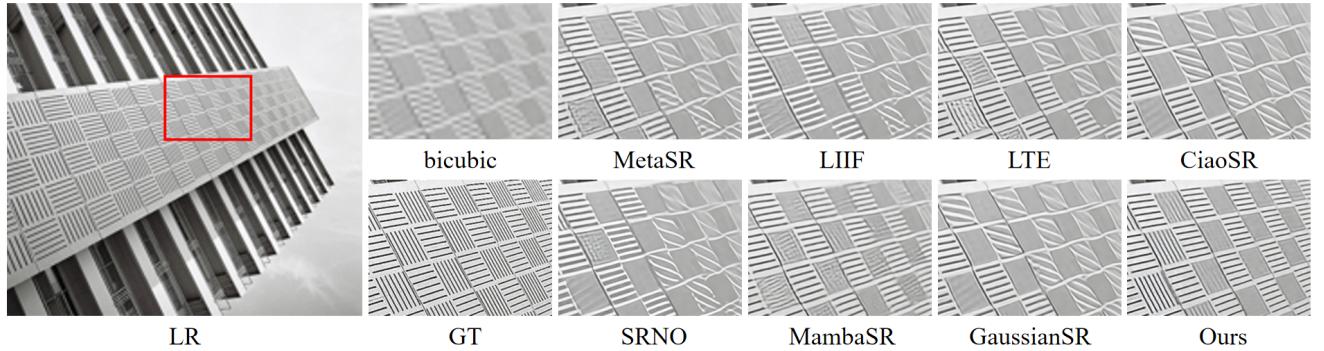


Figure 8. More qualitative comparison. The visual quality of our method outperforms existing methods. Please zoom in for a better view.

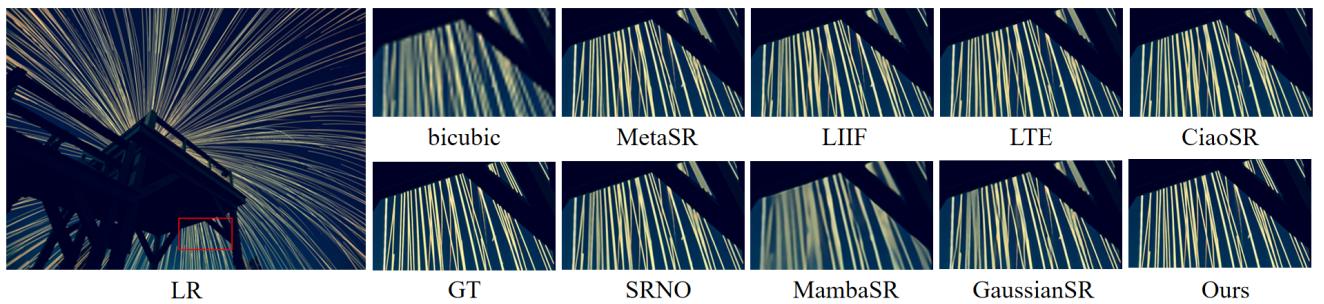


Figure 9. More qualitative comparison. The visual quality of our method outperforms existing methods. Please zoom in for a better view.