UNVEILING HIDDEN DETAILS: A RAW DATA-ENHANCED PARADIGM FOR REAL-WORLD SUPER-RESOLUTION

Anonymous authorsPaper under double-blind review

000

001

002

004

006

008 009 010

011 012

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

033

035

037

040

041

042

043

044

045

046

047

048

051

052

ABSTRACT

Real-world image super-resolution (Real SR) aims to generate high-fidelity, detailrich high-resolution (HR) images from low-resolution (LR) counterparts. Existing Real SR methods primarily focus on generating details from the LR RGB domain, often leading to a lack of richness or fidelity in fine details. In this paper, we pioneer the use of details hidden in RAW data to complement existing RGB-only methods, yielding superior outputs. We argue that key image processing steps in Image Signal Processing, such as denoising and demosaicing, inherently result in the loss of fine details in LR images, making LR RAW a valuable information source. To validate this, we present RealSR-RAW, a comprehensive dataset comprising over 10,000 pairs with LR and HR RGB images, along with corresponding LR RAW, captured across multiple smartphones under varying focal lengths and diverse scenes. Additionally, we propose a simple yet efficient and general RAW adapter to effectively integrate LR RAW data into existing CNNs, Transformers, and Diffusion-based Real SR models by extracting fine-grained details from RAW data to enhance performance. Extensive experiments demonstrate that incorporating RAW data significantly enhances detail recovery and improves Real SR performance across ten evaluation metrics, including both fidelity and perception-oriented metrics, under real-world and wild-captured scenarios. Our findings open a new direction for the Real SR task, with the dataset and code being made available to support future research.

1 Introduction

Real-world image super-resolution (Real SR), a fundamental task in image processing, is designed to enhance the resolution and quality of low-resolution (LR) images (Mou et al., 2022; Liu et al., 2022; Zhou et al., 2020; Wu et al., 2024; Liang et al., 2024; Yang et al., 2023; Chen et al., 2024; Guo et al., 2024; Neshatavar et al., 2024; Jiang et al., 2024a). Numerous studies have developed specialized CNNs, Transformers, and Diffusion models to learn pixel relationships in LR images and generate high-resolution (HR) images with finer details (Chen et al., 2022; Yu et al., 2024; Liu et al., 2023; Zhang et al., 2023b; Sun et al., 2023; Zhang et al., 2024). However, these existing approaches primarily focus on the detail-limited RGB domain. As is well known, SR is an ill-posed problem, making it difficult to recover rich details and high-fidelity results by relying solely on detail-limited LR RGB images (Chen et al., 2023a; Huang et al., 2020; Wang et al., 2021a; Peng et al., 2024a; Luo et al., 2024; Yan et al., 2024; Li et al., 2024), as shown in Figure 1.

During the camera imaging process, photons reflected from physical objects are captured by CMOS or CCD sensors to produce RAW images, which cannot be directly perceived by the human visual system (Blahut, 2010; Prasanna & Rai, 2014). A complex image signal processing (ISP) pipeline, involving a number of operations, is then applied to generate a visually observable RGB image (Pitas, 2000), as illustrated in Figure 2. However, certain modules within the ISP pipeline, such as denoising and demosaicing, inevitably lead to the loss of image details (Ignatov et al., 2020). In Figure 2 (b) and Sec. 3, we visualize the residual images of bypassing denoising and before and after the denoising and demosaicing, assessing information loss with feedback from human users and Multimodal Large Language Models (MLLMs). We can observe that both users and MLLMs agree that in the vast majority of scenarios, both denoising and demosaicing in ISP can lead to a loss of detail. This analysis

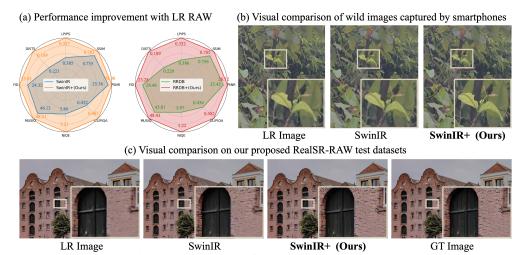


Figure 1: (a) Equipped with LR RAW, the performance of existing RGB-only Real SR models is significantly improved. (b-c) LR RAW also aids Real SR models in generating superior high-fidelity details that are hard to learn in the LR RGB space, thereby significantly enhancing visual quality.

reveals that some fine details are indeed lost during ISP, which exacerbates the challenges of the Real SR task. This raises an important question: Can the LR RAW images, containing rich and original details information, be utilized to assist Real SR in producing more detail-rich and high-fidelity HR images?

Our answer is **absolutely**. We compare three learning objectives: LR RGB \rightarrow HR RGB (i.e., using LR RGB images to generate HR RGB images), LR RAW \rightarrow HR RGB, and LR RGB + RAW → HR RGB, concluding that the latter, where LR RAW complements LR RGB, delivers the best performance. Since existing Real SR datasets lack paired LR RAW and LR and HR RGB images, we introduce RealSR-RAW, a dataset containing over 10,000 image pairs, including LR RAW, paired LR and HR RGB images. Captured using multiple smartphones across diverse scenes and cameras with different focal lengths, this dataset enables a thorough evaluation of LR RAW's effectiveness. We experiment with three representative Real SR models—CNNs, Transformers, and Diffusion-based methods—and show that simply incorporating LR RAW data largely enhances performance. To maximize the benefits of RAW data, we also design a general RAW adapter to integrate LR RAW information seamlessly into these frameworks by adaptively suppressing noise in LR RAW and aligning the distribution of RAW features to RGB. The results are striking: our approach yields up to 1.109 dB and 0.038 improvements in PSNR and SSIM, consistently producing images with richer, more high-fidelity details, as shown in Figure 1. Our proposed dataset and baseline establish a solid foundation for future research, offering valuable resources for the research community to build upon and further advance the state-of-the-art in Real SR. The contributions of this paper can be summarized:

- We introduce RealSR-RAW, the first Real SR dataset containing over 10,000 high-quality paired LR and HR RGB images, along with corresponding LR RAW data.
- For the first time, we explore the effectiveness of LR RAW data as a detail supplement to boosting Real SR models, opening a new avenue for advancements in the field.
- To fully leverage LR RAW data, we propose a novel, general RAW adapter that efficiently captures
 useful information and aligns the distribution of RAW features to the RGB domain, resulting in
 significant improvements across multiple real-world benchmarks and metrics.

2 Related Work

Real-world image super-resolution. Real-world image super-resolution (Real SR) is an ill-posed problem in image processing, aiming to generate detail-rich and visual pleasing high-resolution images from low-resolution scenes (Lugmayr et al., 2020; Li et al., 2022; Lugmayr et al., 2019; Ji et al., 2020; Mou et al., 2022; Liu et al., 2022; Fritsche et al., 2019b; Zhou et al., 2020; Wang et al., 2021a; Chen et al., 2019). Numerous works have meticulously designed various architecture using CNNs (Wang et al., 2018; 2021a), Transformers (Chen et al., 2023b; Liang et al., 2021b),

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133 134

135 136 137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159 160

161

and Diffusion models (Sun et al., 2023; Yue et al., 2024) to enhance SR performance. For instance, the CNN-based RRDB network (Wang et al., 2018) has been widely adopted in many SR architectures (Wang et al., 2021b; Zhang et al., 2021a; Fritsche et al., 2019a). Liang *et al.* were the first to apply the powerful swin transformer to SR, achieving notable performance (Liang et al., 2021b). Yue *et al.* introduced ResShift, which improves efficiency and performance by generating image residuals through diffusion model (Yue et al., 2024). On the other hand, many researchers also proposed to collect or synthesize paired LR and HR RGB images to enhance the generalization ability of Real SR models in real-world scenarios (Wei et al., 2020a; Cai et al., 2019; Peng et al., 2024b; Zhang et al., 2023a; Tang et al., 2022). However, existing Real SR methods mainly focus on RGB images with limited details and suffer from addressing this ill-posed problem, thereby leading to over-smooth and low-fidelity details.

RAW image enhancement. With the rapid development of smartphone and photography technology, numerous works have focused on enhancing directly to original RAW images (Jiang et al., 2024b; Huang et al., 2022; Lu & Jung, 2022; Conde et al., 2022; 2024b; Yue et al., 2022; Liu et al., 2021). For instance, Conde et al. organized a RAW SR competition focused on learning the mapping from LR RAW to HR RAW (Conde et al., 2024c). However, the LR RAW used in the competition was simulated from HR RAW through degradation simulation, which limited the network's ability to generalize to real-world scenarios. Yi et al. proposed using diffusion models to establish the mapping from low-quality RAW images captured by smartphones to high-quality RGB images from DSLRs (Yi et al., 2024). Chen et al. proposed a model to directly reconstruct normal-exposure RGB images from low-light RAW images (Chen et al., 2018). Xu et al. first synthesized LR RAW and RGB images from HR RAW and then proposed to learn color transformation from LR RGB and performed enhancement in the RAW space to generate HR images (Xu et al., 2019). Burst Image Super-Resolution was proposed to generate a high-resolution RGB image directly from a series of LR RAW images captured by burst photography (Bhat et al., 2021). To the best of our knowledge, we are the first to collect real paired images with LR RGB, HR RGB, and LR RAW and explore the benefits of RAW as a detail supplement to boost the representation capability of image details for Real SR.

3 WHY LR RAW DATA CAN BOOST REAL SR?

Image signal processing. In the camera imaging process, photons are captured by CMOS or CCD sensors, which measure light intensity and produce a Bayer RAW image. Since RAW data is in Bayer format and only contains a single color channel per pixel, it cannot be directly interpreted by the human visual system. To convert RAW data into a perceptually meaningful RGB image, a complex Image Signal Processing (ISP) (Prasanna & Rai, 2014; Blahut, 2010) pipeline is applied. While the exact composition of ISP pipelines can vary significantly across different cameras and devices, certain core operations are universally implemented. For example, demosaicing reconstructs full-color RGB images from the mosaic-like pattern of the Bayer filter, while denoising reduces noise introduced by sensor limitations, high ISO levels, and photon shot noise. Color correction is employed to map the device-specific sensor response to a standardized color space, while white balance adjustment further refines this process by compensating for lighting conditions, and neutralizing color casts caused by the ambient light's color temperature. Another essential operation is defective pixel correction, which addresses sensor irregularities by identifying and interpolating faulty pixel data to maintain image consistency. Collectively, these steps play a pivotal role in converting sensor data into high-quality RGB images. However, certain processes within the ISP pipeline, such as denoising (Tian et al., 2020; Fan et al., 2019) and demosaicing (Li et al., 2008; Li, 2005), inevitably result in the loss of fine details in the final RGB image, as discussed in following. This loss poses significant challenges for Real SR methods operating solely in the RGB domain, making it difficult to reconstruct detail-rich and high-fidelity HR images from the degraded LR RGB data.

Detail loss during image signal processing. To avoid copyright concerns related to commercial ISPs, we analyze the problem of detail loss using an open-source available ISP, OpenISP¹, and the widely-used RAW processing library, RawPy², on our collected datasets. Specifically, we employ two analysis methods: bypassing and step-by-step analysis to explore individual ISP modules and compare the resulting images, as well as to analyze the image differences before and after processing through specific modules. Given that image details are mainly characterized as high-frequency signals,

¹https://github.com/cruxopen/openISP

²https://pypi.org/project/rawpy

(a) Can the detail-rich LR RAW assist Real SR in generating better image details? LR RAW Real SR Pipeline Denoise, Demosaicing, Color SR Correction, etc. Models (Details may be lost.) HR RGB **Detail Supplement** (b) Analysis of detail loss in ISP RGB Image (Bypass) ISP MLLMs Users

Figure 2: Existing RealSR methods focus on LR RGB images, as shown in (a). However, LR RGB images often suffer from detail loss due to ISP, as shown in (b), which exacerbates the challenges of RealSR. Therefore, we think: *Can the detail-rich LR RAW information assist Real SR in generating better image details?*

we focus on modules that potentially impact high-frequency information, such as the denoising and demosaicing process, to investigate detail loss during the ISP. More analyses of other modules in ISP are presented in Appendix E. For a comprehensive evaluation, we involve both human volunteers and Multimodal Large Language Models (MLLMs) to assess information loss.

Bypass analysis. Using the RawPy library, we process RAW data both bypassing and non-bypassing the denoising module to generate two RGB images for comparison. As shown in Figure 2, the image bypassing denoising exhibits a certain level of noise but retains more image details. We further compute the residual between the two images, revealing more structural details alongside some noise, as visualized in Figure 2(b) and Figure 9 in the appendix. This suggests that detail loss occurs during denoising. To verify this systematically, we randomly select 100 RAW images from the our dataset and repeat the bypass/non-bypass operations. We present the paired RGB images and residuals to ten volunteers, asking: {USER: Please determine if the residual image on the right contains the structural content information of the image on the left? Answer Yes or No.} Additionally, we utilize the MLLM model LLava1.5 (Liu et al., 2024) to evaluate a larger test set consisting of 1000 images. The results indicate that in 95% of the scenarios, ten volunteers agree that the residuals contain detailed structural information, with LLava corroborating this in 94% of the cases.

Step-by-step analysis. We also perform a step-by-step analysis of the denoising and demosaicing processes using OpenISP to explore potential detail loss. Specifically, we visualize the images before and after the denoising and demosaicing and then visualize the residual images to analyze the difference introduced during this process. As shown in Figure 10 and 11 in Appendix, the results show that these residuals retain substantial structural information. To further assess this, ten volunteers and MLLMs are invited for evaluation as above analysis. The results indicate that in **99%** and **98%** of the scenarios, volunteers recognize detailed structural information in the denoising and demosaicing residuals, respectively, while LLava reaches the same conclusion in **98.2%** and **97.9%** of the scenes.

From the above analysis, it is clear that detail loss occurs throughout the ISP pipeline. As a result, performing Real SR in the LR RGB domain poses significant challenges in recovering detail-rich and high-fidelity images due to the ill-posed nature of the task. To address this, we propose leveraging LR RAW to enhance Real SR and achieve better reconstruction of finer image details.

4 Dataset and Method

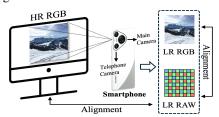
4.1 COLLECTING THE REALSR-RAW DATASET

As shown in Table 1, current real-world super-resolution datasets, such as DIV2K (Agustsson & Timofte, 2017), UHD4K (Zhang et al., 2021b), RealSR (Cai et al., 2019), and DRealSR (Wei et al.,

Table 1: Comparison with existing Real SR data. For the first time, we collect over 10,000 scenes with paired LR RAW, LR RGB, and HR RGB.

Dataset	with HR	with LR	with RAW	Number
DIV2K	✓	Х	Х	800
UHD4K	/	×	×	8,099
RealSR	/	/	X	559
DRealSR	✓	✓	X	2,000
Ours	√	✓	✓	11,726

Figure 3: Illustration of data collection and alignment.



2020a), are limited to RGB images and offer a relatively small number of paired samples, which hampers their diversity and broader applicability. More comparisons between existing datasets are provided in Appdenix D. To unlock the potential of LR RAW data, we present RealSR-RAW, the first large-scale dataset comprising over 10,000 diverse scenes with paired LR RAW, LR RGB, and HR RGB images. Specifically, we first gather high-quality clean images with a resolution of 4K and above in PNG format from the open-source platform Unsplash³. To ensure compliance, we contact Unsplash's official team to receive support and remove any images with potential copyright or ethical concerns. These high-quality images are then displayed on ultra high-definition monitors, where we capture LR RGB and LR RAW images using the main and telephoto cameras of HUAWEI Mate 50 Pro and P70 phones at different focal lengths. The original high-resolution images are used as ground truth HR images. Finally, we apply a two-stage alignment process: first aligning the LR RGB images to their corresponding LR RAW counterparts, and then aligning the HR RGB images to the LR data using estimated homography matrices and optical flow, as shown in Figure 3. We also perform color correction to ensure color-consistent pairs. In total, we collect 11,726 image pairs, which are divided into a training subset and a test benchmark. The resolution of the LR RGB and RAW images range from approximately 1K to 2K, while the HR RGB images range from 2K to 4K, with a scaling factor of 2. More details of data collection are provided in Appendix A. Our dataset will be made open-source to facilitate community research.

4.2 REAL SR WITH LR RAW CONCATENATION

Popular Real SR methods reconstruct HR image, \mathcal{HR}_{RGB} , from a LR RGB image, \mathcal{LR}_{RGB} , using a dedicated SR model. The SR model typically consists of a shallow feature extraction module, L_s , a feature enhancement module, L_e , and a feature-to-image mapping layer, L_f :

$$\mathcal{HR}_{RGB} = L_f \left(L_e \left(\mathcal{LR}_{RGB} \right) \right) . \tag{1}$$

Building on this formulation, a straightforward strategy to introduce RAW images is to concatenate the LR RAW image \mathcal{LR}_{RAW} with the LR RGB image \mathcal{LR}_{RGB} of the same size as the input, which can be expressed as:

$$SR_{RGB} = L_f \left(L_e \left(L_s \left(\mathcal{L}R_{RGB} \| \mathcal{L}R_{RAW} \right) \right) \right). \tag{2}$$

where \parallel is the concatenation operation. We are surprised to find that this simple approach largely improves the performance of the Real SR model, as demonstrated in Section 5.2, highlighting the effectiveness of incorporating RAW data.

4.3 REAL SR WITH RAW ADAPTER

Considering that LR RAW is in Bayer format and contains an amount of noise, directly concatenating LR RAW and RGB images can lead to distribution mismatches and noise interference. To address these issues, we propose a general and efficient RAW adapter that facilitates the fusion of LR RGB and RAW in the feature space, fully leveraging the potential of RAW information. Also, this adapter can be seamlessly integrated into various Real SR models, as illustrated on the left of Figure 4.

In detail, as shown on the right of Figure 4, we first use shallow feature extractors L_s^{RGB} and L_s^{RAW} to process the LR RGB and LR RAW, producing \mathcal{F}_{RGB} and \mathcal{F}_{RAW} :

$$\mathcal{F}_{RGB}, \mathcal{F}_{RAW} = L_s^{RGB} \left(\mathcal{LR}_{RGB} \right), L_s^{RAW} \left(\mathcal{LR}_{RAW} \right). \tag{3}$$

³https://unsplash.com/

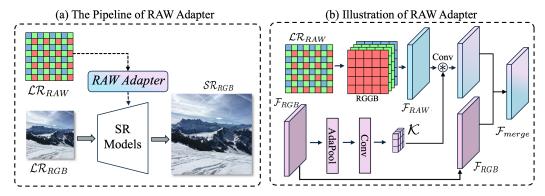


Figure 4: (a) The proposed RAW adapter seamlessly integrates into various popular Real SR methods to boost their representation capability of detail. (b) Illustration of the proposed RAW adapter.

Specifically, L_s^{RAW} unpacks the Bayer format into RGGB channels, applies convolutional blocks to extract features, and utilizes transposed convolution to upsample the resolution, matching it with the RGB features. Next, adaptive kernels are generated from \mathcal{F}_{RGB} to convolve with \mathcal{F}_{RAW} , producing \mathcal{F}_{RAW}' , which is aligned with the RGB features. These adaptive kernels, \mathcal{K} , are obtained by performing adaptive pooling and convolution on \mathcal{F}_{RGB} to perceive the distribution of RGB while modulating learnable kernels, \mathcal{K}_{learn} .

$$\mathcal{K} = Conv\left(AdaPool\left(\mathcal{F}_{RGB}\right)\right) \cdot \mathcal{K}_{learn}.$$
(4)

Finally, we concatenate \mathcal{F}_{RAW}' with \mathcal{F}_{RGB} , followed by a convolution to produce the fused result, \mathcal{F}_{merge} , and carry out the reconstruction HR image $\mathcal{SR}_{RGB} = L_f \left(L_e \left(\mathcal{F}_{merge} \right) \right)$. This process is expressed as:

$$\mathcal{F}_{RAW}^{'} = \mathcal{F}_{RAW} * \mathcal{K}, \mathcal{F}_{merge} = Conv\left(\mathcal{F}_{RAW}^{'} \| \mathcal{F}_{RGB}\right). \tag{5}$$

This design offers two key advantages. First, it adaptively fuses RAW features based on the distribution of individual RGB images, greatly enhancing model flexibility. Second, using noise-free RGB features to generate the kernels improves the extraction of useful details from RAW data while mitigating the influence of noise. As demonstrated in Table 5, the proposed RAW adapter significantly elevates model performance compared with the simple concatenation.

5 EXPERIMENTS AND ANALYSIS

5.1 IMPLEMENTATION

Training details. To evaluate the impact of LR RAW data, we compare the traditional RGB-only LR RGB input with our proposed RAW adapter using LR RGB + LR RAW input for Real SR under consistent experimental settings. All experiments are conducted at the ×2 super-resolution scale using the L1 loss function for training and evaluation unless otherwise specified. Further training details and evaluation on perceptual-oriented GAN loss are provided in Section 5.2 and Appendix F. Real SR models. We conduct experiments on three popular and representative Real SR models, including the CNN-based RRDB network (RRDB) (Wang et al., 2018; 2021b), the transformer-based model SwinIR (Liang et al., 2021a), and the diffusion-based model ResShift (Yue et al., 2024). Metrics. To comprehensively evaluate the quality of generated images, we employ a total of ten widely-used and popular image quality assessment metrics for evaluation, including four reference-based metrics: PSNR↑ (Huynh-Thu & Ghanbari, 2008), SSIM↑ (Wang et al., 2004), LPIPS↓ (Zhang et al., 2018), and DISTS↓ (Ding et al., 2020); and six no-reference metrics: FID↓ (Heusel et al., 2017), MUSIQ↑ (Ke et al., 2021), NIQE↓ (Mittal et al., 2012), CLIP-IQA↑ (Radford et al., 2021), NIMA↑ (Talebi & Milanfar, 2018), and MANIQA↑ (Yang et al., 2022). Note that ↑ and ↓ indicate

5.2 QUANTITATIVE AND QUALITATIVE RESULTS

that higher and lower values respectively represent better image quality.

RealSR-RAW benchmark. To demonstrate the improvements that RAW images can bring to Real SR, we utilize three popular Real SR models and compare the different learning mappings: LR RGB

enhancing visual quality.

Table 2: The model and model+ represent the Real SR model with traditional mapping LR RGB \rightarrow HR RGB, and our proposed RAW adapter (LR RGB + RAW \rightarrow HR RGB), respectively.

Dataset	Models	PSNR↑	SSIM↑	LPIPS↓	DISTS↓	FID↓	MUSIQ↑	NIQE↓	CLIP-IQA↑
	SwinIR	25.367	0.755	0.385	0.221	24.326	46.137	5.890	0.432
	SwinIR+	25.982	0.783	0.337	0.189	23.029	48.217	5.216	0.487
	Gain	0.615	0.028	0.048	0.032	1.297	2.080	0.674	0.055
M50-M	RRDB	25.426	0.756	0.386	0.220	24.462	45.811	5.979	0.434
	RRDB+	26.126	0.785	0.332	0.189	23.287	48.437	5.227	0.482
	Gain	0.700	0.029	0.054	0.031	1.175	2.626	0.752	0.048
	SwinIR	25.181	0.749	0.322	0.206	3.976	39.634	6.156	0.430
	SwinIR+	25.588	0.766	0.296	0.185	3.239	40.143	5.815	0.455
3.450 F	Gain	0.407	0.017	0.026	0.021	0.737	0.509	0.341	0.025
M50-T	RRDB	25.279	0.753	0.316	0.203	4.032	40.047	6.022	0.447
	RRDB+	25.720	0.770	0.290	0.182	3.192	40.665	5.764	0.469
	Gain	0.441	0.017	0.026	0.021	0.840	0.618	0.258	0.022
	SwinIR	24.642	0.776	0.300	0.173	3.263	46.745	4.646	0.508
	SwinIR+	25.744	0.815	0.251	0.145	2.391	49.076	4.339	0.539
D70 14	Gain	1.102	0.039	0.049	0.028	0.872	2.331	0.307	0.031
P70-M	RRDB	24.836	0.781	0.295	0.169	3.221	46.886	4.630	0.520
	RRDB+	25.945	0.819	0.242	0.142	2.387	49.406	4.321	0.556
	Gain	1.109	0.038	0.053	0.027	0.834	2.520	0.309	0.036
	SwinIR	24.753	0.735	0.356	0.220	8.077	38.593	6.305	0.417
	SwinIR+	25.108	0.749	0.334	0.203	5.603	39.412	6.056	0.431
D70 F	Gain	0.355	0.014	0.022	0.017	2.474	0.819	0.249	0.014
P70-T	RRDB	24.829	0.737	0.354	0.220	7.874	39.224	6.269	0.426
	RRDB+	25.185	0.751	0.332	0.204	5.605	39.917	6.029	0.437
	Gain	0.356	0.014	0.022	0.016	2.269	0.693	0.240	0.011

Table 3: Performance comparison of different learning mappings for ResShift on the M50-M dataset.

Methods	PSNR↑	SSIM↑	LPIPS↓	DISTS↓	FID↓	$MUSIQ \!\!\uparrow$	CLIP-IQA↑
ResShift+	24.809	0.732	0.330	0.173	24.310	47.528	0.447
	25.071	0.761	0.312	0.161	23.790	47.882	0.449

 \rightarrow HR RGB, and our proposed RAW adapter (LR RGB + RAW \rightarrow HR RGB). Note that since the official implementation of ResShift only supports the super-resolution factor of \times 4, its performance is also evaluated at this scale. "M50" refers to the Mate 50 Pro phone, while "P70" denotes the Pura 70 phone. "M" and "T" indicate the main and telephoto cameras, respectively. As shown in Table 2 and 3, our method largely surpasses traditional LR RGB \rightarrow HR RGB approach across all benchmarks and metrics. For instance, compared to traditional Real SR, our method improves PSNR by 1.109 dB and LPIPS by 0.053 for the RRDB model on the P70-M dataset. Furthermore, we are surprised that simply inserting the RAW image into model input also achieves considerable gains, as shown in Table 5, demonstrating the significant potential of RAW images for Real SR. As shown in Figure 1 and 5, we can observe that, compared to LR RGB \rightarrow HR RGB, our proposed method assists the RealSR model in extracting more image details from RAW data, generating higher fidelity and detail-rich HR images. More results (e.g., comparison with StableSR) are provided in Appendix B. **Real-world captured LR images.** To evaluate the effectiveness of RAW data in real-world scenarios, we use the main and telephoto cameras of the Mata 50 Pro to capture 192 and 224 pairs of LR RGB and LR RAW images, respectively. RRDB models pre-trained on Mata 50 Pro datasets are applied for evaluation. Since real-world test images in the wild lack ground truth, we employ five widely used no-reference metrics for assessment. As shown in Table 4, incorporating LR RAW significantly improves the Real SR model's performance across all metrics. Figure 1(c) and Figure 13 in the appendix illustrate that our method is also capable of generating HR images with richer textures. **User study.** We conduct a user study using 10 randomly selected real LR images captured by the Mata 50 Pro, evaluating the performance of RRDB and SwinIR. Ten volunteers are invited to rate the quality of the generated images on a scale of 1 to 10. In Figure 6, RRDB+ and SwinIR+, enhanced by our RAW adapter, achieve higher average scores of 7.92 and 7.99, respectively, due to improved

Validation of GAN loss. We also compare Real SR model performance on P70-M when trained using

detail representation from RAW data. These results demonstrate the effectiveness of our method in



Figure 5: Visual comparison of RRDB and SwinIR on our RealSR-RAW dataset.

Devices	Models	MUSIQ↑	NIQE↓	CLIP-IQA↑	NIMA↑	MANIQA↑
	RRDB	31.940	7.867	0.300	3.687	0.253
	RRDB+	34.313	7.850	0.310	3.767	0.260
M	SwinIR SwinIR+	32.088 35.060	7.815 7.677	0.299 0.309	3.707 3.708 3.842	0.252 0.262
	RRDB	47.313	7.346	0.431	4.181	0.284
	RRDB+	48.406	7.314	0.445	4.359	0.290
T	SwinIR	45.965	7.532	0.414	4.250	0.275
	SwinIR+	46.562	7.480	0.420	4.305	0.282

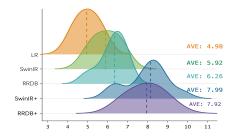


Table 4: Performance improvement of in the real images captured by M50.

Figure 6: User study.

the commonly adopted perceptual-oriented GAN loss. Following Wang et al. (2021b), the total loss function is combined with L1, GAN, and VGG loss function. As shown in Table 6, the RAW adapter still consistently enhances performance across all image quality metrics under perceptual-oriented GAN loss. For example, there is a 1.051 dB and 0.03 improvement in the PSNR and LPIPS metric, confirming that the RAW adapter improves the model's ability to perceive finer details.

Model complexity. To demonstrate the efficiency of our RAW adapter, we compare model parameters, FLOPs, and inference time with different mappings. As shown in Table 5, it is evident that our method achieves noticeable performance improvements with minimal additional computational overhead. Compared to only using LR RGB and directly concatenating LR RAW images, our RAW adapter is capable of better extracting detailed information from RAW images and integrating it into the RGB feature space, with only a slight increase in computational complexity.

5.3 DISCUSSION AND ANALYSIS

Why not LR RAW \rightarrow HR RGB? Considering that LR RAW contains rich information, an intuitive approach might be to directly map LR RAW to HR RGB, using an SR model to generate high-quality HR RGB images from a single LR RAW input. However, generating RGB images from RAW data typically requires complex image processing operations within the ISP, making it difficult for a single SR model to handle the entire LR RAW \rightarrow HR RGB mapping. For example, HR RGB images adhere to a well-defined color space, which is challenging for a model to reproduce without explicit color correction and adjustment. To validate this, we conduct experiments with RRDB on three different mappings on the P70, as shown in Appendix Table 7. As expected, the LR RAW \rightarrow HR RGB results show color shifts due to the lack of color adjustment, leading to lower performance compared to traditional Real SR methods like LR RGB \rightarrow HR RGB. In contrast, our proposed RAW adapter effectively extracts detailed information from RAW images to enhance RealSR performance in the RGB space. Additionally, the lack of large-scale, high-quality LR RAW \rightarrow HR RGB datasets may have hindered further exploration in this area. Nonetheless, this approach holds significant research

Table 5: Performance and computational complexity of RRDB with different mapping. The input size is $3\times224\times224$. LR RGB + RAW \rightarrow HR represent our RAW adapter.

Mapping	Param	FLOPs	Time	PSNR↑	LPIPS↓
LR RGB \rightarrow HR LR RGB \parallel LR RAW \rightarrow HR LR RGB + LR RAW \rightarrow HR	9.58M	482.9G		25.913	0.386 0.339 0.332

and practical potential, which we plan to investigate in the future.

Why not LR RAW \rightarrow HR RAW \rightarrow ISP? We identify three main challenges with this pipeline: (a) It is difficult for a single SR model to learn clean mappings from noisy LR RAW inputs. As a result, any remaining noise or artifacts introduced by the SR model will be amplified during ISP, leading to degraded image quality. (b) The RAW space lacks sufficient image/model priors, such as those available in the Stable Diffusion models (Rombach et al., 2022), making it harder to design a powerful RAW-based model. In contrast, the RGB space can leverage these priors for better reconstruction, which is a key insight behind our approach that combines the strengths of both RGB and RAW spaces. (c) Increasing image resolution in the RAW space before ISP significantly raises the computational load of ISP, making this pipeline impractical for edge devices like smartphones and cameras. More detailed analyses and comparisons of existing methods, such as RAW SR and RGB-only RealSR, will be presented in the Appendix A. Improvement gaps between main and telephoto cameras. As shown in Table 2, different SR models exhibit varying performance improvements between the main and telephoto cameras on P70 and Mate 50 Pro phones. For instance, in the RRDB model with our proposed RAW adapter, the PSNR improvement on the P70's main camera is 1.109 dB, while it is only 0.356 dB on the telephoto camera. These discrepancies may arise from differences in sensor quality. Manufacturers typically prioritize enhancing the main camera, as it is the most frequently used, resulting in a higher-quality sensor that captures more detailed information in RAW images. Consequently, our RAW adapter is more effective at extracting detail from the main camera, leading to greater performance gains in Real SR models.

Cross-Lens generalization ability. We conduct cross-lens experiments to evaluate the generalization capability of the RAW adapter under different lens conditions. Specifically, we perform cross-lens tests between the main cameras of the M50 and P70 smartphones, using RRDB pretrained on the M50 and P70 to evaluate the test sets of P70 and M50, respectively. As shown in Appendix Table 8, the RAW adapter improves performance in both M50→P70 and P70→M50 scenarios. In the P70→M50 test, our RAW adapter boosts PSNR by 0.691 dB, SSIM by 0.028, and LPIPS by 0.106. These results demonstrate that the proposed RAW adapter exhibits strong generalization across different lenses.

To thoroughly examine the effectiveness of our proposed RAW adapter, we provide additional discussions, comparisons, in-depth analysis, and extensive visual comparisons in the Appendix. These supplementary materials offer further insights into performance improvements across various scenarios and highlight the adapter's robustness.

6 Conclusion

In this paper, we explore the potential of leveraging LR RAW data as a detailed supplement to unveiling hidden details to enhance real-world super-resolution, overcoming the limitations of traditional RGB-only Real SR methods. We also introduce the RealSR-RAW dataset for community research, consisting of over 10,000 high-quality paired images, including LR RGB, HR RGB, and LR RAW data. Furthermore, we propose a novel RAW adapter that adaptively suppresses noise in RAW data and aligns RAW features with the RGB domain, improving the detail recovery of various existing Real SR models and producing high-fidelity, detail-rich HR images. Extensive experiments demonstrate that our RAW adapter significantly enhances the visual quality of current Real SR methods across all metrics. We hope that our dataset and findings will open new avenues for Real SR research.

In the future, we aim to design more advanced SR models to fully harness the detailed information in RAW data and integrate it with RGB for improving Real SR and other low-level vision tasks. Additionally, we plan to expand our datasets by collecting more RAW data from a wider range of devices, enhancing both the quality and quantity of the data. Furthermore, since other metadata within the ISP is available during camera deployment, we believe that utilizing this information alongside RAW images presents a promising opportunity for further improving the quality of generated images.

REFERENCES

- Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, July 2017. 4
- Goutam Bhat, Martin Danelljan, Luc Van Gool, and Radu Timofte. Deep burst super-resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9209–9218, 2021. 3
- Richard E Blahut. Fast algorithms for signal processing. Cambridge University Press, 2010. 1, 3
- Jianrui Cai, Hui Zeng, Hongwei Yong, Zisheng Cao, and Lei Zhang. Toward real-world single image super-resolution: A new benchmark and a new model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3086–3095, 2019. 3, 4, 16
- Chang Chen, Zhiwei Xiong, Xinmei Tian, Zheng-Jun Zha, and Feng Wu. Camera lens super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1652–1660, 2019. 2
- Chen Chen, Qifeng Chen, Jia Xu, and Vladlen Koltun. Learning to see in the dark. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3291–3300, 2018. 3
- Du Chen, Jie Liang, Xindong Zhang, Ming Liu, Hui Zeng, and Lei Zhang. Human guided ground-truth generation for realistic image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14082–14091, 2023a. 1
- Haoyu Chen, Wenbo Li, Jinjin Gu, Jingjing Ren, Sixiang Chen, Tian Ye, Renjing Pei, Kaiwen Zhou, Fenglong Song, and Lei Zhu. Restoreagent: Autonomous image restoration agent via multimodal large language models. *arXiv preprint arXiv:2407.18035*, 2024. 1
- Honggang Chen, Xiaohai He, Linbo Qing, Yuanyuan Wu, Chao Ren, Ray E Sheriff, and Ce Zhu. Real-world single image super-resolution: A brief review. *Information Fusion*, 79:124–145, 2022.
- Xiangyu Chen, Xintao Wang, Jiantao Zhou, Yu Qiao, and Chao Dong. Activating more pixels in image super-resolution transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22367–22377, 2023b. 2
- Marcos V Conde, Steven McDonagh, Matteo Maggioni, Ales Leonardis, and Eduardo Pérez-Pellitero. Model-based image signal processors via learnable dictionaries. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 481–489, 2022. 3
- Marcos V Conde, Florin Vasluianu, and Radu Timofte. Bsraw: Improving blind raw image superresolution. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 8500–8510, 2024a. 16
- Marcos V Conde, Florin Vasluianu, and Radu Timofte. Toward efficient deep blind raw image restoration. In 2024 IEEE International Conference on Image Processing (ICIP), pp. 1725–1731. IEEE, 2024b. 3
- Marcos V Conde, Florin-Alexandru Vasluianu, Radu Timofte, Jianxing Zhang, Jia Li, Fan Wang, Xiaopeng Li, Zikun Liu, Hyunhee Park, Sejun Song, et al. Deep raw image super-resolution. a ntire 2024 challenge survey. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6745–6759, 2024c. 3
- 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. 6
- Linwei Fan, Fan Zhang, Hui Fan, and Caiming Zhang. Brief review of image denoising techniques. *Visual Computing for Industry, Biomedicine, and Art*, 2(1):7, 2019. 3

- Manuel Fritsche, Shuhang Gu, and Radu Timofte. Frequency separation for real-world super-resolution. In 2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW), pp. 3599–3608. IEEE, 2019a. 3
 - Manuel Fritsche, Shuhang Gu, and Radu Timofte. Frequency separation for real-world superresolution. In 2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW), pp. 3599–3608. IEEE, 2019b. 2
 - Xingbei Guo, Ziping Ma, Qing Wang, and Pengxu Wei. Towards real-world continuous superresolution: Benchmark and method. In 2024 IEEE International Conference on Multimedia and Expo (ICME), pp. 1–6. IEEE, 2024. 1
 - 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 equilibrium. *Advances in neural information processing systems*, 30, 2017. 6
 - Haofeng Huang, Wenhan Yang, Yueyu Hu, Jiaying Liu, and Ling-Yu Duan. Towards low light enhancement with raw images. *IEEE Transactions on Image Processing*, 31:1391–1405, 2022. 3
 - Yan Huang, Shang Li, Liang Wang, Tieniu Tan, et al. Unfolding the alternating optimization for blind super resolution. *Advances in Neural Information Processing Systems*, 33:5632–5643, 2020. 1
 - Quan Huynh-Thu and Mohammed Ghanbari. Scope of validity of psnr in image/video quality assessment. *Electronics letters*, 44(13):800–801, 2008. 6
 - Andrey Ignatov, Luc Van Gool, and Radu Timofte. Replacing mobile camera isp with a single deep learning model. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pp. 536–537, 2020. 1, 18
 - Xiaozhong Ji, Yun Cao, Ying Tai, Chengjie Wang, Jilin Li, and Feiyue Huang. Real-world super-resolution via kernel estimation and noise injection. In *proceedings of the IEEE/CVF conference* on computer vision and pattern recognition workshops, pp. 466–467, 2020. 2
 - Aiwen Jiang, Zhi Wei, Long Peng, Feiqiang Liu, Wenbo Li, and Mingwen Wang. Dalpsr: Leverage degradation-aligned language prompt for real-world image super-resolution. *arXiv preprint arXiv:2406.16477*, 2024a. 1
 - Siyuan Jiang, Senyan Xu, and Xingfu Wang. Rbsformer: Enhanced transformer network for raw image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6479–6488, 2024b. 3
 - Junjie Ke, Qifei Wang, Yilin Wang, Peyman Milanfar, and Feng Yang. Musiq: Multi-scale image quality transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 5148–5157, 2021. 6
 - Bingchen Li, Xin Li, Hanxin Zhu, Yeying Jin, Ruoyu Feng, Zhizheng Zhang, and Zhibo Chen. Sed: Semantic-aware discriminator for image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 25784–25795, 2024. 1
 - Wenbo Li, Kun Zhou, Lu Qi, Liying Lu, and Jiangbo Lu. Best-buddy gans for highly detailed image super-resolution. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 1412–1420, 2022. 2
 - Xin Li. Demosaicing by successive approximation. *IEEE Transactions on Image Processing*, 14(3): 370–379, 2005. 3
 - Xin Li, Bahadir Gunturk, and Lei Zhang. Image demosaicing: A systematic survey. In *Visual Communications and Image Processing 2008*, volume 6822, pp. 489–503. SPIE, 2008. 3
 - Jie Liang, Radu Timofte, Qiaosi Yi, Shuaizheng Liu, Lingchen Sun, Rongyuan Wu, Xindong Zhang, Hui Zeng, Lei Zhang, Yibin Huang, et al. Ntire 2024 restore any image model (raim) in the wild challenge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6632–6640, 2024. 1

- Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 1833–1844, 2021a. 6
 - Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 1833–1844, 2021b. 2, 3
 - 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*, 2022. 1, 2
 - Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 26296–26306, 2024. 4
 - Xiaohong Liu, Kangdi Shi, Zhe Wang, and Jun Chen. Exploit camera raw data for video super-resolution via hidden markov model inference. *IEEE Transactions on Image Processing*, 30: 2127–2140, 2021. 3
 - Yihao Liu, Hengyuan Zhao, Jinjin Gu, Yu Qiao, and Chao Dong. Evaluating the generalization ability of super-resolution networks. *IEEE Transactions on pattern analysis and machine intelligence*, 2023. 1
 - Yucheng Lu and Seung-Won Jung. Progressive joint low-light enhancement and noise removal for raw images. *IEEE Transactions on Image Processing*, 31:2390–2404, 2022. 3
 - Andreas Lugmayr, Martin Danelljan, and Radu Timofte. Unsupervised learning for real-world super-resolution. In 2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW), pp. 3408–3416. IEEE, 2019. 2
 - Andreas Lugmayr, Martin Danelljan, and Radu Timofte. Ntire 2020 challenge on real-world image super-resolution: Methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pp. 494–495, 2020. 2
 - Ziwei Luo, Fredrik K Gustafsson, Zheng Zhao, Jens Sjölund, and Thomas B Schön. Photo-realistic image restoration in the wild with controlled vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6641–6651, 2024. 1
 - Anish Mittal, Rajiv Soundararajan, and Alan C Bovik. Making a "completely blind" image quality analyzer. *IEEE Signal processing letters*, 20(3):209–212, 2012. 6
 - Chong Mou, Yanze Wu, Xintao Wang, Chao Dong, Jian Zhang, and Ying Shan. Metric learning based interactive modulation for real-world super-resolution. In *European Conference on Computer Vision*, pp. 723–740. Springer, 2022. 1, 2
 - Reyhaneh Neshatavar, Mohsen Yavartanoo, Sanghyun Son, and Kyoung Mu Lee. Icf-srsr: Invertible scale-conditional function for self-supervised real-world single image super-resolution. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1557–1567, 2024.
 - 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*, 2024a. 1
 - 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*, 2024b. 3
 - I Pitas. Digital image processing algorithms and applications. *John Wiley & Sons Inc google schola*, 2:133–138, 2000. 1
 - Mahesh K Prasanna and Shantharama C Rai. Image processing algorithms-a comprehensive study. *International Journal of Advanced Computer Research*, 4(2):532, 2014. 1, 3

- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021. 6
 - Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. Highresolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022. 9
 - Haoze Sun, Wenbo Li, Jianzhuang Liu, Haoyu Chen, Renjing Pei, Xueyi Zou, Youliang Yan, and Yujiu Yang. Coser: Bridging image and language for cognitive super-resolution. *arXiv preprint arXiv:2311.16512*, 2023. 1, 3
 - Hossein Talebi and Peyman Milanfar. Nima: Neural image assessment. *IEEE Transactions on Image Processing*, 27(8):3998–4011, 2018. doi: 10.1109/TIP.2018.2831899. 6
 - Chengzhou Tang, Yuqiang Yang, Bing Zeng, Ping Tan, and Shuaicheng Liu. Learning to zoom inside camera imaging pipeline. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 17552–17561, 2022. 3, 18
 - Chunwei Tian, Lunke Fei, Wenxian Zheng, Yong Xu, Wangmeng Zuo, and Chia-Wen Lin. Deep learning on image denoising: An overview. *Neural Networks*, 131:251–275, 2020. 3
 - Longguang Wang, Yingqian Wang, Xiaoyu Dong, Qingyu Xu, Jungang Yang, Wei An, and Yulan Guo. Unsupervised degradation representation learning for blind super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10581–10590, 2021a. 1, 2
 - 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 *Proceedings of the European conference on computer vision (ECCV) workshops*, pp. 0–0, 2018. 2, 3, 6
 - 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*, pp. 1905–1914, 2021b. 3, 6, 8, 16, 18
 - 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. 6
 - 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*, pp. 101–117. Springer, 2020a. 3, 4, 16
 - 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, pp. 101–117. Springer, 2020b. 15
 - Rongyuan Wu, Tao Yang, Lingchen Sun, Zhengqiang Zhang, Shuai Li, and Lei Zhang. Seesr: Towards semantics-aware real-world image super-resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 25456–25467, 2024. 1
 - Xiangyu Xu, Yongrui Ma, and Wenxiu Sun. Towards real scene super-resolution with raw images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1723–1731, 2019. 3, 17
 - Qingsen Yan, Axi Niu, Chaoqun Wang, Wei Dong, Marcin Woźniak, and Yanning Zhang. Kgsr: A kernel guided network for real-world blind super-resolution. *Pattern Recognition*, 147:110095, 2024.

- Sidi Yang, Tianhe Wu, Shuwei Shi, Shanshan Lao, Yuan Gong, Mingdeng Cao, Jiahao Wang, and Yujiu Yang. Maniqa: Multi-dimension attention network for no-reference image quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1191–1200, 2022. 6
- Tao Yang, Rongyuan Wu, Peiran Ren, Xuansong Xie, and Lei Zhang. Pixel-aware stable diffusion for realistic image super-resolution and personalized stylization. *arXiv preprint arXiv:2308.14469*, 2023. 1
- Mingxin Yi, Kai Zhang, Pei Liu, Tanli Zuo, and Jingduo Tian. Diffraw: Leveraging diffusion model to generate dslr-comparable perceptual quality srgb from smartphone raw images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 6711–6719, 2024. 3
- Fanghua Yu, Jinjin Gu, Zheyuan Li, Jinfan Hu, Xiangtao Kong, Xintao Wang, Jingwen He, Yu Qiao, and Chao Dong. Scaling up to excellence: Practicing model scaling for photo-realistic image restoration in the wild. *arXiv preprint arXiv:2401.13627*, 2024. 1
- Huanjing Yue, Zhiming Zhang, and Jingyu Yang. Real-rawvsr: Real-world raw video super-resolution with a benchmark dataset. In *European Conference on Computer Vision*, pp. 608–624. Springer, 2022. 3
- Zongsheng Yue, Jianyi Wang, and Chen Change Loy. Resshift: Efficient diffusion model for image super-resolution by residual shifting. *Advances in Neural Information Processing Systems*, 36, 2024. 3, 6
- Kai Zhang, Jingyun Liang, Luc Van Gool, and Radu Timofte. Designing a practical degradation model for deep blind image super-resolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4791–4800, 2021a. 3
- Kaihao Zhang, Dongxu Li, Wenhan Luo, Wenqi Ren, Bjorn Stenger, Wei Liu, Hongdong Li, and Yang Ming-Hsuan. Benchmarking ultra-high-definition image super-resolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021b. 4
- Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 586–595, 2018. 6
- Ruofan Zhang, Jinjin Gu, Haoyu Chen, Chao Dong, Yulun Zhang, and Wenming Yang. Crafting training degradation distribution for the accuracy-generalization trade-off in real-world super-resolution. In *International Conference on Machine Learning*, pp. 41078–41091. PMLR, 2023a.
- Wenlong Zhang, Xiaohui Li, Xiangyu Chen, Xiaoyun Zhang, Yu Qiao, Xiao-Ming Wu, and Chao Dong. Seal: A framework for systematic evaluation of real-world super-resolution. In *The Twelfth International Conference on Learning Representations*, 2023b. 1
- Wenlong Zhang, Xiaohui Li, Guangyuan Shi, Xiangyu Chen, Yu Qiao, Xiaoyun Zhang, Xiao-Ming Wu, and Chao Dong. Real-world image super-resolution as multi-task learning. *Advances in Neural Information Processing Systems*, 36, 2024. 1
- Yuanbo Zhou, Wei Deng, Tong Tong, and Qinquan Gao. Guided frequency separation network for real-world super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pp. 428–429, 2020. 1, 2

A More details and analysis of data collection

A.1 DETAILS OF DATA COLLECTION PROCESS

Image Capturing. During the data collection process, we use different smartphones to capture HR images displayed on the UHD monitor. First, we switch to "Pro" mode in the smartphone to obtain LR RAW and LR PNG data. We then select the main and telephoto cameras for capturing, while adjusting the position of the smartphone to ensure that the captured images include as much of the display as possible without capturing any background information. Additionally, we secure the camera tripod and adjust the camera level to maintain horizontal alignment as much as possible. An intelligent smartphone script facilitates the automated capturing of images and switching of images on the display. The total data collection process spans two weeks. Furthermore, to avoid moiré artifacts, we use professional high-resolution displays with high pixel density (see Fig. 5 of the main text). At the same time, multiple professional image processing experts have thoroughly examined our results to ensure that there are no moiré patterns.

Image Alignment. Upon completing the collection process, we obtain several sets of LR RAW, LR PNG images, and the HR PNGs displayed on the screen. These three types of images are initially misaligned, so to precisely align them, we employ a two-step alignment method, as follows:

- 1. The first step involves aligning the LR RAW and LR PNG. Due to the smart stabilization technologies in commercial smartphones, LR PNG is often cropped from LR RAW, resulting in some misalignment, with LR RAW covering a slightly larger scene area. Specifically, we first process the LR RAW using a traditional ISP (as provided in RawPy's ISP) without cropping to obtain a set of LR PNG'. At this point, we use the alignment algorithm in the cv2 library to align LR PNG to LR PNG', resulting in a new LR PNG, denoted as LR PNG*. Now, LR PNG* and LR RAW are completely aligned in spatial positioning, allowing us to proceed to the next alignment step.
- 2. In the second step, we primarily align the LR PNG* with the HR PNG. Specifically, we employ homography matrices and optical flow for spatial alignment. Additionally, we perform color correction to ensure that the colors between LR PNG* and HR PNG are aligned. During this alignment process, the LR PNG* is spatially aligned with the HR PNG, and at this point, LR PNG* is cropped. It is important to note that we record the cropping coordinates at this stage. Once aligned, LR PNG* and HR PNG form a fully aligned image pair, which we denote as the LR PNG, HR PNG image pair.
- 3. During the second step, we record the cropping coordinate information *p* for LR PNG*. Given that LR PNG* and LR RAW are spatially aligned from the first step, we simply introduce the coordinate information *p* into LR RAW and perform the cropping to achieve complete alignment between LR PNG* and LR RAW. At this point, the LR PNG, LR RAW, HR PNG images are fully aligned.

In the end, multiple professional image processing experts have thoroughly examined our results to ensure that there are no spatial and color misalignments.

Details of used monitor. For data collection, we used the EIZO ColorEdge CG319X monitor. The detailed specifications are as follows:

- Measured contrast ratio: 1500:1
- Bit depth: Supports 10-bit simultaneous display via a 24-bit LUT, rendering 1.07 billion colors
- Color gamut: Covers 99% Adobe RGB and 98% DCI-P3
- Tone curve support: PQ (Perceptual Quantization) and HLG (Hybrid Log-Gamma) for HDR workflows
- Sub-pixel layout: Standard RGB stripe arrangement on an IPS panel

These specifications ensure high fidelity in the displayed HR reference images and reproducibility of the dataset collection process.

Challenges of RealSR-RAW Dataset. We ensure precise alignment between LR and HR images, verified by multiple image processing experts. However, smartphone-captured images typically suffer from greater detail loss compared to those captured by DSLRs Wei et al. (2020b), making our dataset

Table 6: Performance comparison of RRDB under different mappings, trained using GAN loss.

Models	PSNR↑	SSIM↑	LPIPS↓	DISTS↓	MUSIQ↑
	23.002 24.053		0.188 0.158	0.102 0.087	57.376 57.868

Table 7: Performance comparison of the RRDB backbone across different learning mappings.

	P70)-M	P70)-T
Mapping	PSNR↑	SSIM↑	PSNR↑	SSIM↑
LR RGB→HR LR RAW→HR	24.833 22.938	0.777 0.739	24.826 23.205	0.734 0.709
$LR~RAW + RGB \rightarrow HR$	25.960	0.819	25.192	0.751

more challenging for super-resolution tasks. Specifically, the PSNR of our LR images is 23.56 dB, significantly lower than the 32.67 dB of DSLR-captured images in the DRealSR dataset.

Number of Collected Datasets. We provide a detailed overview of the collected datasets, including specific quantities of training and testing data. As shown in Table 9, we collected a total of 11,726 paired samples. The training and testing data for different focal lengths using Mate 50 Pro and P70 are also detailed in Table 9.

A.2 DICUSSIONS AND ANALYSIS OF DATA COLLECTION PROCESS

Why not use different focal lengths of smartphones for capturing paired images? Considering that current smartphones already incorporate super-resolution processes, and the telephoto lenses on smartphones often have lower quality, the images captured at the telephoto end tend to suffer from blurriness and lack of detail. Therefore, they are not suitable to serve as HR images.

Why not use smartphones to capture LR and high-quality cameras for HR images? Existing data acquisition methods often employ high-quality cameras to capture paired images at different focal lengths (Cai et al., 2019; Wei et al., 2020a). However, the major drawback of this approach is the potential for significant discrepancies due to the time lag between the two captures, which restricts the relatively static object and greatly limits its diversity. Additionally, this method incurs substantial manual effort, resulting in datasets with only a few thousand or even just a few hundred pairs, which is insufficient in the era of deep learning. Similarly, using high-quality cameras to capture HR images also confines us to photograph static objects, making the method less applicable for broader use cases. Therefore, we adopt a new pipeline to capture a large paired dataset with diverse scenes.

Comparison with the BSRAW dataset (Conde et al., 2024a). Conde et al. proposed a method in (Conde et al., 2024a) to generate LR RAW by simulating degradation from HR RAW. While this method allows for the creation of a large amount of paired data, it also leads to insufficient generalization capability for real-world scenarios. In contrast, our dataset not only ensures authenticity but also emphasizes that LR RAW information can serve as a valuable supplement for enhancing detail representation in RealSR networks. Therefore, genuine LR PNG and HR PNG paired data are also crucial, which BSRAW cannot provide. Although RGB images can be obtained from RAW through open-source ISPs, commercial ISPs used in real-world scenarios remain largely inaccessible. These commercial ISPs are what users typically interact with and hold more practical value. Thus, our data acquisition method offers a perspective on collecting data under commercial ISPs.

Comparison with Real-ESRGAN (Wang et al., 2021b). Although Real-ESRGAN (Wang et al., 2021b) introduces a degradation modeling approach similar to ISP simulation, it can only synthesize paired LR and HR images in the RGB domain and cannot synthesize LR RAW data. Additionally, networks trained on synthetic data often lack generalization capabilities for real-world scenarios. Specifically, we tested the RRDB network trained on Real-ESRGAN (using the official open-source code and pretrained network) in our real-world scenarios, as shown in Table 10. We can see that the RRDB network trained on synthetic data performs poorly in real-world scenarios due to insufficient

Table 8: Cross-lens generalization performance. The M50 \rightarrow P70 indicates the generalization performance on the main camera of P70 test data using the pre-trained model of the M50. It can be observed that the proposed RAW adapter, still significantly outperforms traditional LR RGB \rightarrow HR RGB in cross-lens scenarios.

Cross-Lens	Model	PSNR↑	SSIM↑	LPIPS↓	DISTS↓	MUSIQ↑	NIQE↓
M50→P70	RRDB	22.654	0.685	0.450	0.213	30.892	6.236
	RRDB +	22.825	0.705	0.419	0.196	35.123	5.477
P70→M50	RRDB	24.512	0.721	0.466	0.261	37.724	7.001
	RRDB +	25.203	0.749	0.360	0.206	52.395	5.014

Table 9: Details of our collected data. Number of training and testing samples used in this study.

Concertabono	Mat	e 50 Pro		P70	Total
Smartphone	Main Camera	Telephoto Camera	Main Camera	Telephoto Camera	Total
Train	2,600	2,800	2,694	2,800	10,894
Test	220	202	218	192	832

realism in its degradation model. However, its performance improves significantly when trained with our proposed real-world dataset. Finally, equipping the RRDB with the RAW adapter will lead to substantial improvements.

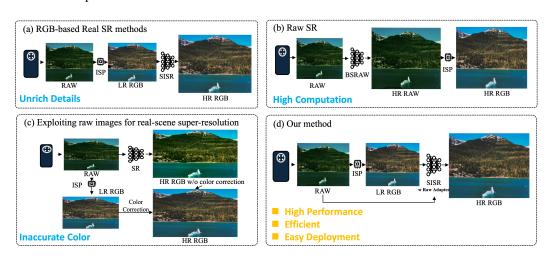


Figure 7: Comprasion with existing methods.

B MORE COMPARISONS AND ANALYSES OF EXISTING METHODS

In this section, we detail the differences between our method and existing approaches, as illustrated in Figure 7. We discuss three existing methods: (a) traditional RGB-based RealSR, (b) RAW SR methods represented by BSRAW, and (c) current methods combining RAW and RGB (Xu et al., 2019). Specifically, the main shortcomings of method (a) has been described in detail in the main text. Due to information loss within the ISP, this method often struggles to reconstruct high-quality, high-fidelity details. Method (b) performs SR in the RAW space; however, the high noise level in RAW space makes SR challenging and the high-resolution RAW images impose significant power consumption in subsequent ISP processes, which is impractical in the real world. The shortcomings of this method are mainly twofold. On the data side, it synthesizes LR RAW and RGB data from HR RAW, which reduces its ability to generalize to real-world scenarios. In contrast, our data is captured directly from real environments. Additionally, the method attempts to directly recover HR

Table 10: Performance of RRDB using different datasets and training methods. Here, RRDB* indicates training with the official Real-ESRGAN (Wang et al., 2021b) synthetic data, RRDB represents training with our proposed real-world data, and RRDB+ shows the results with our proposed RAW adapter.

Datasets	Model	PSNR↑	SSIM↑	LPIPS↓	FID↓
P70-M	RRDB* (Wang et al., 2021b) RRDB RRDB+	22.274 24.836 25.945	0.715 0.781 0.819	0.261 0.295 0.242	10.258 3.221 2.387
P70-T	RRDB* (Wang et al., 2021b) RRDB RRDB+	23.780 24.829 25.185	0.693 0.737 0.751	0.379 0.354 0.332	15.113 7.874 5.605

RGB from LR RAW, which is highly challenging, and using a 3x3 color matrix complicates effective color correction. In summary, our proposed method offers a new perspective for RealSR, providing a more powerful, high-performance, deployable, and efficient approach.

Comparison with StableSR. In addition to RealESRGAN, we compare StableSR on the P70 dataset. Both RealESRGAN and StableSR rely on simulated degradation, which limits generalization to real captured scenes. As shown in Table 11, RRDB trained on our dataset adapts better to real scenes, and equipping it with the RAW Adapter achieves the best results. This further highlights the advantage of leveraging RAW data.

Table 11: Performance comparison on P70-M.

Model	PSNR↑	SSIM↑	LPIPS↓	FID↓
StableSR	21.964	0.677	0.221	7.56
RRDB (RealESRGAN)	22.274	0.715	0.261	10.258
RRDB (w/o RAW Adapter)	24.836	0.781	0.295	3.221
RRDB (w/ RAW Adapter)	25.945	0.819	0.242	2.387

Comparison with Zoom-to-Learn Dataset. The Zoom-to-Learn dataset Tang et al. (2022) constructs LR-HR pairs by capturing scenes at different focal lengths using a zoom lens, resulting in approximately 500 real-world pairs. While this approach provides real data, it has two key limitations.

- (a) The shooting and collection process is highly time-consuming and labor-intensive, which restricts the dataset to only 500 pairs, making it insufficient for training modern deep learning models. In contrast, our pipeline enables efficient collection of over 10,000 pairs using professional-grade displays and smartphones with a rigorous alignment procedure.
- (b) The captured content in Zoom-to-Learn is largely static to ensure alignment, limiting scene diversity. Our method, however, allows for the inclusion of dynamic objects such as humans and animals, thereby enhancing dataset diversity.

As shown in Table 4 of the main text, the proposed dataset improves performance across multiple image quality metrics on images captured with both the main and telephoto cameras of the Mate 50 Pro. Figure 1(b) further demonstrates noticeable enhancements in detail restoration, validating the real-world generalization capability of our dataset. In addition, recent work on ISP-free pipelines Ignatov et al. (2020) also highlights the importance of detail loss introduced by ISP processing, further supporting our motivation to leverage RAW data for RealSR.

C MORE ANALYSES OF THE PROPOSED RAW ADAPTER

Kernel Adaptivity and Effectiveness. Our plug-and-play RAW Adapter seamlessly integrates into various RealSR models, enhancing both quantitative performance and visual quality with minor overhead, as shown in the Table of main text. To further explore the adaptivity and effectiveness of RAW Adapter, we illustrate the adaptivity of kernels learned from different RGB scenes, as shown

in Figure 8. The distributions learned from RGB in different scenes are different, which means that our method can adaptively extract information from RAW images. Additionally, we validate the effectiveness of this design through experiments: (1) removing kernels and (2) learning kernels from RAW data. We found that the above methods lead to PSNR drops of 0.195 dB and 0.171 dB, respectively, confirming the advantage of learning kernels from RGB.

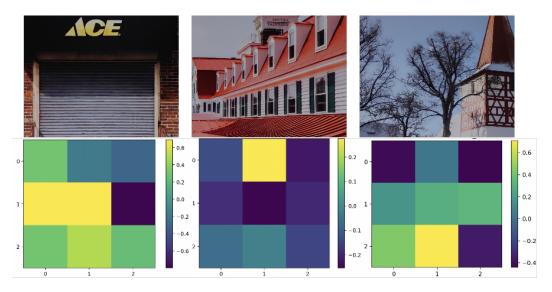


Figure 8: Demonstration of the adaptivity of kernels learned from different RGB scenes.

Ablation on Upsampling and Kernel Generation After unpacking RAW in Bayer format, the resolution is half that of the PNG format, requiring upsampling in the feature space. We compared several upsampling methods on the RRDB backbone with the M50-M benchmark: PixelShuffle, Bicubic, and (Ours) transposed convolution. As shown in Table 12, PixelShuffle and Bicubic result in PSNR drops of 0.109 dB and 0.205 dB, respectively, compared to our method, demonstrating that transposed convolution achieves superior performance. Moreover, to examine whether kernels should be generated from RGB or RAW, we conducted experiments with removing kernels, generating kernels from RAW, and (Ours) generating from RGB. Results in Table 13 show that removing kernels leads to a PSNR drop of 0.295 dB, and using RAW results in a drop of 0.271 dB compared to our method. This confirms the necessity of learning kernels from RGB for aligning RAW and RGB feature spaces.

Table 12: Ablation study on upsampling strategies (RRDB, ×2, M50-M).

Setting	PSNR↑	SSIM↑	LPIPS↓
PixelShuffle	25.953	0.768	0.351
Bicubic	25.921	0.774	0.356
Ours	26.126	0.785	0.332

Table 13: Ablation study on kernel generation (RRDB, $\times 2$, M50-M).

Setting	PSNR↑	SSIM↑	LPIPS↓
Remove Kernels	25.831	0.763	0.361
From RAW	25.955	0.771	0.349
Ours	26.126	0.785	0.332

Additional Ablation on Fusion Strategies and Learned Kernel. We further compare different fusion strategies for RGB and RAW features on the RRDB backbone using the M50-M benchmark: (S1) direct concatenation of RAW and RGB features, (S2) AdaIN-like channel modulation, and (Ours)

the proposed kernel-based fusion. As shown in Table 14, direct concatenation fails to fully exploit RAW due to the feature space gap, and modulation also struggles to capture detailed information from RAW features. In contrast, our method learns a kernel in the RGB space that aligns the distribution of RAW features to the RGB domain, effectively capturing useful information and achieving the best performance. We also study the role of the learned kernel in the RAW Adapter. Three settings are considered: (S3) removing the learned kernel (equivalent to direct concatenation), (S4) learning the kernel from RGB features but applying it only within RGB, and (Ours) applying it to RAW features for cross-domain alignment. As shown in Table 15, removing the kernel limits performance, and restricting the kernel to RGB reduces performance by 0.242 dB in PSNR. This demonstrates that using a kernel learned from RGB to align RAW features is both effective and efficient.

Table 14: Ablation study on feature fusion strategies (RRDB, ×2, M50-M).

Setting	PSNR↑	SSIM↑	LPIPS↓
S1: Concatenation	25.761	0.761	0.368
S2: Modulation	25.976	0.779	0.341
Ours	26.126	0.785	0.332

Table 15: Ablation study on learned kernel (RRDB, $\times 2$, M50-M).

Setting	PSNR↑	SSIM↑	LPIPS↓
S3: Remove Learned Kernel	25.761	0.761	0.368
S4: Kernel applied to RGB only	25.884	0.771	0.359
Ours	26.126	0.785	0.332

More Comparison of Computational Overheads. We compare the PSNR performance, number of parameters, FLOPs, runtime, and memory usage of RRDB and RRDB with our design. The Table 16 shows that our method introduces negligible additional GPU memory usage and computational complexity. However, it achieves significant improvements in image quality, as demonstrated in Tables 2-5 and Figures 1 and 5 of the main text.

More Comparison on ×8. To further demonstrate the superiority of our method across different scales, we conduct ×8 super-resolution (SR) experiments. As shown in Table 16, our method achieves a significant 0.82 dB improvement in PSNR compared to the baseline RRDB, with nearly negligible increases in computational overhead, including parameters, FLOPs, runtime, and memory usage. Specifically, while the parameter count and FLOPs only increase slightly (from 9.59M to 9.66M and 550.4G to 554.1G, respectively), the runtime remains nearly identical (0.1042s vs. 0.1048s). Additionally, the memory usage increases marginally from 12.91G to 12.99G. These results verify the practical effectiveness of the proposed method in achieving improved image quality with minimal computational cost.

Table 16: Comparison of computational overhead and performance across different scales.

Scale	Method	PSNR	Param	FLOPs	Time	Memory
×2	RRDB	25.426	9.57M	482.9G	0.0795s	10.13G
	RRDB(Ours)	26.126	9.64M	485.6G	0.0799s	10.15G
×8	RRDB			550.4G		12.91G
	RRDB(Ours)	23.234	9.66M	554.1G	0.1048s	12.99G

D MORE VISUAL RESIDUAL ANALYSIS

Details of the residual image from denoising and demosaicing. Here, we provide more details about the residual image resulting from denoising and demosaicing. As an integral part of the ISP pipeline, demosaicing cannot be bypassed because we need to obtain an RGB three-channel image from a single-channel Bayer format, which necessitates demosaicing. Therefore, we perform a step-by-step analysis. Specifically, since the output of demosaicing is three-channel while the input is

single-channel, we replicate the single-channel input to create a three-channel image. The input and output of demosaicing have the same resolution, allowing us to directly subtract them to obtain the residual image. However, the denoising process within the ISP can be bypassed, so we conduct two types of analyses for the denoising process. Both analyses indicate that denoising and demosaicing result in the loss of image details.

Bypass analysis. Here, we provide more residual analysis visualizations, as shown in Figure 9. It can be seen that the residuals on the right of all scenes contain most of the image detail information, leading us to conclude that denoising can lead to a loss of detail.



Figure 9: Visualizations of bypass analysis method. On the right of each scene is the denoised sRGB image, and on the left are the residuals with and without bypass denoising.

Step-by-step analysis. Here, we provide more residual analysis visualizations of denoising and demosaicing, as shown in Figure 10 and 11. It can be seen that the residuals on the right of all scenes contain most of the image detail information, leading us to conclude that denoising and demosaicing can lead to a loss of detail.

It can be seen that the residuals on the right of all scenes contain most of the image detail information, leading us to conclude that denoising and demosaicing can lead to a loss of detail.



Figure 10: Visualizations of step-by-step analysis method. On the left of each scene is the denoised sRGB image, and on the right are the residuals after and before denoising.



Figure 11: Visualizations of step-by-step analysis method. On the left of each scene is the denoised sRGB image, and on the right are the residuals after and before demosaicing.

E MORE ANALYSES OF OTHER MODULES IN ISP

Analyses of Image Stabilization. In Section 3, we demonstrate that the denoising and demosaicing modules within the ISP can degrade image details. Additionally, considering that modern smartphones are often equipped with image stabilization systems, utilizing an internal gyroscope to estimate a motion trajectory warp matrix in the YUV space and apply it to images, we further investigate whether this process results in detail loss by simulating both a warp and an unwarp matrix. We select 100 images from the Mate 50 Pro test set for processing and analyze the original and warped-unwarped images. By subtracting these, we derive residual images for analysis and visualization, as shown in Figure 12.

Ten volunteers participate in an evaluation, being asked: *USER: Please determine if the residual image on the right contains the structural content information of the image on the left. Answer Yes or No.* The results indicate that in 97% of the scenarios, volunteers agree that the residuals contain detailed structural information.



Figure 12: This visualization illustrates the loss of detail introduced by Image Stabilization. Examination of the residual images reveals that a significant amount of image detail is retained.

F MORE DETAILS OF TRAINING MODEL

We use the open-source and widely-used BasicSR framework to conduct experiments on three representative RealSR methods: RRDB, SwinIR, and ResShift. We utilize the public BasicSR for



Figure 13: visual comparison of the real images captured by smartphone in the wild.



Figure 14: visual comparison of the RRDB model on the Mate 50 Pro phone.

training and evaluate Real-SR methods with a total of 16 NVIDIA V100 GPUs. The training details for each method are as follows:

Training details of RRDB. During the training of RRDB, the input size is set to 128×128 , with the resolution of the ground truth set to 256×256 . The batch size per GPU is 4, and we utilize a total of 4 V100 GPUs for training RRDB.

Training details of SwinIR. For SwinIR, to facilitate rapid validation, we select the Small version of SwinIR, keeping the resolution of the input image and ground truth consistent with SwinIR. The batch size per GPU is 4, and we use a total of 8 V100 GPUs for training SwinIR.

Training details of ResShift. For ResShift, we maintain the settings consistent with the official release. While training ResShift, since our collected dataset is for $2\times$ super-resolution and the official open-source ResShift only supports $4\times$ super-resolution, we enlarge the HR images by two times to construct a $4\times$ super-resolution for training.

Our proposed RAW adapter aims to extract detailed information from RAW images to assist Real SR, thus we maintain the training scenarios of Real SR completely consistent, with only the mapping differing: one is the traditional Real SR, LR RGB \rightarrow HR RGB, and our method is LR RGB+RAW \rightarrow HR RGB.

G MORE COMPARISON

G.1 MORE VISUAL COMPARISON

Here, we present more visual comparisons of real images captured by smartphones in the wild, as shown in Figure 13. It can be observed that our proposed method achieves richer detail and superior visual quality in real-world scenarios.

Here, we present more visual comparisons of the RRDB model on the Mate 50 Pro phone, as shown in Figure 14. It can be observed that our proposed method achieves higher fidelity in image details, closely resembling the ground truth.

G.2 LARGER-SCALE USER STUDY

We conducted a larger-scale user study with 100 randomly selected real LR images captured by the Mate 50 Pro. Thirty volunteers rated the quality of generated images from RRDB and SwinIR models on a scale of 1 to 10. As shown in Table 17, compared to the original RRDB and SwinIR, which received average scores of 6.37 and 6.57, our enhanced RRDB+ and SwinIR+ models achieved higher scores of 7.81 and 8.02, respectively, due to improved detail representation from RAW data. These results further confirm the effectiveness of the RAW Adapter in enhancing visual quality.

Table 17: Larger-scale user study (Mate 50 Pro, 30 participants).

Methods	LR	RRDB	RRDB+ (Ours)	SwinIR	SwinIR+ (Ours)
Score	5.12	6.37	7.81	6.57	8.02

H USE OF LARGE LANGUAGE MODELS

Large Language Models (LLMs) were used only for grammar checking and text polishing. All research ideas, methods, and analyses are solely by the authors.