# Exposure-Limited Image Enhancement with Generative Diffusion Prior

Baiang Li, Sizhuo Ma, Yanhong Zeng, Xiaogang Xu, Youqing Fang, Zhao Zhang, *Senior Member, IEEE*, Jian Wang\* and Kai Chen\*

Abstract—Many consumer cameras are equipped with 8-bit image sensors, which often struggle to capture scenes with a High Dynamic Range (HDR). This limitation can result in overexposed or underexposed regions, a loss of fine details due to low bit-depth compression, skewed color distributions, and noticeable noise in dark areas. Traditional Standard Dynamic Range (SDR) image enhancement methods typically focus on color mapping by expanding the color range and adjusting brightness. However, they often fail to restore details in *dynamic range extremes*, *i.e.* regions where pixel values approach the minimum or maximum limits. We define "exposure-limited image enhancement" as the process of enhancing images with large missing areas due to exposure issues within the SDR space, which differs from existing "mis-exposed image enhancement" methods primarily aimed at correcting color distributions. To enhance these exposure-limited images and overcome the limitations of current models, we propose a novel two-stage approach. In the first stage, we remap color and brightness to a suitable range while preserving existing details. In the second stage, we use a diffusion prior to generate content in severely overexposed or underexposed regions, which are otherwise lost during capture. Notably, this generative refinement module can also serve as a plug-and-play component alongside existing enhancement methods. Extensive experiments demonstrate that our method significantly improves image quality and detail, outperforming state-of-the-art techniques in dynamic range extremes. The project page is at https://Sagiri0208.github.io.

Index Terms—Computational Photography, Image Enhancement, Generative Methods

# **1** INTRODUCTION

EAL-WORLD scenes often feature broad dynamic 2 Rranges, yet many smartphone cameras are equipped 3 with 8-bit image sensors with limited dynamic ranges. As a 4 result, these sensors cannot simultaneously capture details 5 in both the bright sun and shaded leaves. One popular 6 strategy for increasing the dynamic range is exposure brack-7 eting [6], which merges multiple low-exposure and high-8 exposure shots into a single high dynamic range (HDR) 9 image. However, this approach must be enabled at capture 10 time, extends the shooting duration, and requires computa-11 tionally intensive motion compensation. 12

Deep neural networks have been developed to restore 13 and enhance overexposed and underexposed regions from 14 a single image. This is usually done by either reconstruct-15 ing an HDR image [4] followed by a tone mapping step, 16 or by directly predicting an enhanced image [1]. How-17 ever, these methods still struggle to consistently deliver a 18 truly satisfying visual experience [7], [8]. The challenge is 19 20 particularly pronounced in *dynamic range extremes*—areas where the pixel values are close to minimum or the max-21 imum possible values, as illustrated in Figure 1. Although 22

Convolutional Neural Networks (CNNs) and transformers excel at tone mapping (e.g., low-light enhancement [9] and denoising [10]) in dark regions, they are not capable of reconstructing large areas of content where data is essentially lost at capture. Consequently, existing methods often produce blurry, unnatural content in such regions (Figure 1(c)). In this work, we aim to extend the capabilities of *exposure-limited image enhancement*, which we define as the integration of tone mapping, noise reduction for low light conditions, detail compensation due to low bit-depth, and *generating* image details obscured or completely lost due to the camera's restricted dynamic range.

Large-scale generative models trained on extensive textimage pairs [11], [12] such as Stable Diffusion [13] excel at synthesizing realistic images and can potentially serve as a powerful tool for generating the details in dynamic range extremes. Motivated by the recent success in leveraging diffusion prior in low-level vision tasks, we propose a two-stage framework tailored to exposure-limited SDR image enhancement. This framework takes an SDR image as input and directly predicts an output SDR image with enhanced color distribution and details. In the first stage, we perform a global adjustment using Latent-SwinIR<sub>c</sub> (LS), a transformerbased model [14] designed to harmonize color distributions and re-map extremely bright or dark areas into a visually pleasing range<sup>1</sup>. This is accomplished through a carefully defined color-mapping loss computed over color histograms.

In the second stage, the initially color-adjusted image is

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

43

44

45

46

47

48

49

50

B. Li, Y. Zeng, Y. Fang and K. Chen are with Shanghai AI Laboratory, Shanghai, China. Also, B. Li is with Hefei University of Technology, Hefei 230601, China.(e-mails: ztmotalee@gmail.com, zengyh1900@gmail.com, fangyouqing@pjlab.org.cn, chenkai@pjlab.org.cn).

X. Xu is with The Chinese University of Hong Kong, Hong Kong 100871, China; affiliated with Zhejiang University, Hangzhou 310058, China. (email: xiaogangxu00@gmail.com)

S. Ma and J. Wang are with Snap Research, Snap Inc., New York, NY 10036, USA. (e-mails: sizhuoma@gmail.com, jwang4@snapchat.com)

Z. Zhang is with Hefei University of Technology, Hefei 230601, China. (email: cszzhang@gmail.com)

<sup>\*</sup> J. Wang and K. Chen are the co-corresponding authors. Besides, J. Wang initialized the project.

<sup>1.</sup> This process is analogous to tone mapping, except that our input is also an SDR image.



Fig. 1: (a) Real-world scenarios have broad dynamic ranges. However when captured by normal low bit mobile camera, the chosen exposure may often face over-saturated bright regions or heavily quantized dark areas with strong noise. (b) Existing works on mis-exposed/low light image enhancement [1], [2], multi-exposure HDR reconstruction [3] and singleexposure HDR reconstruction [4], [5] can enhance SDR images with no or small area of missing content, they are not designed to deal with *large dynamic range extreme areas*. (c) As a result, such methods often create blurry and unnatural imagery at dynamic range extremes. In this work, we aim to enhance these exposure-limited situations by analyzing and decomposing the complex task into the following challenges: (a) color and brightness mapping, (b) denoising of dark areas, (c) detail enhancement in low bit-depth regions, and (d) content generation for saturated or near-black regions.

refined by our diffusion-based model (named Sagiri). Even 52 though the images are mapped to normal color distribution, 53 content lost due to exposure limitations remains missing. 54 Sagiri utilizes the powerful generative capabilities of pre-55 trained diffusion models with ControlNet [15] to further en-56 hance the partially restored content and synthesizes realistic 57 image details in the dynamic range extremes. Trained using 58 a carefully crafted synthetic degradation strategy, Sagiri can 59 also function as a plug-and-play module to improve results 60 from existing SDR image enhancement and HDR recon-61 struction methods. Additionally, we introduce an adaptive 62 regional processing technique during sampling, allowing 63 users to guide content generation through custom prompts 64 65 such as text or pixel masks. It is important to note that we intentionally limit our scope to enhancing SDR images to 66 leverage the powerful diffusion models pretrained on SDR 67 images. Extending this technique to the HDR domain could 68 enable the use of flexible tone-mapping operators, which we 69 leave for future work 70

Our key contributions are summarized as follows: 71

- We present LS-Sagiri, a novel two-stage approach 72
- 73 for exposure-limited image enhancement. Stage 1 74
  - adjusts the overall color and brightness, while Stage 2

refines and generates missing details.

- We integrate a powerful generative diffusion based model to generate realistic content in saturated or black regions, while also enhancing fine details in areas affected by low bit-depth. Additionally we adopt a two-step training strategy to make it a plug-andplay module, effectively enhancing existing methods.
- Comprehensive experiments demonstrate that our approach delivers superior visual and quantitative performance and can flexibly enhance results from a wide range of existing techniques.

## 2 **RELATED WORK**

In this section, we review related work on similar tasks and highlight their differences from the exposure-limited image enhancement problem we address. Additionally, we discuss image inpainting methods in supplementary material, which, although relevant, cannot directly substitute for the proposed Sagiri network.

#### **Mis-Exposed and Low Light Image Enhancement** 21 93

80 81 82

83

84

85

86

87

88

89

90

91

92

75

76

77

78



Fig. 2: Method overview. (a) Given a degraded input, the Stage 1 model performs color mapping to adjust the entire image to a more balanced color distribution. (b) In the Stage 2, the color-adjusted image is concatenated with a random noise map and sent to the parallel VAE encoder. After a shape-adjusting convolutional layer, the encoded latent feature map is then sent to the decoder in the denoising U-Net. (c) A default unknown region mask (the region where the pixel has the maximum or minimum value is selected) is used during training to indicate the dynamic range extremes. (d) During inference, the users can define their own mask. Note that the input and output of each stage is both SDR image. Please zoom in on the images to observe the differences.

Mis-Exposed Image Enhancement. Rather than attempting 94 to recover or generate HDR images, some research focuses 95 on correcting the exposure of single SDR images to produce 96 outputs directly within the SDR domain. Wang et al. [1] 97 98 observe that an image's local color distributions often contain both overexposed and underexposed areas and propose 99 a method to directly enhance these regions within the LDR 100 domain. Li et al. [16] further notice that these error exposed 101 regions display opposite color tone distribution shifts and 102 further propose a model to handle the shift. These methods 103 are functionally similar to our first-stage model Latent-104 SwinIR<sub>c</sub> and do not have strong generative capabilities to 105 fill in the missing details in *dynamic range extremes*. 106

Low Light Image Enhancement. Low light image enhance-107 ment is another related task that focuses on color map-108 ping and noise reduction for low-light images within the 109 SDR domain. Early methods like [17], [18] face limitations 110 when addressing complex real-world scenarios, which later 111 proposed deep learning methods have improved upon. 112 LLNet [2] utilizes a deep autoencoder to transform un-113 derexposed images into enhanced versions, simultaneously 114 reducing noise and increasing brightness. Retinex-Net [19] 115 builds on the classical Retinex theory by decomposing im-116 ages into reflectance and illumination components, allowing 117 for targeted brightness adjustment and effective denoising. 118 More advanced CNN and Transformer-based methods, such 119 as [20], [21], [22], [23] and generative methods include [24], 120 121 [25] have been proposed recently. However, these methods lack the capability to generate content in large blank regions 122 and are not designed to address overexposure, thus they are 123 124 not included in our comparison.

# 2.2 Single Image HDR Reconstruction

Non-Generative Methods. As the reconstruction of HDR images from a single input primarily involves color mapping and pixel value prediction without large exposure-limited region, many methods approach this task as a restoration problem without employing generative networks. Among these, SingleHDR [4] inverts the SDR image formation pipeline to recover HDR information. However, pipelinebased methods can lead to error accumulation at each intermediate step. HDRUNet [5] learns an end-to-end mapping for single-image HDR reconstruction, featuring denoising and dequantization. RawHDR [7] targets raw images by learning exposure masks to address challenging high dynamic range regions. Le et. al explores another direction by reconstructing images through multi-exposure generation [26], predicting images at various exposure levels from the input image and subsequently fusing them using standard HDR merging strategies. Other representative methods include [27], [28], [29], [30]. Although these methods have advanced HDR image reconstruction, their primary focus is not on addressing large exposed regions, resulting in weaker performance when applied to our specific task.

Generative Methods. Recent generative models offer new possibilities for HDR reconstruction. Fei et al. [31] in-148 troduce a diffusion-based framework for unsupervised restoration and enhancement, utilizing hierarchical guid-150 ance and patch-level operations to produce high-quality 151 results. However, it requires multiple LDR inputs and involves a lengthy inference process. GlowGAN [8] uses a 153 generative adversarial network to learn HDR content from in-the-wild LDR images without explicit supervision but

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

147

149

152

154

struggles with large overexposed regions. In contrast, our 156 second-stage Sagiri model leverages a diffusion prior to 157 effectively address this challenge. 158

#### 3 Method 159

Our method aims to restore and enhance exposure-limited 160 images using a two-stage framework. In the first stage, we 161 employ a restoration model to adjust brightness and color 162 distributions, yielding an image whose color statistics more 163 closely resemble ground truth (GT) data. However, this 164 restoration model alone has limited generative capacity and 165 struggles to reconstruct details lost in *dynamic range extremes*. 166 To address this, we introduce a second-stage model that 167 leverages a diffusion prior to generate and refine missing 168 or severely degraded regions. 169

Specialized Loss Functions. To effectively guide the learn-170 ing process at each stage, we design distinct loss func-171 tions suited to each model's strengths. In Stage 1, our 172 color reconstruction loss prioritizes brightness adjustment and 173 color alignment by matching the predicted image's color 174 histogram to that of the target image. In Stage 2, we define a 175 content enhancement loss focused on generating high-fidelity 176 textures and shifting the content distribution closer to that 177 of detail-rich references. This encourages the model to syn-178 thesize missing details effectively. 179

#### 3.1 Color Mapping 180

Our first-stage model is based on SwinIR [14] but with mod-181 ified pre- and post-processing. Specifically, we apply a pixel 182 unshuffle operation to downsample the original low-quality 183 input by a factor of 8. Subsequently, a  $3 \times 3$  convolutional 184 layer extracts shallow features in color space. These features 185 pass through Residual Swin Transformer Blocks (RSTB) 186 for processing. Nearest-neighbor upsampling and another 187  $3 \times 3$  convolutional layer are then repeated three times to 188 return the features to the original resolution. We refer to 189 this modified model as Latent-SwinIR<sub>c</sub> (seen in top left of 190 the Figure 2), which focuses on global color and brightness 191 adjustment while offering preliminary content recovery. The 192 core processing is performed in a downsampled latent space 193 to reduce the computational cost, as the network must take 194 the entire image as input rather than tiles to capture the 195 color distribution accurately. To improve color mapping 196 and brightness adjustment, we introduce the following loss 197 function: 198

$$L_{\rm color} = \lambda_1 L_{\rm mse} + \lambda_2 L_{\rm cd} + \lambda_3 L_{\rm fdp}, \tag{1}$$

where  $L_{\rm mse}$ ,  $L_{\rm cd}$ , and  $L_{\rm fdp}$  are the MSE, Color Distribution, 199 and Frequency Domain Preservation losses, respectively. 200 The coefficients  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  balance each component. 201 *Further details are provided in the supplementary material.* 202

### 3.2 Conditional Generation 203

After the first stage, the restored image may still contain 204 regions of poor visual quality, particularly in dynamic range 205 206 extremes. Traditional methods often struggle to synthesize high-fidelity details in these areas [1], [7]. We harness the 207 generative power of a pretrained diffusion model to over-208 209 come this limitation (Figure 2).

First, we encode the Stage 1 output via a Variational Autoencoder (VAE) [32] to obtain a latent representation. This latent is then combined with noise and fed into a parallel encoder module, which mirrors the encoder architecture in the U-Net denoiser. The outputs from different encoder blocks serve as latent controls, concatenated with the U-Net's decoder features. Newly introduced parameters are initialized to zero, while the pretrained denoising U-Net remains frozen except for  $1 \times 1$  convolutional layers added before each concatenation.

During training, we mark the pixels that exceed the dynamic range (which will has value of 0 or reach the maximum in a SDR image) as "unknown" regions (See Figure 2 (c)). To ensure that the detail generation is focused to the unknown regions while being semantically and aesthetically harmonious with the known regions, at the *t*-th step, the known regions are preserved by directly diffusing the initial latent feature, while the unknown regions are inferred from the model's denoised output [33]:

$$x_{t-1}^{\text{known}} \sim \mathcal{N}\left(\sqrt{\overline{\alpha}_t} x_0, (1 - \overline{\alpha}_t)I\right),$$
 (2)

$$x_{t-1}^{\text{unknown}} \sim \mathcal{N}(\mu_{\theta}(x_t, t), \Sigma_{\theta}(x_t, t)),$$
 (3)

$$x_{t-1} = m_{\text{latent}} \odot x_{t-1}^{\text{known}} + (1 - m_{\text{latent}}) \odot x_{t-1}^{\text{unknown}}, \quad (4)$$

where the known pixels  $x_{t-1}^{\text{known}}$  are directly diffused from  $x_0$ , and the unknown pixels  $x_{t-1}^{\text{unknown}}$  are sampled from the diffusion model, blended by a mask  $m_{\text{latent}}$ . After t steps, the final denoised latent is decoded by the LDM decoder.

To guide Sagiri toward realistic content generation, we introduce the following loss:

$$L_{\rm content} = \lambda_4 L_{\rm mse} + \lambda_5 L_{\rm ssim} + \lambda_6 L_{\rm fdp}, \tag{5}$$

where  $L_{\rm mse}$ ,  $L_{\rm ssim}$ , and  $L_{\rm fdp}$  ensure structural fidelity, realistic textures, and frequency consistency. Coefficients  $\lambda_4, \lambda_5, \lambda_6$  balance these terms. Details are provided in the 237 supplementary material.

## 3.3 Training Strategy

**Overall Pipeline.** We train Latent-SwinIR<sub>c</sub> on the HDR-Real dataset [4] to learn appropriate color and brightness mappings. For Sagiri, we first pretrain it on the large-scale Places365 dataset [34] to expand its capability of generating diverse scenes, and then finetune it on HDR-Real.

Degradation Generation in Pre-training. To bridge the domain gap, we simulate Latent-SwinIR<sub>c</sub> outputs during Sagiri's pre-training on Places365 by crafting various realistic degradation patterns. Specifically, we generate a "degradation mask" using random lines of varied thickness, followed by dilation and Gaussian blurring. This mask is used to blend the original image with a heavily blurred version, producing blur-like artifacts mimicking over- or underexposed areas. This approach improves Sagiri's ability to serve as a *plug-and-play* module not only for Latent-SwinIR<sub>c</sub> but also for other methods we compare with (Figure 4).

Unknown Region Mask. During Sagiri's pre-training, we do not apply the unknown region mask; the model is encouraged to autonomously identify and handle low-quality areas. During finetuning, however, we apply the binary mask to emphasize the unique challenge of inpainting over-/under-exposed regions.

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

238

239

240

241

242

243

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260



Fig. 3: **Performance of LS-Sagiri**. Previous restoration-based methods can only restore exposure-limited areas to blurry content. Although Glow-GAN [8] is a generative method, it fails to handle *large* overexposed regions, often rendering them black. In contrast, our method can generate realistic content based on existing information and specified mask areas.



Fig. 4: The first row is the result obtained using our degradation strategy, while the second row is the reference images. We aim to simulate the degradation caused by other models in dynamic range extremes during SDR enhancement and train Sagiri to handle these situations effectively.

# 262 4 EXPERIMENTS

## 263 4.1 Training and Inference Settings

Training. We train Latent-SwinIR<sub>c</sub> on the HDR-Real training 264 set [4] for 150,000 iterations using a batch size of 16. Our 265 Sagiri model is initialized from pretrained Stable Diffusion 266 v2.1. We first pretrain Sagiri on 250,000 randomly selected 267 images from the Places365 dataset [34] for 70,000 steps, then 268 fine-tune it on the HDR-Real training set for an additional 269 20,000 steps. All training stages use the Adam optimizer [35] 270 with a learning rate of  $1 \times 10^{-4}$ , conducted on four NVIDIA 271 A100 GPUs. The loss coefficients  $\lambda_1$ ,  $\lambda_2$ ,  $\lambda_3$ ,  $\lambda_4$ ,  $\lambda_5$ , and  $\lambda_6$ 272 are set to 10, 1, 0.1, 1, 1, and 0.01, respectively. The total 273 274 training process takes 3 days using 4 NVIDIA A100 GPUs. Inference. During inference, our model processes an SDR 275 input image, applying an unknown region mask that detects 276 pixels with values of 0 or 255. We use just 30 steps of DDPM 277 sampling [36]. Note that all the visualizations are from the 278 test set or real-world examples. 279

Datasets. We evaluate performance on HDR-Real [4], 280 NTIRE [37], HDR-Eye [38], Eye-over, and Eye-under. The 281 latter two are constructed by uniformly adjusting the expo-282 sure values of HDR-Eye to create overexposed or underex-283 284 posed images. Since existing datasets lack a large number of test images with severe content loss, these newly created 285 sets better assess our method's effectiveness in extreme 286 287 conditions.

Prompt Usage in Training and Inference. We employ 288 CogVLM [39] to summarize input images into prompts. 289 During Sagiri's fine-tuning on HDR-Real, these prompts are 290 generated using the ground-truth images to align the model 291 with prompt inputs. For inference on the HDR-Real testing set [4], prompts are produced from low-quality inputs. For 293 the HDR-Eye [38], Eye-over, Eye-under, and NTIRE [37] 294 datasets, we do not provide prompts. This design aims to 295 evaluate Sagiri's adaptability under different scenarios. 296

## 4.2 Results

**Baseline Methods.** We compare our approach with Single-HDR [4], LCDPNet [1], HDRUNet [5], and GlowGAN [8]. For fairness, we uniformly apply the Reinhard tone mapping function [40] across the training and testing sets to generated pairs of SDR images for training and evaluating SDR methods (LCDPNet and ours), and as a common postprocessing step for evaluating HDR methods (SingleHDR, HDRUNet and GlowGAN). Additional comparisons with GDP [31] are provided in the supplementary material.

Metrics. To assess both our complete LS-Sagiri pipeline 307 and the generalizability of Sagiri as a refinement module 308 for other models, we rely on no-reference metrics such 309 as BRISQUE [41], NIQE [42], MANIQA [43], and CLIP-310 IQA [44], which focus on visual quality. We do not use 311 PSNR, SSIM [45], or LPIPS [46] because prior studies [37], 312 [47], [48] have shown that these metrics can be unreliable 313 for generative models and this challenging generative task. 314 **Performance of Latent-SwinIR**<sub>c</sub>. Figures 7(a–f) demonstrate 315 the ability of Latent-SwinIR<sub>c</sub> to correct image color distri-316 butions. Rows 1–3 show LQ images captured in extremely 317 dark conditions, a challenging scenario for image enhance-318 ment. SingleHDR reveals more details yet suffers from low 319 contrast. LCDP-Net and HDRUNet fail to adequately adjust 320 the brightness, leaving the result underexposed. GlowGAN 321 struggles to restore quantized details, leading to images 322 lacking essential textures. Thanks to our novel loss formula-323 tion, Latent-SwinIR<sub>c</sub> achieves superior brightness and color 324 mapping without compromising detail, delivering high-325 quality, natural-looking images. 326

297

298

299

300

301

302

303

304

305

TABLE 1: Quantitative results on HDR-Real [4], NTIRE [37], HDR-Eye [38], Eye-over and Eye-under datasets. The latter two datasets are made by uniformly adjusting the exposure value of HDR-Eye dataset to synthesize datasets with large areas at dynamic range extremes. In addition to comparing the performance of our pipeline with existing methods, we plugged Sagiri into each model to see performance improvements. The results show that (1) Sagiri enhances the performance of each method, and (2) LS-Sagiri achieves the best overall results.

Datasets		OR-Real		NTIRE				
Metrics	BRISQUE↓	NIQE↓	MANIQA↑	CLIP-IQA↑	BRISQUE↓	NIQE↓	MANIQA↑	CLIP-IQA↑
SingleHDR [4]	23.597	20.839	0.367	0.387	22.730	<b>21.399</b>	0.250	0.411
SingleHDR+Sagiri	<b>19.855</b>	<b>20.326</b>	<b>0.556</b>	<b>0.649</b>	<b>10.211</b>	21.622	<b>0.385</b>	<b>0.676</b>
LCDPNet [1]	30.704	20.660	0.344	0.383	19.237	<b>20.978</b>	0.267	0.415
LCDPNet+Sagiri	<b>24.464</b>	<b>20.318</b>	<b>0.542</b>	<b>0.641</b>	<b>9.951</b>	21.622	<b>0.385</b>	<b>0.674</b>
HDRUNet [5]	41.521	21.388	0.341	0.361	52.898	22.752	0.229	0.377
HDRUNet+Sagiri	<b>24.935</b>	<b>20.704</b>	<b>0.503</b>	<b>0.609</b>	<b>21.353</b>	<b>21.749</b>	<b>0.397</b>	<b>0.650</b>
GlowGAN [8]	36.727	21.774	<b>0.470</b>	0.503	21.769	<b>24.053</b>	<b>0.403</b>	0.478
GlowGAN+Sagiri	<b>22.840</b>	<b>21.602</b>	0.443	<b>0.554</b>	<b>15.549</b>	24.078	0.354	<b>0.511</b>
Latent-SwinIR <sub>c</sub>	35.407	21.457	0.291	0.303	31.298	22.000	0.224	0.392
LS-Sagiri	<b>19.725</b>	<b>20.309</b>	<b>0.569</b>	<b>0.670</b>	<b>9.724</b>	<b>21.652</b>	<b>0.395</b>	<b>0.671</b>

Datasets	HDR-Eye			Eye-over			Eye-under		
Metrics	BRISQUE↓	MANIQA↑	CLIP-IQA↑	BRISQUE↓	MANIQA↑	CLIP-IQA↑	BRISQUE↓	MANIQA↑	CLIP-IQA↑
SingleHDR [4]	18.338	0.452	0.466	20.573	0.447	0.428	33.675	0.244	0.244
SingleHDR+Sagiri	<b>15.092</b>	<b>0.570</b>	<b>0.697</b>	<b>14.969</b>	<b>0.557</b>	<b>0.676</b>	<b>13.477</b>	<b>0.339</b>	<b>0.523</b>
LCDPNet [1]	20.672	0.453	0.475	26.374	0.398	0.365	54.493	0.311	0.335
LCDPNet+Sagiri	14.137	<b>0.543</b>	<b>0.665</b>	<b>14.973</b>	<b>0.478</b>	<b>0.638</b>	<b>37.825</b>	<b>0.382</b>	<b>0.552</b>
HDRUNet [5]	27.672	0.418	0.390	24.545	0.454	0.410	72.920	0.364	0.403
HDRUNet+Sagiri	<b>14.846</b>	<b>0.555</b>	<b>0.662</b>	<b>15.905</b>	<b>0.560</b>	<b>0.668</b>	<b>40.954</b>	<b>0.460</b>	<b>0.610</b>
GlowGAN [8]	<b>16.042</b>	<b>0.506</b>	<b>0.536</b>	<b>16.930</b>	<b>0.503</b>	<b>0.561</b>	46.667	<b>0.356</b>	<b>0.483</b>
GlowGAN+Sagiri	19.775	0.430	0.473	20.040	0.401	0.466	<b>37.745</b>	0.286	0.432
Latent-SwinIR <sub>c</sub>	25.870	0.329	0.286	25.345	0.321	0.286	45.168	0.256	0.252
LS-Sagiri	<b>14.777</b>	<b>0.538</b>	<b>0.675</b>	<b>14.667</b>	<b>0.535</b>	<b>0.669</b>	<b>12.066</b>	<b>0.462</b>	<b>0.660</b>

In Rows 4–8, five images with decreasing exposure times are tested. As the exposure decreases, existing methods exhibit fading contrast (SingleHDR), dull brightness (LCDP-Net and HDRUNet), or further detail loss (GlowGAN). In contrast, Latent-SwinIR<sub>c</sub> consistently preserves robust color and brightness distributions, demonstrating remarkable robustness to varying exposure levels.

**Performance of LS-Sagiri.** Figures 7(g) and 3 present the 334 performance of the complete LS-Sagiri pipeline. Despite 335 the overall robust performance of Latent-SwinIR<sub> $c_1$ </sub> it is no-336 table that the results still contain certain degradations, such 337 as blurry areas in exposure-limited regions, as illustrated 338 in Figure 7. In Figure 7(g), Sagiri clearly refines Latent-339 SwinIR<sub>c</sub>'s outputs by synthesizing realistic details, notably 340 improving the perceived image quality. Figure 3 features 341 inputs with large overexposed regions, where SingleHDR 342 introduces blurry areas, LCDP-Net and HDRUNet fail to 343 regulate brightness, and GlowGAN cannot recover heavily 344 quantized details. Only LS-Sagiri succeeds in filling over-345 saturated regions with coherent, realistic details. Table 1 346 shows that LS-Sagiri achieves the top performance across 347 nearly all metrics, confirming its efficacy in enhancing LDR 348 349 images and its broad applicability to various datasets.

We note that Latent-SwinIR $_c$  alone does not always yield top metric scores. This may stem from the limitations of noreference image quality metrics, which do not thoroughly account for global brightness distribution. Future research is needed to explore more comprehensive metrics to better evaluate such enhancements.

Sagiri as a Plug-and-Play Module. Beyond refining Latent-SwinIR<sub>c</sub> outputs, Sagiri can seamlessly integrate with other models. Figure 8 shows how Sagiri corrects dynamic range extremes even when the initial enhanced images vary significantly in quality, ultimately improving the overall perceptual quality. Table 1 supports this versatility: Sagiri significantly boosts nearly every baseline's output. Minor exceptions (e.g., GlowGAN) are detailed in the supplementary material.

**Further evaluations.** A user study on the subjective quality of the methods and a comparison with inpainting methods can be found in the supplementary material.

## 4.3 Ablation Studies

**Importance of the Two-Stage Model.** To verify whether our two-stage pipeline is essential, we tested Sagiri alone for both color restoration and detail enhancement. As shown in Figure 5 (Top), Sagiri alone struggles with color restoration and brightness adjustment, demonstrating the necessity of Latent-SwinIR<sub>c</sub> as the first stage.

**Effect of Pre-training and Prompts.** Figure 5 (Middle) shows that our pre-training approach and prompt-guided generation significantly improve visual quality, while the content reconstruction loss (ConRLoss) further enhances structural integrity. In addition, Figure 5 (Bottom) illustrates how users can control (1) the region to be generated by

354

355

357

358

359

360

361

362

363

364

365

366

368

369

370

371

372

373

374

375

376

377

378

379



Fig. 5: Ablation studies. (Top) We enforce Sagiri to learn both color distribution correction and details generation, which leads to weak color mapping capabilities. (Middle) Ablation of pretraining, text prompt and content reconstruction loss. Prompts generated by CogVLM [39]: "A white waterfall is flowing down from the cliff, surrounded by rocks and trees." (Bottom) We use different user-defined unknown region mask and different prompts on Sagiri to refine SingleHDR's [4] results. Left: We manually select the red box. Right: We select the entire image. Prompt a: "The sky is filled with clouds." Prompt b: "The sun is setting, and the sky is filled with clouds." Please zoom in to see more details.



Fig. 6: Use different prompts to control the generated results. Prompt a: "A building with a **red** brick exterior, white columns, and a **black** door..." Prompt b: "A building with a **black** brick exterior, white columns, and a **red** door...". Furthermore, The model has poor responsiveness to prompts that do not fit the current context, as we found. Prompt c: "The **sun** is setting in the forest, and the trees are **black**." Prompt d: "The **moon** is setting in the forest, and the trees are **green**". Please zoom in to see more details.

customizing the unknown-region mask, and (2) the content
to be generated by supplying a user-defined text prompt.
More results featuring prompt-guided generation in ex treme regions are shown in Figure 6.

## 385 5 CONCLUSION

We propose a novel pipeline for exposure-limited image 386 enhancement, anchored by our robust and flexible Sagiri 387 model. The pipeline comprises two stages: Stage 1 (Latent-388 SwinIR<sub>c</sub>, LS) rectifies brightness and color distributions, 389 while Stage 2 (Sagiri) synthesizes content in missing or 390 severely degraded regions and refines overall details. By 391 design, Sagiri is also compatible as a plug-and-play module, 392 allowing it to enhance outputs from diverse restoration 393 394 methods. Our experiments confirm the superior performance of the two-stage LS-Sagiri framework and demonstrate Sagiri's remarkable ability to generate realistic details. **Besides, please note that given the challenging task,** which is open end, and the generative character of stable diffusion based method, our approach is mainly designed for beautify the given input, not to restrictly following the reference image. Currently, Sagiri produces only SDR images due to limitations of the Stable Diffusion on which it is built. A promising future direction is to extend Sagiri to generate HDR outputs, thus giving users greater flexibility in applying customized tone mapping.

395

396

397

399

400

401

402

403

404



Fig. 7: (a-f) Performance of Latent-SwinIRc (LS). Existing methods often struggle to recover content in regions with extreme dynamic range. In contrast, Latent-SwinIRc, thanks to its uniquely designed loss function, captures a more balanced color distribution. This is evident in the Row 1–3. Moreover, in Row 4–8, where the input low-quality (LQ) images in column (a) exhibit a gradual decrease in exposure, columns (b-f) show that the performance of existing methods deteriorates with decreasing exposure. Although SingleHDR provides results closest to our method, it still produces low-contrast "hazy" outputs at low exposures. In (f), LS demonstrates robust preservation of color and brightness despite the decreasing exposure levels. Additionally, (g) highlights the Sagiri model's excellence in generating detailed content across large regions, thereby enhancing the overall quality. Please Zoom in the figures for details.



Fig. 8: **Sagiri as a plug-and-play module.** Although the images generated by the baselines significantly differ from each other, Sagiri shows strong versatility and improves the visual quality of almost all of them. Additionally, the combination of LS-Sagiri surpasses the performance of other models integrated with Sagiri, confirming the robustness and adaptability of our framework.

#### REFERENCES 406

427

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

465

466

467

468

469

470 471

- [1] H. Wang, K. Xu, and R. W. Lau, "Local color distributions prior 407 for image enhancement," in ECCV, 2022. 408
- [2] K. G. Lore, A. Akintayo, and S. Sarkar, "Llnet: A deep autoen-409 coder approach to natural low-light image enhancement," Pattern 410 Recognition, vol. 61, pp. 650-662, 2017. 411
- Z. Liu, Y. Wang, B. Zeng, and S. Liu, "Ghost-free high dynamic 412 [3] range imaging with context-aware transformer," in ECCV, 2022. 413
- Y.-L. Liu, W.-S. Lai, Y.-S. Chen, Y.-L. Kao, M.-H. Yang, Y.-Y. 414 [4] Chuang, and J.-B. Huang, "Single-image hdr reconstruction by 415 learning to reverse the camera pipeline," in CVPR, 2020. 416
- X. Chen, Y. Liu, Z. Zhang, Y. Qiao, and C. Dong, "Hdrunet: Single 417 [5] 418 image hdr reconstruction with denoising and dequantization," in CVPR, 2021. 419
- P. E. Debevec and J. Malik, "Recovering high dynamic range 420 [6] radiance maps from photographs," in *SIGGRAPH*, 1997. Y. Zou, C. Yan, and Y. Fu, "Rawhdr: High dynamic range image 421
- 422 [7] reconstruction from a single raw image," in ICCV, 2023. 423
- 424 [8] C. Wang, A. Serrano, X. Pan, B. Chen, K. Myszkowski, H.-P. Seidel, Theobalt, and T. Leimkühler, "Glowgan: Unsupervised learning 425 of hdr images from ldr images in the wild," in IĈCV, 2023. 426
- C. Li, C. Guo, L. Han, J. Jiang, M.-M. Cheng, J. Gu, and C. C. Loy, [9] "Low-light image and video enhancement using deep learning: 428 A survey," IEEE transactions on pattern analysis and machine intelli-429 gence, 2021. 430
  - [10] S. W. Zamir, A. Arora, S. Khan, M. Hayat, F. S. Khan, and M.-H. Yang, "Restormer: Efficient transformer for high-resolution image restoration," in Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2022, pp. 5728-5739.
  - [11] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," NeurIPS, 2020.
  - [12] J. Song, C. Meng, and S. Ermon, "Denoising diffusion implicit models," in ICLR, 2021.
  - [13] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in CVPR, 2022.
  - [14] J. Liang, J. Cao, G. Sun, K. Zhang, L. Van Gool, and R. Timofte, 'Swinir: Image restoration using swin transformer," in ICCV, 2021.
  - [15] L. Zhang, A. Rao, and M. Agrawala, "Adding conditional control to text-to-image diffusion models," in ICCV, 2023.
  - [16] Y. Li, K. Xu, G. P. Hancke, and R. W. Lau, "Color shift estimation-and-correction for image enhancement," in *Proceedings of the* IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2024, pp. 25389-25398.
  - [17] D. J. Jobson, Z.-u. Rahman, and G. A. Woodell, "Properties and performance of a center/surround retinex," IEEE transactions on image processing, vol. 6, no. 3, pp. 451-462, 1997.
  - [18] S. M. Pizer, E. P. Amburn, J. D. Austin, R. Cromartie, A. Geselowitz, T. Greer, B. ter Haar Romeny, J. B. Zimmerman, and K. Zuiderveld, "Adaptive histogram equalization and its variations," Computer vision, graphics, and image processing, vol. 39, no. 3, pp. 355-368, 1987.
  - [19] C. Wei, W. Wang, W. Yang, and J. Liu, "Deep retinex decomposition for low-light enhancement," arXiv preprint arXiv:1808.04560, 2018.
- [20] C. Guo, C. Li, J. Guo, C. C. Loy, J. Hou, S. Kwong, and R. Cong, 461 462 "Zero-reference deep curve estimation for low-light image enhancement," in Proceedings of the IEEE/CVF conference on computer 463 464 vision and pattern recognition, 2020, pp. 1780-1789.
  - [21] L. Ma, T. Ma, R. Liu, X. Fan, and Z. Luo, "Toward fast, flexible, and robust low-light image enhancement," in *Proceedings of the* IEEE/CVF conference on computer vision and pattern recognition, 2022, pp. 5637-5646.
  - [22] X. Xu, R. Wang, C.-W. Fu, and J. Jia, "Snr-aware low-light image enhancement," in Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2022, pp. 17714–17724.
- [23] B. Li, H. Zheng, Z. Zhang, Y. Zhao, Z. Zhao, and H. Zhang, 472 "Dynamic grouped interaction network for low-light stereo image 473 enhancement," in Proceedings of the 31st ACM International Confer-474 ence on Multimedia, 2023, pp. 2468-2476. 475
- 476 [24] G. Kim, D. Kwon, and J. Kwon, "Low-lightgan: Low-light enhancement via advanced generative adversarial network with 477 478 task-driven training," in 2019 IEEE International conference on image processing (ICIP). IEEE, 2019, pp. 2811–2815. [25] X. Yi, H. Xu, H. Zhang, L. Tang, and J. Ma, "Diff-retinex: Re-479
- 480 thinking low-light image enhancement with a generative diffusion 481

model," in Proceedings of the IEEE/CVF International Conference on Computer Vision, 2023, pp. 12302-12311.

- [26] P.-H. Le, Q. Le, R. Nguyen, and B.-S. Hua, "Single-image hdr reconstruction by multi-exposure generation," in WACV, 2023.
- [27] J. Cai, S. Gu, and L. Zhang, "Learning a deep single image contrast enhancer from multi-exposure images," IEEE Transactions on Image Processing, vol. 27, no. 4, pp. 2049–2062, 2018.
- [28] Y. Endo, Y. Kanamori, and J. Mitani, "Deep reverse tone mapping," ACM Transactions on Graphics (Proc. of SIGGRAPH ASIA 2017), vol. 36, no. 6, Nov. 2017.
- [29] G. Eilertsen, J. Kronander, G. Denes, R. K. Mantiuk, and J. Unger, "Hdr image reconstruction from a single exposure using deep cnns," ACM transactions on graphics (TOG), vol. 36, no. 6, pp. 1-15, 2017.
- [30] M. S. Santos, R. Tsang, and N. Khademi Kalantari, "Single image hdr reconstruction using a cnn with masked features and perceptual loss," ACM Transactions on Graphics, vol. 39, no. 4, 7 2020.
- [31] B. Fei, Z. Lyu, L. Pan, J. Zhang, W. Yang, T. Luo, B. Zhang, and B. Dai, "Generative diffusion prior for unified image restoration and enhancement," in CVPR, 2023.
- [32] D. P. Kingma, M. Welling et al., "An introduction to variational autoencoders," Foundations and Trends® in Machine Learning, 2019.
- [33] A. Lugmayr, M. Danelljan, A. Romero, F. Yu, R. Timofte, and L. Van Gool, "Repaint: Inpainting using denoising diffusion probabilistic models," in CVPR, 2022.
- [34] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, "Places: A 10 million image database for scene recognition," IEEE TPAMI, 2017.
- [35] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," arXiv preprint arXiv:1412.6980, 2014.
- [36] A. Q. Nichol and P. Dhariwal, "Improved denoising diffusion probabilistic models," in ICML, 2021.
- J. Gu, H. Cai, C. Dong, J. S. Ren, R. Timofte, Y. Gong, S. Lao, S. Shi, [37] J. Wang, S. Yang et al., "Ntire 2022 challenge on perceptual image quality assessment," in CVPR, 2022.
- [38] H. Nemoto, P. Korshunov, P. Hanhart, and T. Ebrahimi, "Visual attention in ldr and hdr images," in 9th International Workshop on Video Processing and Quality Metrics for Consumer Electronics (VPQM), 2015.
- [39] W. Wang, Q. Lv, W. Yu, W. Hong, J. Qi, Y. Wang, J. Ji, Z. Yang, L. Zhao, X. Song et al., "Cogvlm: Visual expert for pretrained language models," arXiv preprint arXiv:2311.03079, 2023.
- [40] E. Reinhard and K. Devlin, "Dynamic range reduction inspired by photoreceptor physiology," IEEE transactions on visualization and computer graphics, vol. 11, no. 1, pp. 13-24, 2005.
- [41] A. Mittal, A. K. Moorthy, and A. C. Bovik, "No-reference image quality assessment in the spatial domain," IEEE TIP, 2012.
- [42] A. Mittal, R. Soundararajan, and A. C. Bovik, "Making a "completely blind" image quality analyzer," IEEE Signal processing letters, 2012.
- [43] S. Yang, T. Wu, S. Shi, S. Lao, Y. Gong, M. Cao, J. Wang, and Y. Yang, "Maniqa: Multi-dimension attention network for noreference image quality assessment," in *CVPR*, 2022. [44] J. Wang, K. C. Chan, and C. C. Loy, "Exploring clip for assessing
- the look and feel of images," in AAAI, 2023.
- [45] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," ÎEEE ŤIP, 2004.
- [46] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in CVPR, 2018.
- [47] G. Jinjin, C. Haoming, C. Haoyu, Y. Xiaoxing, J. S. Ren, and D. Chao, "Pipal: a large-scale image quality assessment dataset for perceptual image restoration," in ECCV, 2020.
- [48] Y. Blau and T. Michaeli, "The perception-distortion tradeoff," in CVPR, 2018.

482

545

546