TEXT BOOSTS GENERALIZATION: A PLUG-AND-PLAY CAPTIONER FOR REAL-WORLD IMAGE RESTORATION

Anonymous authors

Paper under double-blind review

Abstract

Generalization has long been a central challenge in real-world image restoration. While recent diffusion-based restoration methods, which leverage generative priors from text-to-image models, have made progress in recovering more realistic details, they still encounter "generative capability inactivation" when applied to out-of-distribution data. To address this, we propose using text as an auxiliary invariant representation to reactivate the generative capabilities of these models. We begin by identifying two key properties of text input in diffusion-based restoration: richness and relevance, and examine their respective influence on model performance. Building on these insights, we introduce Res-Captioner, a module that generates enhanced textual descriptions tailored to image content and degradation levels, effectively mitigating response failures. Additionally, we present RealIR, a new benchmark designed to capture diverse real-world scenarios. Extensive experiments demonstrate that Res-Captioner significantly boosts the generalization ability of diffusion-based restoration models, while remaining fully plug-and-play.



Figure 1: State-of-the-art methods like SUPIR (Yu et al., 2024a) face "generative capability inactivation" on out-of-distribution (OOD) data. Our Res-captioner reactivates their generative capabilities by providing detailed and accurate descriptions.

1 INTRODUCTION

Diffusion-based image restoration methods (Yu et al., 2024a; Sun et al., 2024; Wu et al., 2024; Wang et al., 2024b; Lin et al., 2023; Tao Yang & Zhang, 2023; Ai et al., 2024; Yu et al., 2024b; Zhang et al., 2024), powered by pre-trained text-to-image (T2I) models (Rombach et al., 2022; Podell et al., 2024), achieve superior texture and detail recovery compared to GAN-based methods (Zhang et al., 2021; Wang et al., 2021; Liang et al., 2021; 2022a; Chen et al., 2022). However, these models still

face the out-of-distribution (OOD) challenge (Koh et al., 2021), arising from misalignment between
training data and real-world test cases. Real-world degradation simulations (Zhang et al., 2021;
Wang et al., 2021) offer a common mitigation approach, but a domain gap persists (Liu et al., 2023b;
Wang et al., 2024a; Kong et al., 2022), especially pronounced for device-induced degradations. As
depicted in Figure 1, even state-of-the-art methods struggle to restore fine textures under complex
degradations, a limitation we refer to as "generative capability deactivation".

060 We define image restoration as $x = \mathcal{R}(x_{l_q})$, where x and x_{l_q} denote high-quality (HQ) and low-061 quality (LQ) images, respectively, and \mathcal{R} is the restoration model. To tackle domain generalization, 062 researchers propose learning a cross-domain invariant representation $z = \mathcal{G}(x_{l_a})$ (Arjovsky et al., 063 2019; Nguyen et al., 2021; Li et al., 2022a) and then train a prediction network conditioned on z: 064 $x = \mathcal{H}(z)$. However, learning degradation-invariant representations with strong generalization and minimal information loss remains difficult in image restoration (Liu et al., 2022), as decoupling 065 content from degradation in the image modality is challenging (Chen et al., 2024; Tran et al., 2021; 066 Li et al., 2022b). To address this, we propose transforming LQ images into the text modality using an 067 image captioner C: $y = C(x_{lq})$, leveraging recent multi-modal advancements (Liu et al., 2024b;a; 068 Chen et al., 2023). This approach offers two advantages: first, in the text modality, degradation-069 related descriptions y_{deg} can be easily separated, leaving the content-related part $y_{cont} = \{w \mid w \in$ $y, w \notin y_{deg}$ as a degradation-invariant representation of x_{lg} . Second, text naturally activates priors 071 in T2I diffusion models, facilitating enhanced texture recovery (Yu et al., 2024a; Sun et al., 2024; 072 Wu et al., 2024; Tao Yang & Zhang, 2023; Yu et al., 2024b; Zhang et al., 2024). 073

However, due to significant information compression during the image-to-text transformation, re-074 lying solely on y_{cont} cannot fully meet the high-fidelity requirements of image restoration tasks. 075 Therefore, we utilize y_{cont} as an auxiliary invariant representation in conjunction with the LQ image 076 input, expressed as: $x = \mathcal{R}(x_{lq}, y_{cont})$. In our framework, image restoration is treated as a dual-077 conditioned image generation problem. Compared to the text input y_{cont} , the LQ image x_{lq} serves 078 as a much stronger condition, being more closely aligned with the final output x. However, when 079 the degradation domain of the LQ image shifts, the information that the model can extract from 080 x_{lq} largely decreases, leading to the problem of generative capability deactivation (illustrated in 081 Figure 1). To address OOD data, we propose adaptively enhancing the auxiliary invariant representation y_{cont} through our Restoration Captioner (Res-Captioner), compensating for the information 082 083 loss from x_{lq} due to domain shifts.

To this end, we identify two key properties of text input in T2I diffusion-based restoration models: richness and relevance. Richness is primarily reflected in the length of the text; the more detailed the text, the richer the generated textures. Relevance, on the other hand, measures the correlation between the description and the HQ image content, with higher relevance leading to greater fidelity between the restored image and the ground truth. Building on these properties, we develop Res-Captioner, which is designed to accommodate varying degradation types and image clarity levels. Notably, Res-Captioner can be seamlessly integrated into restoration models, enhancing generalization without requiring retraining of the restoration model itself.

Finally, given the limitations of current real-world image restoration benchmarks (Cai et al., 2019;
 Wei et al., 2020), such as the restricted variety of imaging devices, and narrow content diversity,
 we introduce a new benchmark called **RealIR**. RealIR encompasses a broader range of degradation
 sources, clarity levels, and diverse photographic scenarios. Through this benchmark, we demon strate that our Res-Captioner significantly improves the generalizability of diffusion-based methods,
 delivering more detailed and high-fidelity restoration results.

⁰⁹⁸ The contributions of this paper can be summarized as:

099

102

103

- We identify the potential of utilizing text as an ancillary invariant representation to enhance generalizability in image restoration, highlighting two key properties—richness and relevance—and their respective impacts on restoration performance.
- Building on our findings, we develop the Res-Captioner, which generates adaptively enhanced ancillary invariant representations, improving the generalizability of pre-trained diffusion-based restoration models in a plug-and-play fashion.
- We introduce a new restoration benchmark, RealIR, to comprehensively assess generalizability. Using both our benchmark and existing public datasets, we demonstrate the effectiveness of the Res-Captioner across multiple restoration methods.

108 **RESTORATION CAPTIONER** 2 109

2.1**PROPERTIES OF TEXT INPUT**

We start by investigating how the text input y affects the performance of restoration methods built on text-to-image (T2I) models. We identify two key properties of the text: richness and relevance. Richness refers to the amount of information conveyed, often reflected in text length, while relevance measures the degree of correlation between the text and the corresponding high-quality (HQ) image. Additionally, we observe that degradation-related or photography-specific descriptions can negatively affect restoration results, highlighting the importance of extracting content-specific descriptions, denoted as y_{cont} .



Figure 2: Visualization of the text richness property. (Left) The richness of textures and details in the restored results increases with text richness. Text that is too short can result in the "generative capability inactivation" problem. Excessively long text can lead to messy generation and artifacts. (**Right**) We can classify image content into three categories based on the effect of increased text richness: I beneficial, II insensitive, and III detrimental. 146

147 148

149

142

143

144

145

110

111 112

113

114

115

116

117

118

119

2.1.1 **RICHNESS PROPERTY**

150 **Observation 1**. The richness of restored textures and details increases proportionally with the 151 richness of the text description.

152 As illustrated in Figure 2, we observe that for all low-quality (LQ) images, increasing text richness 153 (*i.e.*, text length) consistently enhances texture restoration. To explore this further, we prepare a 154 dataset of 120 HQ images from diverse scenarios and generate the corresponding LQ images using 155 Real-ESRGAN (Wang et al., 2021). GPT-4 is employed to generate detailed descriptions for the 156 HQ images. We then evaluate two representative restoration models, SUPIR (Yu et al., 2024a) and 157 StableSR (Wang et al., 2024b), for verification. The text input is encoded using CLIP (Radford et al., 158 2021) to generate 77 tokens. We then repeatedly append the last 20 tokens, excluding the EOS token, and follow (Xia et al., 2024) to integrate and inject these length-varying tokens into the restoration 159 models to produce the restored results. Texture richness is assessed using two non-reference metrics: 160 MANIQA (Yang et al., 2022) and MUSIQ (Ke et al., 2021). As shown in Figure 3 (a, b), both metrics 161 demonstrate a positive correlation with the number of text tokens, supporting **Observation 1**.



Figure 3: Demonstration of the richness property. (a, b): There is a positive correlation between 170 text richness and the richness of textures in the restored results. (c, d): The optimal text richness (indicated by an asterisk) is proportional to the degree of deviation between the test degradation 172 domain and the training degradation domain. Best viewed zoomed in.

171

175 We attribute this property to the data bias inherent in pre-trained T2I models, where images with 176 richer content are typically paired with more detailed descriptions during training. Similar observa-177 tions have been made in T2I research (Betker et al., 2023; Yang et al., 2024), where longer prompts 178 lead to more enriched scenes. However, in the context of image restoration, this effect primarily 179 enhances texture quality rather than introducing new objects or elements.

180 **Observation 2.** The optimal level of text richness is influenced by factors such as degradation 181 severity, and image content. 182

As discussed, detailed text descriptions improve texture restoration. However, as shown in Figure 2, 183 exceeding the optimal range of text richness may lead to undesirable artifacts or messy generation. 184 For instance, the squirrel's eyes and mouth are misaligned with the LQ image, and the bee shows 185 over-sharpening effects. We posit that the optimal text richness is proportional to the domain gap between training and testing degradations. To validate this, we prepare LQ images either simulated 187 or captured in the wild with different zoom ratios and evaluate the performance of SUPIR and 188 StableSR in relation to text richness. As illustrated in 3 (c, d), as the test degradation increasingly 189 diverges from the training setting (e.g., $4 \times$ Real-ESRGAN degradation), the optimal text richness 190 similarly increases. This is because, as degradation severity intensifies, the useful information the 191 model can extract from LQ images diminishes, necessitating more informative textual inputs to compensate for the information loss. 192

193 We also find that the optimal text richness is influenced by the content of the LQ image. Follow-194 ing (Liang et al., 2022b), we categorize three groups based on the impact of increased text richness 195 on image content: beneficial, insensitive, and detrimental. Category I, "beneficial", includes fine-196 grained textures (e.g., feathers, leaves, sand) and regular structures (e.g., walls, windows), which benefit from longer text input as it activates the model's generative capability. Category II, "in-197 sensitive", consists of smooth areas and large-scale structures (e.g., sky), where text richness has 198 minimal effect. Category III, "detrimental", includes non-rigid structures (e.g., text, crowds), where 199 excessively long text may compromise fidelity. 200

201 202

203

204

2.1.2**RELEVANCE PROPERTY**

205 **Observation 3**. The fidelity of restored textures improves in correlation with the relevance of the 206 text description.

207 To quantitatively characterize the text relevance property, we introduce the concept of the "text-208 replacing ratio". This is defined as the ratio of original words in the text input y that are replaced 209 with non-meaningful words like "the" or "for". As the text-replacing ratio increases, the relevance 210 between the text input and the corresponding HQ image decreases, while the text richness remains 211 unchanged. As shown in Figure 4, we observe that although restored results retain rich textures as 212 the text-replacing ratio increases, they suffer from a decline in fidelity. This is confirmed by the 213 decreasing DISTS scores (lower is better), measured between the HQ image and the restored output. At higher text-replacing ratios, even the overall semantics of the restoration become distorted. For 214 instance, at a ratio of 1.0, a lizard's head is incorrectly restored as flowers. In such cases, the model 215 continues to generate textures but lacks the appropriate guidance to produce accurate ones.

216 217 218 219 220 221 222 222 223 LQ 0.0 Control L

Figure 4: Visualization and demonstration of the text relevance property. Left: The accuracy of textures and details in the restored results decreases as the text-replacing ratio increases, indicating that text relevance contributes to the fidelity of the restoration. **Right**: DISTS increases with a higher text-replacing ratio, further indicating a decrease in the fidelity of the restored results.

2.1.3 HARMFUL DESCRIPTION

Observation 4. Descriptions related to degradation or photography can lead to global or localized
 blurring in the restored images.

233 We discover that degradation-related descriptions such as "blur" or "blurred", and photography 234 terms like "shallow depth of field" or "bokeh effect", may lead to blurred outputs. Even when 235 descriptions like "the background is blurred, while the main subject is sharp" accurately reflect the 236 HQ image, they can cause overall blurring in the restored results. This is likely due to the limited spatial control capabilities of pre-trained T2I models (Avrahami et al., 2023), which amplifies 237 the blurring effect. To validate this, we use GPT-4 to generate two captions of similar length: one 238 without harmful descriptions and another including them. To exclude the effects of text richness 239 and relevance, we duplicate the harmless description, labeled "Without Harmful Description", and 240 combine both harmless and harmful descriptions to create "With Harmful Description". As shown 241 in Figure 5, the description without harmful terms successfully restores clearer and richer details, 242 while the harmful description does not. 243





2.2 ANCILLARY INVARIANT REPRESENTATION ENHANCEMENT

258 Learning degradation-invariant representations from LQ images is highly challenging. To address 259 this, we propose using text free of degradation-related descriptions as an auxiliary invariant represen-260 tation to improve generalization. As discussed in Section 2.1, text plays a crucial role in controlling 261 both the richness and fidelity of textures in restored results. However, existing image captioners (Liu et al., 2024b;a; Chen et al., 2023), which are not specifically designed for image restoration, not only 262 generate harmful descriptions but also fail to adaptively enhance text richness. Consequently, they 263 may contribute to the "generative capability inactivation" problem (Figure 1) in real-world scenar-264 ios. To address this issue, we introduce Res-Captioner, a restoration-specific captioner that generates 265 high-quality text descriptions for real-world LQ images across diverse degradation levels and con-266 tent categories, ensuring adaptive control over both richness and relevance. 267

Training data generation. We first collect HQ images from (Unsplash), ImageNet (Deng et al., 2009), and SAM (Kirillov et al., 2023), filtering out overly smooth ones using Sobel filters based on image gradient standard deviation. This ensures a selection of rich-content, high-clarity HQ images



257

244

245

246

247

248

249 250

251

224

225

226

227 228 229



Figure 6: (a) The generation and annotation process of our training data. (b) Chain-of-Thought captioning of our Res-Captioner. (c) Network structure of our Res-Captioner.

from diverse scenarios. Next, as shown in Figure 6 (a), we leverage five pre-trained latent diffusion
 models (LDM) (Rombach et al., 2022) to generate LQ images that simulate varying imaging devices
 and zoom ratios. Training details for the LDM are provided in the appendix. We also include a
 percentage of Real-ESRGAN-generated (Wang et al., 2021) LQ images to further enhance diversity.

To ensure high relevance while minimizing hallucination, as illustrated in Figure 6 (a), we use GPT-4 to generate descriptions of varying lengths for each HQ image. Several prompting techniques, detailed in the appendix, are applied to avoid degradation-related or photography-specific content. These descriptions are fed into the restoration model, producing multiple restored candidates for each LQ image. Human annotators select the optimal text input that provides the best visual result, balancing texture richness and fidelity. The token length of the selected description is then calculated and combined with the description to form the final caption output in the format <token length, description>. In total, we curate 5,500 LQ image-caption pairs for training our Res-Captioner.

Chain-of-Thought captioning. Our goal is to generate accurate descriptions with appropriate rich-298 ness for LQ images. As discussed in **Observation 1** and **Observation 2**, the optimal text richness, 299 primarily reflected in token length, is determined by image content and degradation levels. Given 300 the diversity of real-world scenarios, we enhance the reasoning and decision-making capabilities 301 of Multimodal Large Language Models (MLLM) by adopting the "Chain of Thought" (CoT) strat-302 egy (Wei et al., 2022) in Res-Captioner. Specifically, as shown in Figure 6 (b), the model first 303 predicts the optimal token number before generating the corresponding caption. As demonstrated in 304 Section 3.4, this approach significantly improves the accuracy of the description length. 305

Network structure. We fine-tune LLaVA-1.5 (Liu et al., 2024a) using low-rank adaptation 306 (LoRA) (Hu et al., 2021) to serve as our Res-Captioner. Since LLaVA is not designed for LQ 307 images, we enhance its ability to perceive image degradations. In addition to the original CLIP 308 visual encoder, we incorporate a degradation-aware visual encoder, as shown in Figure 6 (c). This 309 encoder consists of a pre-trained degradation extractor, known for its sensitivity to various degra-310 dations (Chen et al., 2024; Liu et al., 2023a), and a lightweight adapter for improved degradation 311 extraction. Specifically, the adapter is built from several MLP layers, first compressing the token 312 count to 1 and then expanding it to N tokens (we set N = 36), enabling the encoder to focus on the 313 global degradation representation while ignoring spatially varying content.

314 315

316

283

284

2.3 REALIR BENCHMARK

The current real-world restoration benchmarks (Cai et al., 2019; Wei et al., 2020) are limited by a narrow range of degradation types, insufficient diversity in imaging devices, and constrained content scope. To overcome these limitations, we introduce RealIR, a new benchmark featuring 152 real LQ images from eight imaging devices, including two DSLRs and six mobile phones, capturing images with varying zoom ratios. We also incorporate 53 LQ images sourced from the internet to capture degradations introduced by network transmission, which differ from device-specific degradations. The dataset covers a wide range of content, including portraits, animals, plants, and architectural scenes, enabling comprehensive evaluations of image restoration methods' generalizability.

324 3 **EXPERIMENTS** 325

326 3.1 IMPLEMENTATION DETAILS

328 Our Res-Captioner is built on LLaVA- 1.5^{1} . We train the model with a batch size of 128 over 500 329 steps using an A800 GPU, employing the Adam optimizer (Kingma & Ba, 2014) with a learning rate of 2×10^{-4} . We integrate Res-Captioners into two diffusion-based restoration models: SUPIR (Yu 330 et al., 2024a), built on SDXL (Podell et al., 2024), and StableSR (Wang et al., 2024b), using Stable 331 Diffusion 2.1 (Rombach et al., 2022). Our models operate in a plug-and-play fashion, seamlessly 332 integrating with restoration models based on the same text-to-image (T2I) backbone. As in (Xia 333 et al., 2024), we iteratively process the long text through the CLIP (Radford et al., 2021) text encoder. 334

335 Details of our training data generation and labeling process are provided in Section 2.2. We col-336 lect 5,500 low-quality (LQ) image-caption training pairs for SUPIR. Recognizing that different T2I backbones exhibit distinct text richness characteristics, we collect an additional 500 pairs for Sta-337 bleSR for fine-tuning. To match the resolution requirements of the respective T2I backbones, we 338 resize the short edge of high-quality (HQ) images to 1024 for SUPIR and 512 for StableSR. The 339 parameters for LoRA follow the standard LLaVA settings. 340

341 342

327

3.2 EXPERIMENTAL SETTINGS

343 Test datasets. Our proposed RealIR dataset encompasses diverse content and degradations from 344 real-world scenes, making it ideal for assessing the generalization ability of restoration models. 345 However, due to the absence of ground-truth images in RealIR, we create an additional multi-346 degradation test set comprising 120 LQ-HQ pairs using pre-trained latent diffusion models (LDM). 347 To ensure fair evaluation, the degradations used for LQ generation are distinct from those in our 348 training set. We categorize the LQ-HQ pairs into three degradation levels based on zoom ratio: light 349 (zoom ratio of 3 to 7), moderate (zoom ratio of 8 to 10), and heavy (zoom ratio of 15 to 20). These 350 two test sets enable a comprehensive evaluation of the restored results' detail richness and fidelity across varying degradation levels. Furthermore, we evaluate our approach on established real-world 351 benchmarks such as RealSR (Cai et al., 2019) and DRealSR (Wei et al., 2020), using randomly 352 cropped patches for more comprehensive analysis. 353

354 Compared methods. Our experiments include state-of-the-art (SOTA) real-world image restora-355 tion methods, such as GAN-based approaches like Real-ESRGAN+ (Wang et al., 2021) and 356 DASR (Liang et al., 2022a), as well as diffusion-based models including StableSR (Wang et al., 2024b), SeeSR (Wu et al., 2024), CoSeR (Sun et al., 2024), and SUPIR (Yu et al., 2024a). Addi-357 tionally, we compare our Res-Captioner with leading image captioning models, such as LLaVA-1.5 358 and ShareCaptioner (Chen et al., 2023). 359

360 Evaluation metrics. For test sets without ground truth, such as RealIR, we use non-reference eval-361 uation metrics aligned with human perception, including MUSIQ (Ke et al., 2021), MANIOA (Yang 362 et al., 2022), LIQE (Zhang et al., 2023), and NIQE (Zhang et al., 2015). For datasets with ground truth, we adopt perceptual distance metrics like DISTS Ding et al. (2020) and LPIPS Zhang et al. 363 (2018), alongside the LIQE metric, which leverages large vision-language models for robust evalu-364 ation. Pixel-level metrics such as PSNR and SSIM are no longer considered, as they exhibit weak 365 correlation with human perception, as discussed in related works (Yu et al., 2024a; Sun et al., 2024). 366

367 368

3.3 COMPARISON WITH STATE OF THE ARTS

369 3.3.1 QUANTITATIVE RESULTS 370

371 Our quantitative results are organized into three parts. First, we assess the generalization ability of 372 existing restoration methods in real-world scenarios using the RealIR benchmark, showing that Res-373 Captioner consistently improves their performance. Second, we evaluate the multi-degradation test 374 set and the existing benchmarks, confirming that Res-Captioner not only enhances detail generation 375 but also preserves fidelity across various degradation levels. Lastly, we compare the performance of 376 other image captioners to our Res-Captioner, highlighting its superior effectiveness.

¹https://huggingface.co/liuhaotian/llava-v1.5-13b

Methods		RealIR (Car	neras)		RealIR (Internet)			
Wellous	MUSIQ↑	MANIQA↑	LIQE↑	NIQE↓	MUSIQ↑	MANIQA↑	LIQE↑	NIQE↓
Real-ESRGAN+	58.54	0.1784	2.425	5.646	58.34	0.2048	2.157	5.646
DASR	53.82	0.1487	2.208	6.748	50.84	0.1397	1.594	6.748
CoSeR	56.91	0.1163	2.597	4.042	66.67	0.1842	3.822	4.042
SeeSR	70.19	0.2138	3.768	3.749	72.65	0.2694	4.243	3.749
StableSR	66.15	0.1924	3.466	4.033	67.66	0.2012	3.913	4.033
StableSR w/ Ours	68.04	0.1955	3.615	3.888	70.90	0.2251	4.252	3.888
SUPIR	60.43	0.1651	2.983	3.492	71.94	0.2727	4.425	3.492
SUPIR w/ Ours	71.38	0.2543	4.056	3.389	73.26	0.3055	4.578	3.389

378 Table 1: Quantitative comparisons on our RealIR benchmark. We highlight best values for each 379 metric and the results of Res-Captioner-enhanced models .

Table 2: Quantitative comparisons between the official model and the Res-Captioner-enhanced model under different degradation levels. We show the improvement percentage on each metric.

Methods	Light Degradation			Moderate Degradation			Heavy Degradation		
methods	DISTS↓	LPIPS↓	LIQE↑	DISTS↓	LPIPS↓	LIQE↑	DISTS↓	LPIPS↓	LIQE↑
StableSR StableSR w/ Ours	0.1791 0.1661 7.3%	0.3311 0.3222 2.7%	3.699 3.735 1.0%	0.1864 0.1692 9.2%	0.3209 0.3086 3.8%	3.603 3.857 7.1%	0.2181 0.1918 12.1%	0.4008 0.3773 5.9%	3.047 3.604 18.3%
SUPIR SUPIR w/ Ours	0.1821 0.1680 7.7%	0.3444 0.3178 7.7%	3.148 4.011 27.4%	0.1883 0.1621 13.9%	0.3473 0.3052 12.1%	3.349 4.226 26.2%	0.2159 0.1873 13.3%	0.4106 0.3754 8.6%	2.840 3.991 40.5%

Table 3: Quantitative comparisons on RealSR and DRealSR T datasets. Bold: Best results. 0

able 4:	Quantitative comparisons
f image	captioners on restoration.

Methods	RealSR		DRealSR			Methods	DISTS↓	LPIPS↓	
	DISTS↓	LPIPS↓	LIQE↑	DISTS↓	LPIPS↓	LIQE↑	Llava-1.5	0.1886	0.3600
SUPIR SUPIR w/ Ours	0.2660 0.2474	0.3889 0.3667	3.477 4.081	0.2906 0.2699	0.4741 0.4409	3.655 4.208	ShareCaptioner Res-Captioner	0.1780 0.1725	0.3394 0.3328

406 407 408

390

391

392 393

402

403 404 405

409 **RealIR benchmark.** The results, shown in Table 1, evaluate both real LQ images captured by 410 various cameras and LQ images collected from the internet. Overall, diffusion-based models exhibit superior visual quality compared to GAN-based models, due to their stronger generative capabilities. 411 Notably, when integrated with our Res-Captioner, diffusion models such as StableSR and SUPIR 412 show significant improvements across all metrics. This highlights how our approach fully activates 413 the generative power of T2I-based restoration models for diverse real-world LQ images. 414

415 However, the improvement introduced by Res-Captioner varies across different restoration models. 416 For example, Res-Captioner enhances the LIQE score of StableSR on manually captured RealIR data by approximately 4.3%, while it increases the LIQE score of SUPIR by an impressive 36%. We 417 attribute this discrepancy to the differing generative capabilities of T2I models. In particular, SUPIR, 418 based on SDXL, suffers from the "generative capability deactivation" issue, which is effectively 419 reactivated by our Res-Captioner, unlocking its full potential. 420

421 Fidelity evaluation. We compare the original models with their Res-Captioner-enhanced versions 422 on the multi-degradation test set, as shown in Table 2. For both StableSR and SUPIR, Res-Captioner 423 consistently improves fidelity, demonstrated by significant gains in reference-based metrics like DISTS and LPIPS. Notably, the performance improvements increase with the severity of degra-424 dation. For instance, the DISTS score of enhanced StableSR improves by approximately 7.3%, 425 9.2%, and 12.1% for light, moderate, and heavy degradation, respectively. This trend supports our 426 approach of using text as an auxiliary invariant representation. As test degradation diverges further 427 from the training distribution, the restoration model extracts less useful information from the LQ 428 image, making the supplementary text provided by Res-Captioner increasingly beneficial. 429

Given the relatively simple and light degradation in RealSR and DRealSR, we use SUPIR as the ref-430 erence model for evaluation. Our Res-Captioner significantly improves the performance of SUPIR 431 in Table 3, further demonstrating the robustness of our approach in real-world scenarios.





Figure 7: Qualitative comparisons on in-the-wild images. **Upper**: Comparisons between SOTA restoration methods and Res-Captioner-enhanced SUPIR. **Lower**: Visual quality improvements introduced by Res-Captioner on StableSR and SeeSR.

Comparison of image captioners. We compare our Res-Captioner to LLaVA-1.5 and ShareCap-tioner in a plug-and-play manner, integrating all captions into SUPIR as described in Section 3.1. Results from the multi-degradation test set, shown in Table 4, demonstrate that Res-Captioner pro-vides superior guidance for image restoration. LLaVA-1.5 typically generates shorter captions (aver-age length of 80), while ShareCaptioner consistently produces overly long captions (average length of 200). As noted in **Observation 2**, both overly short and excessively long captions can negatively affect restoration results. In contrast, Res-Captioner dynamically adjusts text richness based on the input image, optimizing restoration quality across varying degradation levels.

Beyond text richness, the generated descriptions differ significantly in quality. As discussed in
 the appendix, other captioners may produce misleading descriptions or hallucinations that degrade
 restoration quality, while our method generates highly relevant, accurate descriptions aligned with
 HQ images, effectively enhancing restoration results.

3.3.2 QUALITATIVE RESULTS

We provide visual comparisons on in-the-wild LQ images in Figure 7. In the upper section, Real-ESRGAN+ struggles with its limited generative capability, failing to recover high-definition textures. Both SUPIR and StableSR experience "generative capability deactivation" when handling
out-of-distribution (OOD) data, leading to large areas of blurring. Although SeeSR responds better
to OOD data, the textures it generates tend to appear overly smooth and unrealistic. In contrast, our
Res-Captioner fully activates the generative potential of the T2I backbone in SUPIR, enabling the
recovery of clearer, more realistic textures, such as detailed flower petals and building facades.

The lower section of Figure 7 illustrates how Res-Captioner improves other restoration models.
When integrated with Res-Captioner, StableSR and SeeSR demonstrate an enhanced ability to recover fine-grained textures and structures, such as goat fur and lantern mesh, significantly outperforming their original versions. Notably, Res-Captioner can be directly applied to restoration models using the same T2I backbone without requiring fine-tuning.

Method	Light Deg	gradation	Moderate	Degradation	Heavy Degradation	
Wiethou	DISTS↓	LPIPS↓	DISTS↓	LPIPS↓	DISTS↓	LPIPS↓
Ours	0.1680	0.3178	0.1621	0.3052	0.1873	0.3754
w/ Min Len.	0.1718	0.3274	0.1753	0.3252	0.2033	0.4009
w/ Max Len.	0.1864	0.3525	0.1770	0.3184	0.1964	0.4039
w/ Low Rel.	0.1738	0.3389	0.1655	0.3061	0.1907	0.3914
w/ Harmful Des.	0.1686	0.3191	0.1678	0.3178	0.1868	0.3883





494 495

496

3.3.3 USER STUDY

497 To further validate Res-Captioner's ability to enhance generalization in real-world scenarios, we 498 conduct a user study on in-the-wild LQ images with 31 experienced researchers. Each participant 499 rates the visual perceptual quality (on a scale of 1 to 10, where higher is better) of results gener-500 ated by Real-ESRGAN+, SeeSR, StableSR, StableSR with Res-Captioner, SUPIR, and SUPIR with 501 Res-Captioner. As illustrated in Figure 8, StableSR and SUPIR show significant performance im-502 provements when paired with Res-Captioner. Notably, SUPIR, when enhanced with Res-Captioner, 503 delivers the highest visual quality among all methods.

504 505

506

3.4 ABLATION STUDY

507 We investigate the impact of the proposed text properties—richness, relevance, and harmful descrip-508 tions—on restoration performance by ablating each aspect in experiments. All models are trained 509 under identical settings, with the only variation being the training data. Additionally, we analyze the effect of our proposed Chain-of-Thought (CoT) captioning and degradation-aware visual encoder 510 on text richness. SUPIR is used as the restoration model in this section. 511

512 **Text richness.** To explore the impact of text richness, we create two training sets using the shortest 513 and longest captions generated by GPT-4, corresponding to the results of "w/ Min Len." and "w/ 514 Max Len." in Table 5. The results show that Res-Captioner achieves the best performance under 515 varying degradation conditions due to its adaptive text richness capability. Moreover, we observe that the "w/ Max Len." model begins to outperform the "w/ Min Len." model as degradation severity 516 increases, which is consistent with our **Observation 2**. 517

518 **Text relevance.** To study this property, we first calculate the length of human-selected optimal 519 captions generated by GPT-4. We then produce same-length low-relevance captions using LLaVA-520 1.5 for training, denoted as "w/ Low Rel.". In contrast, our Res-Captioner ("Ours") achieves superior 521 restoration results, highlighting the importance of high-relevance descriptions for restoration.

522 **Harmful descriptions.** In Section 2.1.3, we identify harmful descriptions that result in blurring 523 in the restored images. Using the optimal text richness, we employ GPT-4 to generate captions 524 incorporating these harmful descriptions. We fine-tune our Res-Captioner with this data, referred 525 to as "w/ Harmful Des." in Table 5. The results show that harmful descriptions negatively affect 526 restoration performance, causing an average 2.7% decrease in LPIPS.

527 **CoT captioning and degradation-aware visual encoder.** We manually annotate the optimal text 528 length, L_o , for 100 LQ images from the RealIR and multi-degradation datasets. To quantify the 529 prediction error, we define the offset level E as: $E = \max(|L_o - L| - 15, 0)/30$, where L is the 530 captioner's output length. The mean offset level for our Res-Captioner is 1.27. Without the CoT 531 captioning, the mean offset level increases by 66.7%, and without the degradation-aware visual 532 encoder, it rises by **31.5%**. These results highlight the effectiveness of our model's design.

533 534

4 CONCLUSION

535 536

537 We leverage text as an auxiliary invariant representation to enhance the generalizability of T2Idiffusion-based restoration models. By focusing on two key properties of text inputs-richness and 538 relevance—we propose Res-Captioner, which significantly improves real-world restoration performance in a plug-and-play manner.

540 REFERENCES

547

551

552

553

- Yuang Ai, Huaibo Huang, Xiaoqiang Zhou, Jiexiang Wang, and Ran He. Multimodal prompt per ceiver: Empower adaptiveness generalizability and fidelity for all-in-one image restoration. In
 CVPR, pp. 25432–25444, 2024.
- Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization.
 arXiv preprint arXiv:1907.02893, 2019.
- Omri Avrahami, Thomas Hayes, Oran Gafni, Sonal Gupta, Yaniv Taigman, Devi Parikh, Dani Lischinski, Ohad Fried, and Xi Yin. Spatext: Spatio-textual representation for controllable image generation. In *CVPR*, pp. 18370–18380, 2023.
 - James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science. https://cdn. openai. com/papers/dall-e-3. pdf*, 2(3):8, 2023.
- 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 *ICCV*, pp. 3086–3095, 2019.
- Chaofeng Chen, Xinyu Shi, Yipeng Qin, Xiaoming Li, Xiaoguang Han, Tao Yang, and Shihui Guo.
 Real-world blind super-resolution via feature matching with implicit high-resolution priors. In
 ACMMM, pp. 1329–1338, 2022.
- Haoyu Chen, Wenbo Li, Jinjin Gu, Jingjing Ren, Haoze Sun, Xueyi Zou, Zhensong Zhang, Youliang
 Yan, and Lei Zhu. Low-res leads the way: Improving generalization for super-resolution by self supervised learning. In *CVPR*, pp. 25857–25867, 2024.
- Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua
 Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
 hierarchical image database. In *CVPR*, pp. 248–255, 2009.
- Keyan Ding, Kede Ma, Shiqi Wang, and Eero P Simoncelli. Image quality assessment: Unifying structure and texture similarity. *TPAMI*, 44(5):2567–2581, 2020.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Junjie Ke, Qifei Wang, Yilin Wang, Peyman Milanfar, and Feng Yang. Musiq: Multi-scale image quality transformer. In *ICCV*, pp. 5148–5157, 2021.
- 578 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint* 579 *arXiv:1412.6980*, 2014.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *ICCV*,
 pp. 4015–4026, 2023.
- Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Bal-subramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Irena Gao, et al. Wilds: A benchmark of in-the-wild distribution shifts. In *ICML*, pp. 5637–5664. PMLR, 2021.
- Xiangtao Kong, Xina Liu, Jinjin Gu, Yu Qiao, and Chao Dong. Reflash dropout in image super resolution. In *CVPR*, pp. 6002–6012, 2022.
- Bo Li, Yifei Shen, Yezhen Wang, Wenzhen Zhu, Dongsheng Li, Kurt Keutzer, and Han Zhao. In-variant information bottleneck for domain generalization. In *AAAI*, volume 36, pp. 7399–7407, 2022a.
- 593 Dasong Li, Yi Zhang, Ka Chun Cheung, Xiaogang Wang, Hongwei Qin, and Hongsheng Li. Learning degradation representations for image deblurring. In ECCV, pp. 736–753. Springer, 2022b.

- 594 Jie Liang, Hui Zeng, and Lei Zhang. Efficient and degradation-adaptive network for real-world image super-resolution. In ECCV, pp. 574–591. Springer, 2022a. 596 Jie Liang, Hui Zeng, and Lei Zhang. Details or artifacts: A locally discriminative learning approach 597 to realistic image super-resolution. In CVPR, pp. 5657-5666, 2022b. 598 Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: 600 Image restoration using swin transformer. In ICCVW, pp. 1833–1844, 2021. 601 602 Xinqi Lin, Jingwen He, Ziyan Chen, Zhaoyang Lyu, Bo Dai, Fanghua Yu, Wanli Ouyang, Yu Qiao, 603 and Chao Dong. Diffbir: Towards blind image restoration with generative diffusion prior. arXiv preprint arXiv:2308.15070, 2023. 604 605 Anran Liu, Yihao Liu, Jinjin Gu, Yu Qiao, and Chao Dong. Blind image super-resolution: A survey 606 and beyond. TPAMI, 45(5):5461-5480, 2022. 607 608 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In CVPR, pp. 26296-26306, 2024a. 609 610 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *NeurIPS*, 611 36, 2024b. 612 613 Yihao Liu, Jingwen He, Jinjin Gu, Xiangtao Kong, Yu Qiao, and Chao Dong. Degae: A new 614 pretraining paradigm for low-level vision. In CVPR, pp. 23292-23303, 2023a. 615 Yihao Liu, Hengyuan Zhao, Jinjin Gu, Yu Qiao, and Chao Dong. Evaluating the generalization 616 ability of super-resolution networks. TPAMI, 2023b. 617 618 A Tuan Nguyen, Toan Tran, Yarin Gal, and Atilim Gunes Baydin. Domain invariant representation 619 learning with domain density transformations. NeurIPS, 34:5264–5275, 2021. 620 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe 621 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image 622 synthesis. In ICLR, 2024. 623 624 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 625 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 626 models from natural language supervision. In ICML, pp. 8748–8763. PMLR, 2021. 627 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-628 resolution image synthesis with latent diffusion models. In CVPR, pp. 10684–10695, 2022. 629 630 Haoze Sun, Wenbo Li, Jianzhuang Liu, Haoyu Chen, Renjing Pei, Xueyi Zou, Youliang Yan, and 631 Yujiu Yang. Coser: Bridging image and language for cognitive super-resolution. In CVPR, pp. 632 25868-25878, 2024. 633 Peiran Ren Xuansong Xie Tao Yang, Rongyuan Wu and Lei Zhang. Pixel-aware stable diffusion for 634 realistic image super-resolution and personalized stylization. In ECCV, 2023. 635 636 Phong Tran, Anh Tuan Tran, Quynh Phung, and Minh Hoai. Explore image deblurring via encoded 637 blur kernel space. In CVPR, pp. 11956–11965, 2021. 638 639 Unsplash. Unsplash dataset. URL https://unsplash.com/data. 640 Hongjun Wang, Jiyuan Chen, Yinqiang Zheng, and Tieyong Zeng. Navigating beyond dropout: An 641 intriguing solution towards generalizable image super resolution. In CVPR, pp. 25532–25543, 642 2024a. 643 644 Jianyi Wang, Zongsheng Yue, Shangchen Zhou, Kelvin CK Chan, and Chen Change Loy. Exploiting 645 diffusion prior for real-world image super-resolution. IJCV, pp. 1–21, 2024b. 646
- 647 Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world blind super-resolution with pure synthetic data. In *ICCVW*, pp. 1905–1914, 2021.

648 649 650	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. <i>NeurIPS</i> , 35: 24824–24837, 2022.
652 653 654	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 <i>ECCV</i> , pp. 101–117. Springer, 2020.
655 656	Rongyuan Wu, Tao Yang, Lingchen Sun, Zhengqiang Zhang, Shuai Li, and Lei Zhang. Seesr: Towards semantics-aware real-world image super-resolution. In <i>CVPR</i> , pp. 25456–25467, 2024.
658 659	Bin Xia, Shiyin Wang, Yingfan Tao, Yitong Wang, and Jiaya Jia. Llmga: Multimodal large language model based generation assistant. 2024.
660 661 662	Ling Yang, Zhaochen Yu, Chenlin Meng, Minkai Xu, Stefano Ermon, and CUI Bin. Mastering text-to-image diffusion: Recaptioning, planning, and generating with multimodal llms. In <i>ICML</i> , 2024.
663 664 665 666	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 <i>CVPR</i> , pp. 1191–1200, 2022.
667 668 669	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. In <i>CVPR</i> , pp. 25669–25680, 2024a.
670 671 672	Yongsheng Yu, Ziyun Zeng, Hang Hua, Jianlong Fu, and Jiebo Luo. Promptfix: You prompt and we fix the photo. <i>arXiv preprint arXiv:2405.16785</i> , 2024b.
673 674	Kai Zhang, Jingyun Liang, Luc Van Gool, and Radu Timofte. Designing a practical degradation model for deep blind image super-resolution. In <i>ICCV</i> , pp. 4791–4800, 2021.
675 676 677	Lin Zhang, Lei Zhang, and Alan C Bovik. A feature-enriched completely blind image quality evaluator. <i>IEEE Transactions on Image Processing</i> , 24(8):2579–2591, 2015.
678 679	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i> , pp. 586–595, 2018.
680 681 682	Weixia Zhang, Guangtao Zhai, Ying Wei, Xiaokang Yang, and Kede Ma. Blind image quality assessment via vision-language correspondence: A multitask learning perspective. In <i>CVPR</i> , pp. 14071–14081, 2023.
684 685 686	Yuhong Zhang, Hengsheng Zhang, Xinning Chai, Zhengxue Cheng, Rong Xie, Li Song, and Wenjun Zhang. Diff-restorer: Unleashing visual prompts for diffusion-based universal image restoration. <i>arXiv preprint arXiv:2407.03636</i> , 2024.
687	
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
699	
700	

702 A APPENDIX

705 A.1 REAL-WORLD LQ IMAGE GENERATION

By reproducing high-definition images, we collect numerous real-world LQ-HQ pairs for training
the real-world LQ generation model. Data is gathered from five different devices, and five LQ
generation models are trained to represent different types of degradation.

We select the latent diffusion model (LDM) (Rombach et al., 2022) as our LQ generator, training it to produce LQ images conditioned on corresponding HQ images. Additionally, the zoom ratio used during image reproduction is incorporated as another part of the conditional information.

For each degradation model, we retain one zoom ratio for the generation of multi-degradation test set, and the rest are used to generate Res-Captioner training data.

715 716

718

704

717 A.2 DETAILS OF TRAINING DATA GENERATION

We use the following prompt to generate captions of varying lengths with GPT-4, while avoiding harmful descriptions through the use of restrictive phrasing.

721 722 Please describe the actual objects in the image in a very detailed 723 manner. Please do not include descriptions related to the focus 724 and bokeh of this image. Please do not include descriptions like 725 to about XXX words.

We generate a total of seven different caption lengths: 80 words, 110 words, 140 words, 200 words, 260 words, 350 words, and 440 words. The interval between lengths increases progressively, as longer captions tend to cause smaller texture changes when recovering with the same richness interval.

731 When training and testing the Res-Captioner, we use the following prompt:

732 Please determine the appropriate caption length and then describe 733 the actual objects in the image in a very detailed manner. Please 734 do not include descriptions related to the focus and bokeh of this 735 image. Please do not include descriptions like the background is 736 blurred.

737 738

739 740

A.3 CONTENT QUALITY COMPARISON OF IMAGE CAPTIONS

In Figure 10 and 11, we compare the caption content quality between our proposed Res-Captioner and existing image captioners. Each caption was carefully examined, with hallucinations and harmful descriptions that could negatively impact restoration clearly marked. We also compared captions generated by Res-Captioner and other image captioners for the same image, highlighting the details missing in the others. Our method clearly demonstrates fewer hallucinations, is free from harmful descriptions, and produces a significant amount of detail closely aligned with the image content, which greatly supports the restoration process.

747 748

A.4 MORE QUALITATIVE RESULTS

749 750

In this section, we provide additional visual comparisons between our method and state-of-the-art (SOTA) methods. As illustrated in Figure 12, when paired with the SUPIR restoration method, which features a powerful generative model backbone, our Res-Captioner shows clear advantages in recovery performance compared to previous SOTA methods. Additionally, our approach significantly improves visual quality when applied to the StableSR restoration method, as demonstrated in Figure 9, highlighting the robustness of our approach across different restoration models.

LQ StableSR StableSR w/ Ours StableSR StableSR w/ Ours LQ

Figure 9: Additional qualitative comparisons of Res-Captioner applied to StableSR on in-the-wild images.

A.5 ANALYSIS ON DEGRADATION TOKEN NUMBER

We analyze the impact of the token length N in our degradation-aware visual encoder on text rich-ness prediction accuracy. As shown in Table 6, the token length is inversely proportional to the text richness offset level E, indicating a positive relationship with text richness prediction accuracy. Balancing computational cost with performance benefits, we select N = 36 as the final degradation token number.

Table 6: Analysis of the token length in the degradation-aware visual encoder.

Token Length N	Offset Level E
4	1.53
9	1.46
16	1.34
25	1.28
36	1.27





use red to indicate some hallucinations and harmful descriptions in the caption. We use green to
 highlight the detailed descriptions provided by one captioner that are missing in the other.



Figure 12: Additional qualitative comparisons with SOTA methods on in-the-wild images.