



JarvisIR: Elevating Autonomous Driving Perception with Intelligent Image Restoration

Yunlong Lin^{1*♣} Zixu Lin^{1*♣} Haoyu Chen^{2*} Panwang Pan^{3*} Chenxin Li⁶
Sixiang Chen² Kairun Wen¹ Yeying Jin⁴ Wenbo Li^{5†} Xinghao Ding^{1†}

¹ Key Laboratory of Multimedia Trusted Perception and Efficient Computing, Ministry of Education of China, Xiamen University, Xiamen, Fujian, China

² The Hong Kong University of Science and Technology (Guangzhou)

³ Bytedance’s Pico ⁴ Tencent ⁵ Huawei Noah’s Ark Lab

⁶ The Chinese University of Hong Kong

Project page: <https://cvpr2025-jarvisir.github.io/>

Abstract

*Vision-centric perception systems struggle with unpredictable and coupled weather degradations in the wild. Current solutions are often limited, as they either depend on specific degradation priors or suffer from significant domain gaps. To enable robust and autonomous operation in real-world conditions, we propose JarvisIR, a VLM-powered agent that leverages the VLM as a controller to manage multiple expert restoration models. To further enhance system robustness, reduce hallucinations, and improve generalizability in real-world adverse weather, JarvisIR employs a novel two-stage framework consisting of supervised fine-tuning and human feedback alignment. Specifically, to address the lack of paired data in real-world scenarios, the human feedback alignment enables the VLM to be fine-tuned effectively on large-scale real-world data in an unsupervised manner. To support the training and evaluation of JarvisIR, we introduce CleanBench, a comprehensive dataset consisting of high-quality and large-scale instruction-responses pairs, including **150K** synthetic entries and **80K** real entries. Extensive experiments demonstrate that JarvisIR exhibits superior decision-making and restoration capabilities. Compared with existing methods, it achieves a **50%** improvement in the average of all perception metrics on CleanBench-Real. Project page: <https://cvpr2025-jarvisir.github.io/>.*

1. Introduction

Vision-centric perception systems often struggle in adverse weather, where images captured in real-world sce-

narios exhibit multiple and coupled degradations. Current adverse weather image restoration methods are primarily categorized into task-specific methods and all-in-one approaches. Both categories struggle with real-world coupled degradations, as shown in Figure 1. Task-specific methods [24, 31, 32, 45, 89] often require prior knowledge of specific degradation types, while real-world degradations are often unknown and coupled. All-in-one methods [13, 18, 30, 37, 49] trained on synthetic datasets in a supervised manner, suffer from a significant domain gap when applied to real-world data. One promising strategy to tackle multiple degradations in the wild is to integrate specialized models that excel in their domains. However, this strategy is highly sensitive to task order, and even minor changes in execution sequence can lead to significant performance degradation. Therefore, autonomously and efficiently coordinating expert models in real-world scenarios is essential for perceptual restoration.

Recently, large language models (LLMs) have exhibited remarkable proficiency in reasoning, decision-making and interaction with environments [26, 29, 53, 85, 98]. These advancements raise an important question: *Could vision-language models (VLMs) act as controllers, managing publicly available specialized restoration models, autonomously planning tasks, and selecting models to facilitate the development of comprehensive restoration systems?* The answer is affirmative, however, constructing such systems is non-trivial and typically requires extensive paired data. In real-world scenarios, while there exists extensive real degraded data, the lack of corresponding labels prevents the implementation of supervised fine-tuning approaches. To tackle this issue and harness large-scale unlabeled data, we design a fine-tuning framework based on human feedback, allowing the VLM to be trained in an unsupervised manner. With this approach, we could create a

* Authors Yunlong Lin and Zixu Lin contributed the most and led the study, while Authors Haoyu Chen and Panwang Pan also made significant contributions. † Corresponding author.

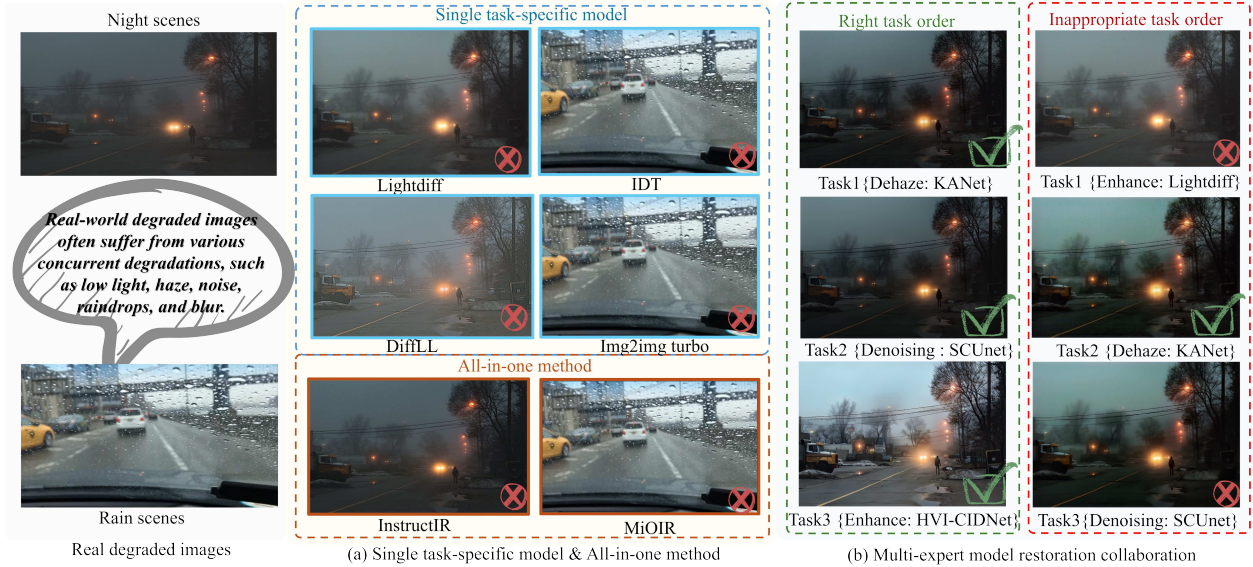


Figure 1. **Limitations of single-task methods, all-in-one methods, and inaccurate task order.** (a) Single-task specific and all-in-one methods fail to address coupled degradation in real-world scenarios. (b) Collaboration among multi-expert models effectively mitigates complex degradation, but is sensitive to the order of tasks. Unlike these approaches, JarvisIR can dynamically schedule different expert models in response to the rapidly changing scenarios and coupled degradation in the wild.

system that performs robustly and reliably in the wild.

In this work, we introduce JarvisIR, a VLM-powered agent integrating VLM (i.e., Llava-Llama3 [46]) with expert restoration models sourced from GitHub and Hugging Face. The development of this system involved two key components: 1) CleanBench, an instruction-following dataset constructed using the self-instruct strategy [73], which includes 150K synthetic and 80K real instruction-response pairs (CleanBench-Real), designed to support both training and evaluation. 2) A supervised fine-tuning (SFT) and human feedback alignment framework for training a VLM as an agent to be reliable and autonomous. Specifically, to enable the VLM to follow user instructions and perceive image degradation, we train it using the synthetic portion of CleanBench via SFT [50]. To enhance system robustness, reduce hallucinations, and improve generalizability in real-world adverse weather, we fine-tune JarvisIR on CleanBench-Real with human feedback. To ensure stability during training and improve overall performance, we propose the MRRHF algorithm, an extension of the ranking responses with human feedback (RRHF) approach [93]. Specifically, to expand the exploration space while maintaining a performance lower bound for JarvisIR, we introduce a hybrid sample generation strategy and regularization term. Furthermore, to comprehensively feedback the quality of system responses during training, we incorporate multiple VLM-based Image Quality Assessment (IQA) models as a unified reward model.

Our contributions can be summarized as follows:

- We introduce JarvisIR, a VLM-powered agent that autonomously manages and coordinates multiple expert restoration models to address coupled weather degradations in real-world environments.
- We present CleanBench, the first high-quality instruction-following dataset specifically curated for developing intelligent restoration systems, containing 150K synthetic and 80K real instruction-response pairs.
- We propose a novel two-stage framework combining supervised fine-tuning and human feedback alignment to enhance system robustness, reduce hallucinations and improve generalizability in the wild.
- Our experiments demonstrate that JarvisIR outperforms strong baselines in terms of decision-making and perception restoration.

2. Related Work

Tool-Augmented LLMs. Recent studies [6, 53, 56, 60–62, 97] highlight the growing potential of large language models (LLMs) for proficient tool usage and decision-making in complex settings. For example, Gorilla [53] facilitates LLMs’ response to Tool calls through dataset construction and fine-tuning. ToolLLM [56] extends this concept to enable interaction with a large number of tools. ToolAlpaca [65] demonstrates the feasibility of generalized tool-use capabilities in smaller LLMs. Toolformer [60] constructs tool-use augmented data to train LLMs to select tools. In the realm of visual tools, various approaches have been proposed to enhance the capabilities of large language

models in handling visual tasks [78, 88], augmented with Hugging Face models [61], Azure models [88], visual foundation models [78].

Alignment of LLMs. Reinforcement Learning from Human Feedback (RLHF) [1, 2, 28, 66] has emerged as a groundbreaking technique for aligning LLMs. The core idea is learning a reward function to reflect human preferences with human annotations and optimize LLMs by RL methods like proximal policy optimization (PPO). During PPO-based optimization, updating LLMs requires the likelihood of an entire generation. However, for LLM agents, human feedback is usually obtained only after the tool response is completed and the function is successfully invoked. Moreover, unlike typical LLM training, our two-stage fine-tuning process integrates both visual and linguistic modalities. Rank Responses to Align Human Feedback (RRHF) [93] has shown promise by using reward models to rank multiple responses, aligning LLMs effectively. This technique allows easy extension to fine-grained tool agents, thereby maximizing the utility of existing reward models.

Image Restoration. Single-task image restoration has achieved significant progress in addressing specific degradation types, such as dehazing [31, 42, 81], low-light enhancement [27, 33, 44], desnowing [8, 14, 17], deraining [12, 16, 82], denoising [7, 94] super-resolution [10, 64, 70, 72], image fusion [22, 23, 43, 74]. However, these task-specific approaches often lack generalizability and adaptability to complex, coupled degradations. To overcome this limitation, adverse weather restoration research aims to develop a unified framework capable of addressing multiple degradation types simultaneously [15, 25, 41, 51]. Another prevailing research line is dedicated to building more intelligent restoration systems. Clarity ChatGPT [77], integrated with advanced visual models, allows users to perform sophisticated image manipulation and enhancement through natural language interactions. RestoreAgent [9] and AgenticIR [101] are contemporaneous independent works that utilize MLLM as a task planner to coordinate multiple restoration tools. Specifically, RestoreAgent [9] involves fine-tuning a vision-language model (VLM) using synthetic datasets to directly produce an execution plan. AgenticIR [101] leverages two off-the-shelf LLMs and VLMs to achieve the scheduling of restoration tools on the synthetic experiment platform. Essentially, both studies focus on building intelligent restoration systems tailored for synthetic degradation conditions. *Conversely, our study aims to develop a robust system for real-world applications, incorporating human feedback to enhance robustness, reduce hallucinations and improve generalizability. Furthermore, our system is built in an unsupervised manner using large-scale, unlabeled real-world data.*

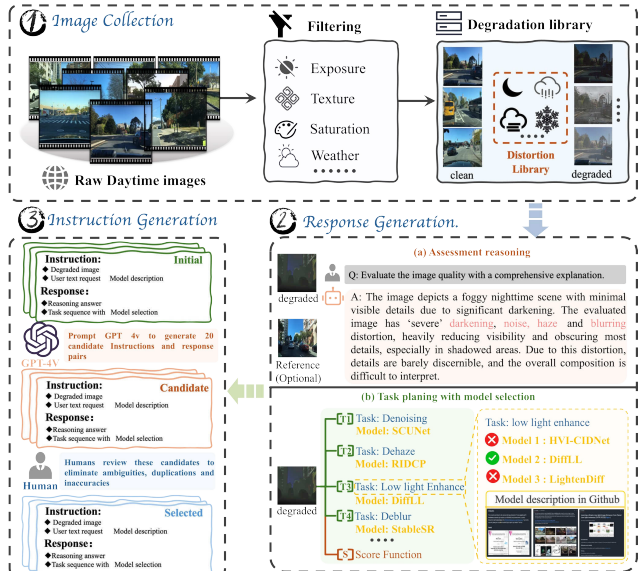


Figure 2. The dataset construction workflow consists of three main steps: 1) Synthesis of degraded images. 2) Generation of Assessment reasoning and the optimal task sequence. 3) Generation of instruction-response pairs for the system.

3. Methodology

In this section, we first describe CleanBench, a comprehensive benchmark consisting of extensive instruction-response pairs used for the training and evaluation of JarvisIR (Sec. 3.1). We then introduce JarvisIR, a VLM agent to call expert restoration models in response to intricate multiple degraded environments in the wild (Sec. 3.2). Finally, we describe the two-stage training framework for JarvisIR, comprising supervised fine-tuning and human feedback alignment.

3.1. CleanBench

High-quality and large-scale datasets are crucial for unleashing the full potential of VLMs. A multimodal instruction sample can be formally represented as a triplet: $\{user\ instruction, degraded\ image, response\}$, where “*user instruction*” specifies the task and describes the restoration tools, “*degraded image*” serves as the visual input to be processed, and the “*response*” provides the ground truth answer. In Figure 2, we outline the construction of our dataset, focusing on the generation of degraded images and the collection of task-specific instructions and responses.

Image Collection. We first collect raw daytime images from various sources, including autonomous driving datasets [4, 59, 92, 102] and natural scenes [5, 34, 39, 42, 47, 87, 99]. Then, Q-instruct [80] serves as a quality filter to extract high-quality samples. To simulate realistic adverse weather scenarios, including rainy, nighttime, snowy, and foggy, we customized the degradation library developed using physical models and image transformation techniques to



Figure 3. Examples of CleanBench-Real dataset.

synthesize degraded images. More detail in supplementary material.

Response Generation. The response from JarvisIR consists of two components: “chain-of-thought” (COT) rationales and the optimal task sequence with model selection. (a) For COT rationales, we distill DepictQA-Wild’s [91] knowledge, which excels in low-level quality reasoning assessment. Specifically, given a degraded image pair, we prompt DepictQA-Wild [91] to assess the quality of the degraded image in terms of clarity, colorfulness, and sharpness, generating detailed degradation and reasoning insights. (b) To determine the optimal task sequence with restoration model selection, we employ an exhaustive search strategy [9] to explore various task permutations and model combinations, scoring each sequence to identify the optimal restoration path.

Task-model Assignment. User instructions include descriptions of available tasks and models, sourced from GitHub or Hugging Face, to formulate task-model assignment as a single-choice problem. Presenting tasks and models as options within a context allows JarvisIR to more effectively identify the appropriate model for each sub-task.

Instruction Generation. Motivated by the self-instruct strategy [73], for each initial user instruction and response, GPT-4V is prompted to generate 20 candidate pairs. We then manually review these candidates to eliminate ambiguity, repetition, and inaccuracies, ultimately selecting 5 instruction-response pairs per degraded image (see supplementary material for details). Ultimately, CleanBench includes a total of 150K instruction-response pairs, which are used in the initial instruction-tuning phase.

CleanBench-Real. To align and evaluate JarvisIR’s performance in real-world scenarios, we introduce CleanBench-Real, comprising 80K unlabeled real degraded images from internet and diverse sources [4, 33, 34, 47, 57, 71, 87, 92]. CleanBench-Real is categorized into four adverse weather scenarios: rainy, night, snowy, and foggy. The degradation in each scenario is complex and intertwined. For example, as presented in Figure 3, an image captured in rain may experience multiple degradations concurrently, including rain, raindrops, defocus blur, and noise (more in sup-

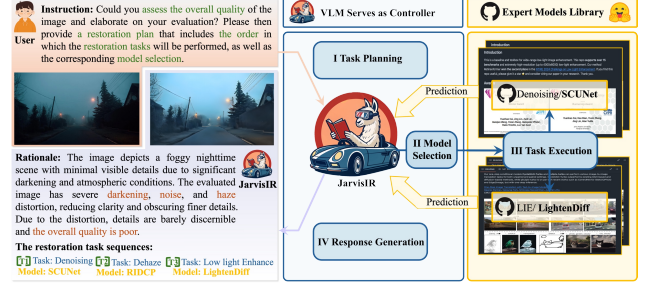


Figure 4. The workflow of JarvisIR. To address real-world coupled weather degradation, we develop JarvisIR, a VLM-powered intelligent system that dynamically schedules expert models for restoration. Initially, JarvisIR assesses the degradation of the input images and parses user instructions to formulate a task plan, selecting the appropriate expert models for each subtask. The selected experts perform their designated tasks and return the results to JarvisIR, which integrates the outcomes and provides the final answer to the user. The design of the figure is inspired by [61].

plementary material). For the division of the training and evaluation sets, we selected 500 images from each of the four CleanBench-Real scenarios to form the evaluation set (2K), while the remaining images are utilized for alignment tuning. Instruction-response pairs are generated in the same way as outlined in CleanBench.

3.2. JarvisIR

JarvisIR is a VLM-powered agent that coordinates multiple expert restoration models to address complex degradation. As illustrated in Figure 4, the workflow of JarvisIR consists of four steps: Task Planning, Model Selection, Task Execution, and Response Generation. To enhance the agent’s decision-making and perception restoration capabilities in real-world scenarios, as depicted in Figure 5, we initially perform supervised fine-tuning (SFT) on CleanBench to obtain an initial version, termed JarvisIR-SFT. Subsequently, the JarvisIR-SFT is further fine-tuned utilizing the MR-RHF algorithm on CleanBench-Real, yielding the JarvisIR-MRRHF model.

3.2.1. JarvisIR-SFT

We employ the standard SFT to get the JarvisIR-SFT model. Formally, the multimodal instruction sample can be denoted in a triplet form $(\mathcal{I}, \mathcal{M}, \mathcal{R})$, where \mathcal{I} , \mathcal{M} , \mathcal{R} represent the user instruction, the degraded image, and the ground truth response, respectively. The VLM predicts an answer \mathcal{A} given the instruction and the degraded image: $\mathcal{A} = f(\mathcal{I}, \mathcal{M}; \theta)$. The training objective is the original auto-regressive objective used to train LLMs [48, 90]:

$$L_{sft} = - \sum_{i=1}^N \log P_{\pi}(\mathcal{R}_i | \{\mathcal{I}_i, \mathcal{M}_i\}, \mathcal{R}_{<i}); \theta), \quad (1)$$

where N is the length of the ground-truth response.

3.2.2. JarvisIR-MRRHF

Intuitively, SFT allows JarvisIR-SFT to achieve favorable performance on synthetic data. Nevertheless, as previously noted, due to the distribution shift, transferring from synthetic training data to real test data, JarvisIR-SFT exhibits increased hallucination, i.e., degraded perception restoration performance and decision-making capability. To improve its generalizability, we further fine-tune JarvisIR on CleanBench-Real with refined ranking responses with human feedback algorithm (MRRHF).

Reward modeling. The reward model evaluates tool-calling outcomes and converts them into structured reward signals to guide the agent’s optimization process. Therefore, selecting an appropriate reward model is crucial. Fortunately, in the image quality assessment (IQA) field, VLM-based IQA models have been developed [80], demonstrating strong performance in evaluating aesthetic quality and image distortion. These IQA models are inherently suitable for serving as reward models. To construct a comprehensive reward model \mathcal{S} , as well as an evaluation system, we integrated multiple IQA models. Specifically, we employ a z-score strategy [9] to standardize the scores assessed by each IQA model separately and then sum the standardized results:

$$\mathcal{S} = \sum_{i=1}^k \frac{s_i - \mu_i}{\sigma_i}, \quad (2)$$

where s_i represents the score assessed by i -th IQA model. μ_i and σ_i represent the mean and standard deviation of s_i , respectively. k indicates the total number of IQA models.

Alignment with MRRHF. We propose an extension to the existing RRHF method that can be used for aligning JarvisIR in a cost-effective manner: 1) A hybrid sample generation strategy that combines offline and online approaches to expand the optimization exploration space while ensuring a performance lower bound. 2) Entropy regularization terms are integrated to foster diversity among agent responses, thereby facilitating exploration during training. Specifically, for a pair of user instruction \mathcal{I}_i and degraded image \mathcal{M}_i , we first adopt offline diverse beam search [68] to get m_1 different responses $\mathcal{R}_{m_1} = \{r_1, r_2, \dots, r_{m_1}\}$ from SFT model π . Similarly, we can obtain $\mathcal{R}_{m_2} = \{r_1, r_2, \dots, r_{m_2}\}$ from policy model ρ (initialized from SFT model π) during training. The combined candidate m responses are denoted as $\mathcal{R}_m = \mathcal{R}_{m_1} \cup \mathcal{R}_{m_2}$. Subsequently, we execute the task sequences specified in candidate responses, calling multiple restoration models to generate restored images. These predictions are then assessed by the reward model \mathcal{S} , yielding scores for each r_i with $\mathcal{S}(r_i) = s_i$. To align with scores $\{s_i\}_m$, we use policy model ρ to give scores p_i for each r_i by:

$$p_i = \frac{\sum_t \log P_\rho(r_{i,t} | \{\mathcal{I}_i, \mathcal{M}_i\}, r_{i,<t}; \theta)}{\|r_i\|}, \quad (3)$$

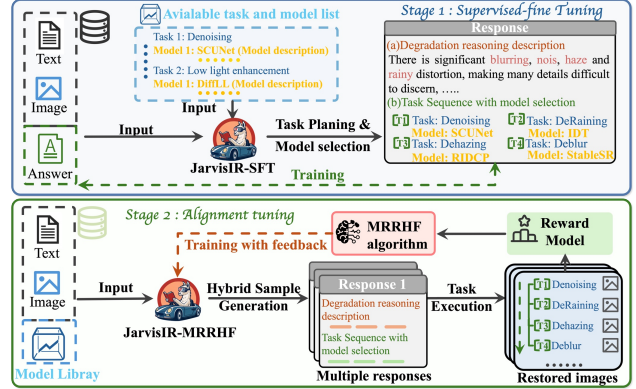


Figure 5. Two-stage training framework of JarvisIR. In the first stage, JarvisIR undergoes supervised fine-tuning on synthetic data from CleanBench to enable it to follow user instructions and recognize image degradation. In the second stage, we further fine-tune JarvisIR on CleanBench-Real using the MRRHF algorithm to improve system robustness, reduce hallucinations, and enhance generalizability under real-world adverse weather conditions.

where p_i is conditional log probability (length-normalized) of r_i under model ρ . The core idea is letting the policy model ρ give larger probabilities for better responses and give smaller probabilities for worse responses. Inspired by PRO [63], we refine the original ranking loss:

$$L_{\text{rank}} = \sum_{s_i < s_j} (s_j - s_i) \max(0, p_i - p_j), \quad (4)$$

and a cross-entropy loss like SFT process is added to learn the response with the highest reward s_i , $i' = \arg \max_i s_i$:

$$L_{ft} = - \sum_t \log P_\rho(r_{i',t} | \{\mathcal{I}_i, \mathcal{M}_i\}, r_{i',<t}; \theta). \quad (5)$$

Furthermore, we define the entropy regularization term as:

$$L_{er} = - \sum_a \rho(a | y) \log \rho(a | y), \quad (6)$$

where y represents the current state of the agent. The overall loss is utilized to optimize the JarvisIR-SFT to derive JarvisIR-MRRHF:

$$L = \lambda_1 L_{\text{rank}} + \lambda_2 L_{ft} + \lambda_3 L_{er}, \quad (7)$$

where λ_1 , λ_2 and λ_3 are constants controlling the relative importance of the different losses, which are empirically set to 0.5, 0.5 and 0.1 in all experiments, respectively.

Discussion of RLHF and RRHF: The training of vanilla RLHF [50] necessitated a policy model, a value model, a reward model, and a reference model, which could be demanding on memory resources. Rank Responses to Align Human Feedback (RRHF) [93] can effectively alleviate



Figure 6. Visual comparisons of various methods on CleanBench-Real. Our approach delivers significant quality improvements, eliminating complex real-world degradation and preserving the most natural details.

the issues of resource-intensive and tedious hyperparameter tuning in RLHF. However, directly fine-tuning JarvisIR using RRHF yields limited improvement to its generalization in real-world scenarios. Although vanilla RRHF employs an off-policy learning strategy that could save time by avoiding the need to generate new responses during training, it has the drawback of relying on a static offline preference dataset for training the policy model. Consequently, the policy might over-optimize for reward on in-distribution data as the model cannot further query the preference oracle during the training process [75]. The RRHF incorporating online sampling like PPO might mitigate this issue, but it demands more GPU resources to store the reference model, thereby significantly decreasing the training speed [93].

4. Experiments

4.1. Experimental Settings

Training Setup. Llava-Llama3-8b [46] serves as the base model for JarvisIR, which undergoes full parameter fine-tuning using the Adam optimizer. During the SFT phase, we fine-tune JarvisIR for 3 epochs with a batch size of 128 and a learning rate of $1e-5$. In the MRRHF tuning phase, we set the diverse beam search size to 3, the diverse beam group to 5, the diversity penalty to 2.0, and the sampling temperature to 0.8. Alignment tuning is performed over 3 epochs with a batch size of 1 and a learning rate of $1e-5$. To speed up training, we select three IQA models—Q-instruct [80], MUSIQ [35] and MANIQA [86]—to construct the unified reward model (Eq. 2). All experiments are conducted on 8 NVIDIA A100 80G GPUs.

Dataset Settings & Metrics. The CleanBench is fully utilized for supervised fine-tuning of Llava-Llama3-8b [46] to obtain JarvisIR-SFT. The training set of CleanBench-Real is used for alignment tuning, yielding JarvisIR-MRRHF. Additionally, JarvisIR’s evaluation is conducted on the validation set of CleanBench-Real, focusing on 1) decision-making ability and 2) perception restoration capability in

Table 1. Comparison of JarvisIR with other strategies on the CleanBench-Real validation set. The “Score” represents the sum of the four normalized metrics. The “Ranking” indicates the given decision’s percentage ranking among all possible decisions. We highlight the **best** and **second-best** results.

Strategy	Score	Ranking(%)
(I) Random Order and Model	1.12	43.2%
(II) Random Order + Predict Model	2.66	34.7%
(III) Random Model + Predict Order	3.08	23.4%
(IV) Pre-defined Order and Model	3.94	22.5%
(V) Human Expert	4.85	18.6%
★JarvisIR-SFT	5.17	14.3%
★JarvisIR-MRRHF	6.21	4.8%

real-world scenarios. Due to the lack of paired clean-degraded data in the real scenarios. Four image quality assessment metrics are used for evaluation: MUSIQ [35], MANIQA [86], CLIP-IQA+ [69], LIQE [95].

Tool Settings. We present the task-specific restoration tools employed in our implementation, including denoising (SCUnet [94]), super-resolution & deblur & compression artifact removal (StableSR-turbo [70] and Real-ESRGAN [72]), deraining (IDT [82], UDR-S2Former [8] and Img2img-turbo [52]), dehazing (RIDCP [81] and KANet [21]), low-light enhancement (Img2img-turbo [52], HVI-CIDNet [83] and LightenDiff [27]) and desnowing (Img2img-turbo [52] and Snowformer [11]). More details are in the supplementary material. Notably, we select lightweight and efficient models instead of the latest state-of-the-art models to simplify the validation process of our proposed paradigm. Incorporating more advanced models could further enhance performance.

4.2. Decision Making Capability

Compared Baselines. We conducted a comparative analysis of JarvisIR against several alternative approaches: (I) Random selection of both the task order and the models, as-

Table 2. Comparison of JarvisIR with All-in-One methods for multi-degraded perception restoration on CleanBench-Real. We highlight the best, second-best and third-best results. Notably, all scenes represent multiple degraded weather conditions, such as haze, low light and blur.

Method	Night Scenes				Rain Scenes			
	MUSIQ \uparrow	MANIQA \uparrow	CLIP-IQA+ \uparrow	LIQE \uparrow	MUSIQ \uparrow	MANIQA \uparrow	CLIP-IQA+ \uparrow	LIQE \uparrow
AirNet [40]	44.26	0.1889	0.4429	1.313	62.61	0.3871	0.5867	3.136
AutoDIR [30]	47.30	0.1885	0.4341	1.403	63.93	0.4002	0.6082	3.312
DA-CLIP [49]	45.86	0.2010	0.4544	1.427	63.28	0.3993	0.5959	3.194
PromptIR [55]	45.45	0.2010	0.4473	1.408	62.85	0.3926	0.5941	3.161
MiOIR [37]	46.93	0.2013	0.4403	1.408	63.07	0.3779	0.5841	3.055
InstructIR [18]	44.03	0.1533	0.3689	1.257	62.93	0.3657	0.5609	3.055
T ³ -DiffWeather [13]	46.79	0.1964	0.4547	1.413	62.67	0.3689	0.5823	3.011
★JarvisIR-SFT	60.77	0.5048	0.5239	3.224	65.03	0.5339	0.6290	4.005
★JarvisIR-MRRHF	67.25	0.5876	0.6336	3.613	70.38	0.7004	0.7127	4.435

Method	Fog Scenes				Snow Scenes			
	MUSIQ \uparrow	MANIQA \uparrow	CLIP-IQA+ \uparrow	LIQE \uparrow	MUSIQ \uparrow	MANIQA \uparrow	CLIP-IQA+ \uparrow	LIQE \uparrow
AirNet [40]	64.23	0.3829	0.6173	2.686	67.32	0.4320	0.6379	3.794
AutoDIR [30]	64.84	0.3966	0.6443	2.928	67.62	0.4305	0.6453	3.824
DA-CLIP [49]	64.78	0.3880	0.6540	2.793	67.71	0.4294	0.6426	3.817
PromptIR [55]	64.54	0.3810	0.6417	2.557	67.34	0.4292	0.6435	3.776
MiOIR [37]	64.93	0.3501	0.5969	2.415	67.28	0.4187	0.6404	3.702
InstructIR [18]	64.82	0.3904	0.6449	2.919	67.98	0.4038	0.6052	3.715
T ³ -DiffWeather [13]	64.58	0.3715	0.6163	2.497	67.72	0.4129	0.6268	3.713
★JarvisIR-SFT	70.45	0.4855	0.6560	3.977	70.24	0.7133	0.7127	4.086
★JarvisIR-MRRHF	74.22	0.7502	0.7805	4.714	73.87	0.8014	0.7918	4.881

suming that task types are accurately determined. (II) Random task order, but models predicted by JarvisIR. (III) Random model selection, but task orders predicted by JarvisIR. (IV) Using a human expert’s predefined order and models for different scenes, assuming the approximate scene degradation can be determined. (V) A human expert manually generates a solution case by case for each image, determining both the task sequence and the appropriate models.

Results. As indicated in Table 1, strategies that involve human expert participation—specifically settings (IV) and (V)—demonstrate strong performance compared to random strategies, ranking within the top 22.5% and 18.6% of all possible strategies, respectively. These results indicate the effectiveness of human experts’ experience in complex decision-making processes. Interestingly, however, our JarvisIR model achieves the highest performance, surpassing even the expert-driven customization strategies. Furthermore, JarvisIR-MRRHF (4.8%) outperforms JarvisIR-SFT (14.3%) in both score and ranking, highlighting that the MRRHF stage in our training framework effectively mitigates hallucination errors in system responses, thereby enabling the generation of more optimal decisions.

4.3. Perception Restoration Ability

Compared All-in-One Methods. We compare JarvisIR with existing advanced all-in-one methods: AirNet [40],

AutoDIR [30], DA-CLIP [49], PromptIR [55], MiOIR [37], InstructIR [18], T³-DiffWeather [13]. For a fair comparison, we repeatedly run these compared methods multiple times to fully leverage their capabilities. Additionally, we supply InstructIR and AutoDIR with explicit prompts detailing degradation scenarios to optimize their performance. **Results.** As shown in Table 2 and Figure 6, JarvisIR outperforms existing All-in-One approaches across all metrics. In Night Scenes, JarvisIR-MRRHF achieves a MUSIQ score of 67.25, which is 42.2% higher than AutoDIR’s score of 47.30. In MANIQA, JarvisIR-MRRHF scores 0.5876, much better than DA-CLIP (0.2010) and MiOIR (0.2013). These results show that JarvisIR autonomously selects optimal task sequences and models, outperforming methods with predefined or random sequences. Additionally, JarvisIR-MRRHF also exceeds the SFT version in all scenes, with notable gains in Rain (70.38 vs. 65.03 MUSIQ) and Fog (74.22 vs. 70.45 MUSIQ). These results demonstrate that JarvisIR fine-tuned with MRRHF can improve generalizability, fewer hallucination errors, and better decision-making ability.

5. Ablation Study

Sample generation strategy. To assess the effectiveness of the hybrid sample generation strategy, we compared it with two variations of the original setting: 1) offline sample gen-

Table 3. Ablation studies on different sample generation strategies and entropy regularization. The ‘‘Reward’’ represents the average reward scores obtained during MRRHF training, spanning from -1 to 1. A negative score indicates a penalty, while a positive score represents a reward. The ‘‘Diversity’’ reflects the average number of unique responses produced during the training process.

Strategy	Reward	Diversity
offline sample generation	0.43	3.63
online sample generation	-0.87	1.27
hybrid sample generation (ours)	0.67	6.55
w/o. entropy regularization	0.50	4.56
w. entropy regularization (ours)	0.67	6.55

eration strategy. 2) online sample generation strategy. The results in Table 3 and Figure 7 yield the following observations: 1) The offline sample generation strategy yields limited performance gains, with a reward score of 0.43 and a diversity score of 3.63. This limitation arises because the sample distribution is restricted to the finite dataset generated by the SFT model using diverse beam search [68]. Consequently, the policy model may over-optimize for in-distribution data, thereby limiting its ability to generalize and achieve higher reward scores. 2) The online sampling strategy initially yields higher reward scores and diversity. However, as training progresses, the model encounters a collapse, leading to a significantly low reward score (-0.87) and decreased diversity (1.27). This instability may result from an excessively large optimization space without adequate constraints during training. When the model reaches a local minimum, it struggles to escape, as the candidate responses generated using diverse beam search [68] are of poor quality, causing the model to produce repetitive and invalid responses. Our hybrid sampling approach combines both online and offline samples, resulting in superior performance with a reward score of 0.67 and the highest diversity score of 6.55. This balanced strategy leverages the advantages of both online and offline sampling, ensuring stable training by providing sufficient exploration space while avoiding the pitfalls associated with purely online sampling. As a result, the hybrid strategy maintains high reward scores and diversity throughout training, outperforming both online and offline strategies.

Entropy regularization. As discussed in Sec. 3.2.2, entropy regularization significantly affects the diversity of system responses during training. The results in Table 3 and Figure 7 show that without this regularization, the reward decreases from 0.67 to 0.50, while the diversity drops from 6.55 to 4.56. This highlights the role of entropy regularization in fostering greater exploration and producing more diverse, high-quality responses.

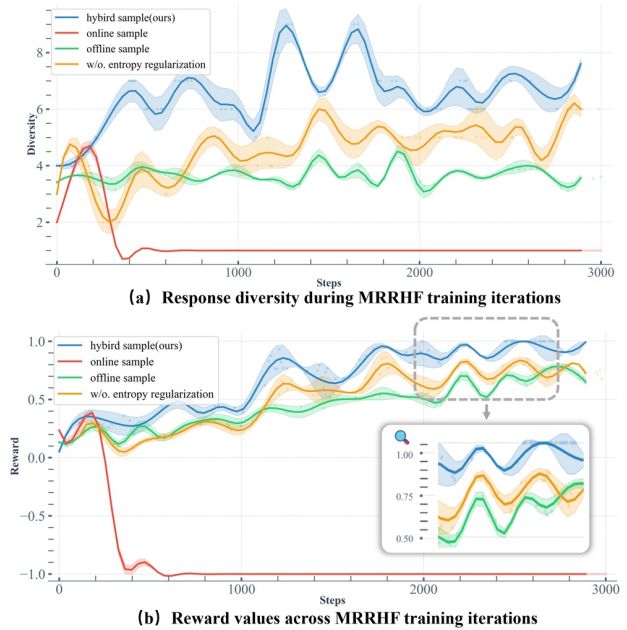


Figure 7. Ablation studies on different sample generation strategies and entropy regularization. (a) Response diversity during MRRHF training iterations. (b) Reward values across MRRHF training iterations.

6. Conclusions

This paper introduces JarvisIR, a VLM-powered intelligent system that leverages Llava-Llama3 to connect distinct restoration expert models. JarvisIR can autonomously schedule different expert models in response to the rapidly changing scenarios and coupled degradation in autonomous driving and natural environments. To enhance system robustness, minimize hallucinations, and improve generalizability, we propose a novel two-stage framework comprising supervised fine-tuning and human feedback alignment. Specifically, we design the human feedback alignment to effectively tune the VLM in an unsupervised manner, leveraging large-scale unlabeled real-world data. To support the training and evaluation of JarvisIR, we present CleanBench, a high-quality, large-scale dataset containing 150K synthetic and 80K real instruction-response pairs. Experiments show that JarvisIR outperforms existing methods, achieving a 50% improvement in the average of all perception metrics on CleanBench-Real.

7. Acknowledgments

This work was supported in part by the National Natural Science Foundation of China under Grant 82172033, Grant U19B2031, Grant 61971369, Grant 52105126, Grant 82272071, and Grant 62271430; and in part by the Dreams Foundation of Jianghuai Advance Technology Center; and in part by the Open Fund of the National Key Laboratory of Infrared Detection Technologies.

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 3
- [2] AI Anthropic. The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card*, 1, 2024. 3
- [3] Tim Brooks, Ben Mildenhall, Tianfan Xue, Jiawen Chen, Dillon Sharlet, and Jonathan T Barron. Unprocessing images for learned raw denoising. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11036–11045, 2019. 16
- [4] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020. 3, 4, 14
- [5] Jianrui Cai, Shuhang Gu, and Lei Zhang. Learning a deep single image contrast enhancer from multi-exposure images. *IEEE Transactions on Image Processing*, 27(4):2049–2062, 2018. 3, 14, 15
- [6] Wenhao Chai, Enxin Song, Yilun Du, Chenlin Meng, Vashisht Madhavan, Omer Bar-Tal, Jeng-Neng Hwang, Saining Xie, and Christopher D Manning. Auroracap: Efficient, performant video detailed captioning and a new benchmark. *arXiv preprint arXiv:2410.03051*, 2024. 2
- [7] Haoyu Chen, Jinjin Gu, Yihao Liu, Salma Abdel Magid, Chao Dong, Qiong Wang, Hanspeter Pfister, and Lei Zhu. Masked image training for generalizable deep image denoising. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1692–1703, 2023. 3
- [8] Haoyu Chen, Jingjing Ren, Jinjin Gu, Hongtao Wu, Xuequan Lu, Haoming Cai, and Lei Zhu. Snow removal in video: A new dataset and a novel method. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 13211–13222, 2023. 3, 6, 15
- [9] Haoyu Chen, Wenbo Li, Jinjin Gu, Jingjing Ren, Sixiang Chen, Tian Ye, Renjing Pei, Kaiwen Zhou, Fenglong Song, and Lei Zhu. Restoreagent: Autonomous image restoration agent via multimodal large language models. *arXiv preprint arXiv:2407.18035*, 2024. 3, 4, 5
- [10] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 25857–25867, 2024. 3
- [11] Sixiang Chen, Tian Ye, Yun Liu, and Erkang Chen. Snowformer: Context interaction transformer with scale-awareness for single image desnowing. *arXiv preprint arXiv:2208.09703*, 2022. 6, 15
- [12] Sixiang Chen, Tian Ye, Jinbin Bai, Erkang Chen, Jun Shi, and Lei Zhu. Sparse sampling transformer with uncertainty-driven ranking for unified removal of raindrops and rain streaks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 13106–13117, 2023. 3
- [13] Sixiang Chen, Tian Ye, Kai Zhang, Zhaohu Xing, Yunlong Lin, and Lei Zhu. Teaching tailored to talent: Adverse weather restoration via prompt pool and depth-anything constraint. In *European Conference on Computer Vision*, pages 95–115. Springer, 2025. 1, 7
- [14] Wei-Ting Chen, Hao-Yu Fang, Cheng-Lin Hsieh, Cheng-Che Tsai, I Chen, Jian-Jiun Ding, Sy-Yen Kuo, et al. All snow removed: Single image desnowing algorithm using hierarchical dual-tree complex wavelet representation and contradict channel loss. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4196–4205, 2021. 3
- [15] Wei-Ting Chen, Zhi-Kai Huang, Cheng-Che Tsai, Hao-Hsiang Yang, Jian-Jiun Ding, and Sy-Yen Kuo. Learning multiple adverse weather removal via two-stage knowledge learning and multi-contrastive regularization: Toward a unified model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17653–17662, 2022. 3
- [16] Xiang Chen, Hao Li, Mingqiang Li, and Jinshan Pan. Learning a sparse transformer network for effective image deraining. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5896–5905, 2023. 3
- [17] Bodong Cheng, Juncheng Li, Ying Chen, and Tiejiong Zeng. Snow mask guided adaptive residual network for image snow removal. *Computer Vision and Image Understanding*, 236:103819, 2023. 3
- [18] Marcos V Conde, Gregor Geigle, and Radu Timofte. High-quality image restoration following human instructions. *arXiv preprint arXiv:2401.16468*, 2024. 1, 7
- [19] Ziteng Cui, Guo-Jun Qi, Lin Gu, Shaodi You, Zenghui Zhang, and Tatsuya Harada. Multitask aet with orthogonal tangent regularity for dark object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2553–2562, 2021. 16
- [20] Yuekun Dai, Chongyi Li, Shangchen Zhou, Ruicheng Feng, Yihang Luo, and Chen Change Loy. Flare7k++: Mixing synthetic and real datasets for nighttime flare removal and beyond. *arXiv preprint arXiv:2306.04236*, 2023. 16
- [21] Yuxin Feng, Long Ma, Xiaozhe Meng, Fan Zhou, Risheng Liu, and Zhuo Su. Advancing real-world image dehazing: perspective, modules, and training. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 6, 15
- [22] Chunming He, Kai Li, Guoxia Xu, Jiangpeng Yan, Longxiang Tang, Yulun Zhang, Yaowei Wang, and Xiu Li. Hqg-net: Unpaired medical image enhancement with high-quality guidance. *TNNLS*, 2023. 3
- [23] Chunming He, Kai Li, Guoxia Xu, Yulun Zhang, Runze Hu, Zhenhua Guo, and Xiu Li. Degradation-resistant unfolding network for heterogeneous image fusion. In *ICCV*, pages 12611–12621, 2023. 3
- [24] Chunming He, Chengyu Fang, Yulun Zhang, Kai Li, Longxiang Tang, Chenyu You, Fengyang Xiao, Zhenhua Guo, and Xiu Li. Reti-diff: Illumination degradation image

- restoration with retinex-based latent diffusion model. *ICLR*, 2025. 1
- [25] Chunming He, Yuqi Shen, Chengyu Fang, Fengyang Xiao, Longxiang Tang, Yulun Zhang, Wangmeng Zuo, Zhenhua Guo, and Xiu Li. Diffusion models in low-level vision: A survey. *TPAMI*, 2025. 3
- [26] Shashank Mohan Jain. Hugging face. In *Introduction to transformers for NLP: With the hugging face library and models to solve problems*, pages 51–67. Springer, 2022. 1
- [27] Hai Jiang, Ao Luo, Xiaohong Liu, Songchen Han, and Shuaicheng Liu. Lightdiffusion: Unsupervised low-light image enhancement with latent-retinex diffusion models. *arXiv preprint arXiv:2407.08939*, 2024. 3, 6, 15
- [28] Songtao Jiang, Yan Zhang, Ruizhe Chen, Yeying Jin, and Zuozhu Liu. Modality-fair preference optimization for trustworthy mllm alignment. *arXiv preprint arXiv:2410.15334*, 2024. 3
- [29] Songtao Jiang, Tuo Zheng, Yan Zhang, Yeying Jin, Li Yuan, and Zuozhu Liu. Med-moe: Mixture of domain-specific experts for lightweight medical vision-language models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 3843–3860, 2024. 1
- [30] Yitong Jiang, Zhaoyang Zhang, Tianfan Xue, and Jinwei Gu. Autodir: Automatic all-in-one image restoration with latent diffusion. *arXiv preprint arXiv:2310.10123*, 2023. 1, 7
- [31] Yeying Jin, Wending Yan, Wenhan Yang, and Robby T Tan. Structure representation network and uncertainty feedback learning for dense non-uniform fog removal. In *Proceedings of the Asian Conference on Computer Vision*, pages 2041–2058, 2022. 1, 3
- [32] Yeying Jin, Wenhan Yang, and Robby T Tan. Unsupervised night image enhancement: When layer decomposition meets light-effects suppression. In *European Conference on Computer Vision*, pages 404–421. Springer, 2022. 1
- [33] Yeying Jin, Beibei Lin, Wending Yan, Yuan Yuan, Wei Ye, and Robby T Tan. Enhancing visibility in nighttime haze images using guided apsf and gradient adaptive convolution. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 2446–2457, 2023. 3, 4
- [34] Yeying Jin, Xin Li, Jiadong Wang, Yan Zhang, and Malu Zhang. Raindrop clarity: A dual-focused dataset for day and night raindrop removal. In *European Conference on Computer Vision*, pages 1–17. Springer, 2025. 3, 4, 14, 15
- [35] Junjie Ke, Qifei Wang, Yilin Wang, Peyman Milanfar, and Feng Yang. Musiq: Multi-scale image quality transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5148–5157, 2021. 6, 17, 18, 19
- [36] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022. 18
- [37] Xiangtao Kong, Chao Dong, and Lei Zhang. Towards effective multiple-in-one image restoration: A sequential and prompt learning strategy. *arXiv preprint arXiv:2401.03379*, 2024. 1, 7
- [38] T Kudo. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226*, 2018. 14
- [39] Boyi Li, Wenqi Ren, Dengpan Fu, Dacheng Tao, Dan Feng, Wenjun Zeng, and Zhangyang Wang. Benchmarking single-image dehazing and beyond. *IEEE Transactions on Image Processing*, 28(1):492–505, 2018. 3, 14, 15
- [40] Boyun Li, Xiao Liu, Peng Hu, Zhongqin Wu, Jiancheng Lv, and Xi Peng. All-in-one image restoration for unknown corruption. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 17452–17462, 2022. 7
- [41] Ruoteng Li, Robby T Tan, and Loong-Fah Cheong. All in one bad weather removal using architectural search. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3175–3185, 2020. 3
- [42] Beibei Lin, Yeying Jin, Wending Yan, Wei Ye, Yuan Yuan, and Robby T Tan. Nighthaze: Nighttime image dehazing via self-prior learning. *arXiv preprint arXiv:2403.07408*, 2024. 3, 14, 15
- [43] Yunlong Lin, Zhenqi Fu, Ge Meng, Yingying Wang, Yuhang Dong, Linyu Fan, Hedeng Yu, and Xinghao Ding. Domain-irrelevant feature learning for generalizable pansharpener. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 3287–3296, 2023. 3
- [44] Yunlong Lin, Zhenqi Fu, Kairun Wen, Tian Ye, Sixiang Chen, Ge Meng, Yingying Wang, Yue Huang, Xiaotong Tu, and Xinghao Ding. Unsupervised low-light image enhancement with lookup tables and diffusion priors. *arXiv preprint arXiv:2409.18899*, 2024. 3
- [45] Yunlong Lin, Tian Ye, Sixiang Chen, Zhenqi Fu, Yingying Wang, Wenhao Chai, Zhaohu Xing, Lei Zhu, and Xinghao Ding. Aglldiff: Guiding diffusion models towards unsupervised training-free real-world low-light image enhancement. *arXiv preprint arXiv:2407.14900*, 2024. 1
- [46] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 2, 6, 14
- [47] Yun-Fu Liu, Da-Wei Jaw, Shih-Chia Huang, and Jenq-Neng Hwang. Desnownet: Context-aware deep network for snow removal. *IEEE Transactions on Image Processing*, 27(6): 3064–3073, 2018. 3, 4, 15
- [48] Gen Luo, Yiyi Zhou, Tianhe Ren, Shengxin Chen, Xiaoshuai Sun, and Rongrong Ji. Cheap and quick: Efficient vision-language instruction tuning for large language models. *Advances in Neural Information Processing Systems*, 36, 2024. 4
- [49] Ziwei Luo, Fredrik K Gustafsson, Zheng Zhao, Jens Sjölund, and Thomas B Schön. Controlling vision-language models for multi-task image restoration. *arXiv preprint arXiv:2310.01018*, 2023. 1, 7
- [50] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022. 2, 5

- [51] Ozan Özdenizci and Robert Legenstein. Restoring vision in adverse weather conditions with patch-based denoising diffusion models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(8):10346–10357, 2023. 3
- [52] Gaurav Parmar, Taesung Park, Srinivasa Narasimhan, and Jun-Yan Zhu. One-step image translation with text-to-image models. *arXiv preprint arXiv:2403.12036*, 2024. 6, 14, 15, 17
- [53] Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. Gorilla: Large language model connected with massive apis. *arXiv preprint arXiv:2305.15334*, 2023. 1, 2
- [54] Fabio Pizzati, Pietro Cerri, and Raoul de Charette. Physics-informed guided disentanglement in generative networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(8):10300–10316, 2023. 17
- [55] Vaishnav Potlapalli, Syed Waqas Zamir, Salman H Khan, and Fahad Shahbaz Khan. Promptir: Prompting for all-in-one image restoration. *Advances in Neural Information Processing Systems*, 36, 2024. 7
- [56] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*, 2023. 2
- [57] Ruijie Quan, Xin Yu, Yuanzhi Liang, and Yi Yang. Removing raindrops and rain streaks in one go. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9147–9156, 2021. 4, 15
- [58] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 14, 20
- [59] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Acdc: The adverse conditions dataset with correspondences for semantic driving scene understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10765–10775, 2021. 3, 14
- [60] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36, 2024. 2
- [61] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face. *Advances in Neural Information Processing Systems*, 36, 2024. 3, 4
- [62] Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Haozhe Chi, Xun Guo, Tian Ye, Yanting Zhang, et al. Moviechat: From dense token to sparse memory for long video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18221–18232, 2024. 2
- [63] Feifan Song, Bowen Yu, Minghao Li, Haiyang Yu, Fei Huang, Yongbin Li, and Houfeng Wang. Preference ranking optimization for human alignment. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 18990–18998, 2024. 5
- [64] Haoze Sun, Wenbo Li, Jianzhuang Liu, Haoyu Chen, Renjing Pei, Xueyi Zou, Youliang Yan, and Yujiu Yang. Coser: Bridging image and language for cognitive super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 25868–25878, 2024. 3
- [65] Qiaoyu Tang, Ziliang Deng, Hongyu Lin, Xianpei Han, Qiao Liang, Boxi Cao, and Le Sun. Toolalpaca: Generalized tool learning for language models with 3000 simulated cases. *arXiv preprint arXiv:2306.05301*, 2023. 2
- [66] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023. 3
- [67] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 14
- [68] Ashwin K Vijayakumar, Michael Cogswell, Ramprasad R Selvaraju, Qing Sun, Stefan Lee, David Crandall, and Dhruv Batra. Diverse beam search: Decoding diverse solutions from neural sequence models. *arXiv preprint arXiv:1610.02424*, 2016. 5, 8, 17
- [69] Jianyi Wang, Kelvin CK Chan, and Chen Change Loy. Exploring clip for assessing the look and feel of images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2555–2563, 2023. 6, 17
- [70] Jianyi Wang, Zongsheng Yue, Shangchen Zhou, Kelvin CK Chan, and Chen Change Loy. Exploiting diffusion prior for real-world image super-resolution. *International Journal of Computer Vision*, 132(12):5929–5949, 2024. 3, 6, 15
- [71] Shuhang Wang, Jin Zheng, Hai-Miao Hu, and Bo Li. Naturalness preserved enhancement algorithm for non-uniform illumination images. *IEEE transactions on image processing*, 22(9):3538–3548, 2013. 4
- [72] Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world blind super-resolution with pure synthetic data. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1905–1914, 2021. 3, 6, 15
- [73] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*, 2022. 2, 4, 24, 25
- [74] Yingying Wang, Yunlong Lin, Xuanhua He, Hui Zheng, Keyu Yan, Linyu Fan, Yue Huang, and Xinghao Ding. Learning diffusion high-quality priors for pan-sharpening: A two-stage approach with time-aware adapter fine-tuning. *IEEE Transactions on Geoscience and Remote Sensing*, 2025. 3
- [75] Zhichao Wang, Bin Bi, Shiva Kumar Pentylala, Kiran Ramnath, Sougata Chaudhuri, Shubham Mehrotra, Xiang-Bo

- Mao, Sitaram Asur, et al. A comprehensive survey of llm alignment techniques: Rlhf, rlaif, ppo, dpo and more. *arXiv preprint arXiv:2407.16216*, 2024. 6
- [76] 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. *Advances in neural information processing systems*, 35:24824–24837, 2022. 18
- [77] Yanyan Wei, Zhao Zhang, Jiahuan Ren, Xiaogang Xu, Richang Hong, Yi Yang, Shuicheng Yan, and Meng Wang. Clarity chatgpt: An interactive and adaptive processing system for image restoration and enhancement. *arXiv preprint arXiv:2311.11695*, 2023. 3
- [78] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023. 3
- [79] Haoning Wu, Zicheng Zhang, Weixia Zhang, Chaofeng Chen, Liang Liao, Chunyi Li, Yixuan Gao, Annan Wang, Erli Zhang, Wenxiu Sun, et al. Q-align: Teaching llms for visual scoring via discrete text-defined levels. *arXiv preprint arXiv:2312.17090*, 2023. 18, 19
- [80] Haoning Wu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Annan Wang, Kaixin Xu, Chunyi Li, Jingwen Hou, Guangtao Zhai, et al. Q-instruct: Improving low-level visual abilities for multi-modality foundation models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 25490–25500, 2024. 3, 5, 6, 18, 19
- [81] Rui-Qi Wu, Zheng-Peng Duan, Chun-Le Guo, Zhi Chai, and Chongyi Li. Ridcp: Revitalizing real image dehazing via high-quality codebook priors. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22282–22291, 2023. 3, 6, 15, 16
- [82] Jie Xiao, Xueyang Fu, Aiping Liu, Feng Wu, and Zheng-Jun Zha. Image de-raining transformer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(11): 12978–12995, 2022. 3, 6, 15
- [83] Qingsen Yan, Yixu Feng, Cheng Zhang, Pei Wang, Peng Wu, Wei Dong, Jinqiu Sun, and Yanning Zhang. You only need one color space: An efficient network for low-light image enhancement. *arXiv preprint arXiv:2402.05809*, 2024. 6, 15
- [84] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *arXiv preprint arXiv:2406.09414*, 2024. 16
- [85] Rui Yang, Lin Song, Yanwei Li, Sijie Zhao, Yixiao Ge, Xiu Li, and Ying Shan. Gpt4tools: Teaching large language model to use tools via self-instruction. *Advances in Neural Information Processing Systems*, 36, 2024. 1
- [86] Sidi Yang, Tianhe Wu, Shuwei Shi, Shanshan Lao, Yuan Gong, Mingdeng Cao, Jiahao Wang, and Yujiu Yang. Maniqa: Multi-dimension attention network for no-reference image quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1191–1200, 2022. 6, 17, 18, 19
- [87] Wenhan Yang, Wenjing Wang, Haofeng Huang, Shiqi Wang, and Jiaying Liu. Sparse gradient regularized deep retinex network for robust low-light image enhancement. *IEEE Transactions on Image Processing*, 30:2072–2086, 2021. 3, 4, 14, 15
- [88] Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. Mm-react: Prompting chatgpt for multimodal reasoning and action. *arXiv preprint arXiv:2303.11381*, 2023. 3
- [89] Tian Ye, Yunchen Zhang, Mingchao Jiang, Liang Chen, Yun Liu, Sixiang Chen, and Erkang Chen. Perceiving and modeling density for image dehazing. In *European conference on computer vision*, pages 130–145. Springer, 2022. 1
- [90] Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. *arXiv preprint arXiv:2306.13549*, 2023. 4
- [91] Zhiyuan You, Jinjin Gu, Zheyuan Li, Xin Cai, Kaiwen Zhu, Tianfan Xue, and Chao Dong. Descriptive image quality assessment in the wild. *arXiv preprint arXiv:2405.18842*, 2024. 4
- [92] Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht Madhavan, and Trevor Darrell. Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2636–2645, 2020. 3, 4, 14
- [93] Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. Rrhf: Rank responses to align language models with human feedback without tears. *arXiv preprint arXiv:2304.05302*, 2023. 2, 3, 5, 6, 17
- [94] Kai Zhang, Yawei Li, Jingyun Liang, Jie Zhang Cao, Yulun Zhang, Hao Tang, Deng-Ping Fan, Radu Timofte, and Luc Van Gool. Practical blind image denoising via swin-conv-unet and data synthesis. *Machine Intelligence Research*, 20(6):822–836, 2023. 3, 6, 15
- [95] Weixia Zhang, Guangtao Zhai, Ying Wei, Xiaokang Yang, and Kede Ma. Blind image quality assessment via vision-language correspondence: A multitask learning perspective. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14071–14081, 2023. 6, 17
- [96] Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*, 2022. 18
- [97] Lirui Zhao, Yue Yang, Kaipeng Zhang, Wenqi Shao, Yuxin Zhang, Yu Qiao, Ping Luo, and Rongrong Ji. Diffagent: Fast and accurate text-to-image api selection with large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6390–6399, 2024. 2
- [98] Zhonghan Zhao, Wenhao Chai, Xuan Wang, Boyi Li, Shengyu Hao, Shidong Cao, Tian Ye, and Gaoang Wang. See and think: Embodied agent in virtual environment. In *European Conference on Computer Vision*, pages 187–204. Springer, 2025. 1
- [99] Shangchen Zhou, Chongyi Li, and Chen Change Loy. Led-net: Joint low-light enhancement and deblurring in the dark.

- In *European conference on computer vision*, pages 573–589. Springer, 2022. [3](#), [14](#), [15](#)
- [100] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. [14](#)
- [101] Kaiwen Zhu, Jinjin Gu, Zhiyuan You, Yu Qiao, and Chao Dong. An intelligent agentic system for complex image restoration problems. *arXiv preprint arXiv:2410.17809*, 2024. [3](#)
- [102] Jannik Zörn, Paul Gladkov, Sofia Dudas, Fergal Cotter, Sofi Toteva, Jamie Shotton, Vasiliki Simaiaki, and Nikhil Mohan. Wayvescenes101: A dataset and benchmark for novel view synthesis in autonomous driving. *arXiv preprint arXiv:2407.08280*, 2024. [3](#)



JarvisIR: Elevating Autonomous Driving Perception with Intelligent Image Restoration

Supplementary Material

This is the supplementary material for the paper: “JarvisIR: Elevating Autonomous Driving Perception with Intelligent Image Restoration.” We provide the following materials in this manuscript:

- Sec.8 More implementation details.
 - Restoration tool settings.
 - Details of Model Setups.
- Sec.9 CleanBench dataset details.
 - Dataset statistics.
 - Details of degradation library.
- Sec.10 More ablation.
 - MRRHF vs. vanilla RRHF.
 - Sample generation strategy and entropy regularization.
 - Effectiveness of differentiated contrast weights.
 - Impact of reasoning for decision-making.
 - Impact of reward model.
- Sec.11 More visual results.
- Sec.12 Limitations, broader impacts and future work.

8. More implementation details

8.1. Restoration tool settings

Table 4 lists the task-specific restoration tools used in our implementation. Notably, some models lack weights corresponding to certain tasks but are inherently adaptable; we collect appropriate data to retrain them. For example, Img2img-turbo [52] is an image-to-image translation method based on SD-turbo that provides night-to-day and rainy-to-day weights but not snow-to-day weights. To enable Img2img-turbo to adapt to snow scenes, we retrain it using the CycleGAN paradigm on the snow subset of the ACDC dataset [59]. Additionally, it is important to note that we are not utilizing the latest state-of-the-art tools, suggesting considerable potential for enhancing our models.

8.2. Details of Model Setups

Model Architecture. In this study, JarvisIR primarily adopts the architecture from Llava-Llama3-8B [46]. Specifically, the input images and instruction texts are first tokenized, then fused, and finally processed by the Large Language Model (LLM) for response generation. (a) Tokenization of input images and instruction texts: We use a frozen CLIP pre-trained ViT-L/14 [58] as the image encoder to convert input images into visual tokens. The instruction texts are tokenized into textual tokens using the SentencePiece tokenizer [38]. To bridge the different em-

bedding spaces of visual and textual tokens, we implement a trainable image projector to map visual tokens into the textual space, following [67, 100]. (b) Token Fusion: We integrate the visual tokens into predefined positions within the textual tokens to achieve token fusion. (c) Response Generation Using LLM: The fused tokens are fed into the LLM to generate the final response. In our experiments, we primarily use Llama3-8B [67]. Even with their advanced features, pre-trained LLMs lack the ability to furnish accurate responses, thorough reasoning regarding degradation, and precise restoration plans without dataset-specific fine-tuning. Therefore, we employ a full parameter fine-tuning technique that efficiently unleashes the potential of LLM to the maximum extent.

Model setup. Since the CLIP pre-trained ViT-L/14 [58] encodes each 14×14 image patch into a visual token, the input image dimensions must be integer multiples of 14. Therefore, we zero-pad the input images to meet this requirement. We encode the image patches into visual tokens using the CLIP pre-trained ViT-L/14 [58], where each token is a 1024-dimensional vector. These visual tokens are subsequently projected by the image projection layer into the LLM’s hidden dimension of 4096.

Training setup. Both the SFT and MRRHF tuning phases utilize the Adam optimizer with learning rate $1e-5$ with cosine decay. The warmup ratio is set to 0.03, the maximum sequence length is 2048, and the weight decay is 4. JarvisIR-SFT undergoes training for three epochs with a batch size of 128, while JarvisIR-MRRHF is trained for three epochs using a batch size of 2. During the MRRHF tuning phase, the diverse beam search settings include a size of 3, 5 beam groups, a diversity penalty of 2.0, and a sampling temperature of 0.8. Training is conducted on 8 GPUs (NVIDIA A100 80G).

9. CleanBench dataset details

9.1. Dataset statistics

CleanBench. In constructing the CleanBench process, we collected large-scale raw daytime images from various sources, including autonomous driving datasets [4, 4, 59, 92] and natural datasets [5, 34, 39, 42, 87, 99]. The CleanBench dataset contains a total of 150K degraded-clean image pairs. For the construction of CleanBench-Real, we gathered 80K real degraded images consisting of night scenes, fog scenes, snow scenes and rain scenes. These data

Table 4. Task-specific restoration tools with descriptions.

Task	Tools	Model Description
Super-resolution	StableSR-turbo [70]	Utilizes pre-trained diffusion models with a time-aware encoder for high-quality super-resolution, deblurring, and artifact removal.
	Real-ESRGAN [72]	Fast GAN for super-resolution, deblurring, and artifact removal, handling complex real-world degradations efficiently.
Denoising	SCUNet [94]	Hybrid UNet-based model combining convolution and transformer blocks, designed for robust denoising under diverse real-world noise conditions.
Compression artifact removal	StableSR-turbo [70]	Utilizes pre-trained diffusion models with a time-aware encoder for high-quality super-resolution, deblurring, and artifact removal.
	Real-ESRGAN [72]	Fast GAN for super-resolution, deblurring, and artifact removal, handling complex real-world degradations efficiently.
Deblurring	StableSR-turbo [70]	Utilizes pre-trained diffusion models with a time-aware encoder for high-quality super-resolution, deblurring, and artifact removal.
	Real-ESRGAN [72]	Fast GAN for super-resolution, deblurring, and artifact removal, handling complex real-world degradations efficiently..
Deraining	IDT [82]	Transformer-based model for de-raining and raindrop removal.
	UDR-S2Former [8]	An uncertainty-aware transformer model for rain streak removal.
	Img2img-turbo-rain [52]	Efficient model based on SD-turbo, designed for fast and effective rain removal in real-world images.
Raindrop removal	IDT [82]	Transformer-based model for de-raining and raindrop removal.
Dehazing	RIDCP [81]	Efficient dehazing model utilizing high-quality codebook priors to handle complex real-world haze.
	KANet [21]	Efficient dehazing network using a localization-and-removal pipeline to handle complex real-world hazy.
Desnowing	Img2img-turbo-snow [52]	Efficient model for removing snow artifacts while preserving natural scene details.
	Snowformer [11]	Transformer-based model for removing snowflakes while preserving natural scene details.
Low-light enhancement	Img2img-turbo-night [52]	Fast and efficient model based on SD-turbo, designed for low-light enhancement in real-world scenarios.
	HVI-CIDNet [83]	Lightweight transformer for low-light and exposure correction, enhancing both image quality and downstream vision tasks efficiently.
	LightenDiff [27]	Diffusion-based framework for low-light enhancement, leveraging Retinex theory and latent-space decomposition for high-quality unsupervised restoration.

come from diverse sources, including the aforementioned autonomous driving datasets. Additionally, to enhance the generalizability of JarvisIR in natural contexts, we incorporated natural adverse weather scenes from internet and public datasets [5, 34, 39, 42, 47, 57, 87, 99].

9.2. Details of degradation library

As described in Sec 3.1 of the manuscript, we simulate realistic adverse weather scenarios—rainy, nighttime, snowy, and foggy conditions—by customizing a degradation library developed with physical models and image transformation techniques to synthesize degraded images. In this section, we detail our degradation implementations, cov-

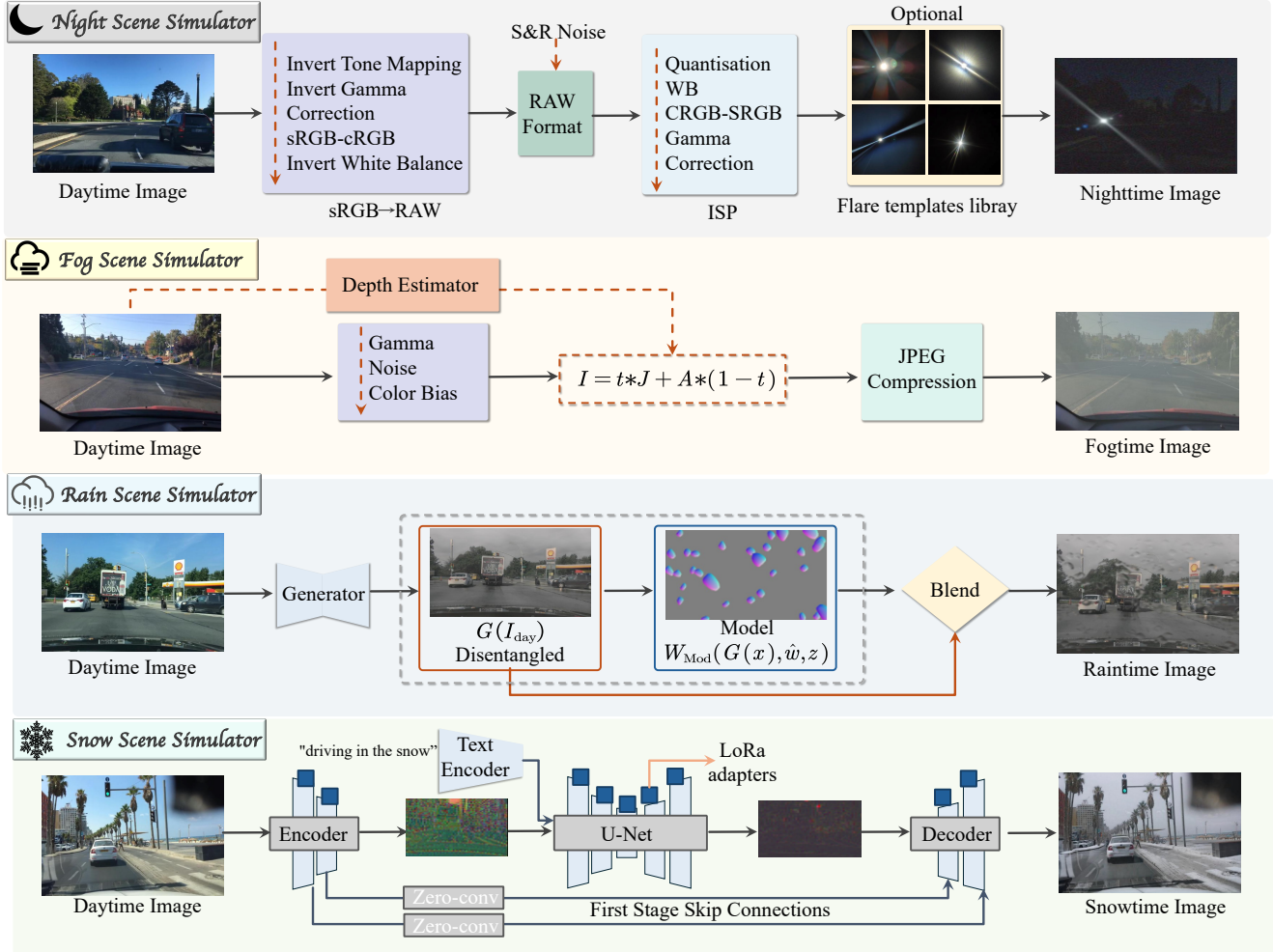


Figure 8. Adverse weather scene simulator. To simulate realistic adverse weather scenarios, including rainy, nighttime, snowy, and foggy, we customized the degradation library developed using physical models and image transformation techniques to synthesize degraded images.

ering the principles, formulas, and severity setups for the Night Scene Simulator, Fog Scene Simulator, Rain Scene Simulator, and Snow Scene Simulator. Examples for each implementation are provided in Figure 11.

Night Scene Simulator. Inspired by the work of [19], we employ a low-light degradation transform to synthesize realistic low-light images, denoted as T_{night} , as illustrated in Figure 8. Specifically, we first convert the daytime image I_{day} into RAW data using the sRGB→RAW process [3]. Next, we linearly attenuate the RAW image and introduce Shot and Read (S&R) noise, which is commonly observed in camera imaging systems [3]. Finally, we apply the Image Signal Processing (ISP) pipeline to convert the low-light sensor data back into sRGB format. Additionally, we incorporate flare degradation using flare templates from the Flare7K++ [20] dataset. The complete low-light degrada-

tion transform T_{night} is given by:

$$T_{night}(I_{day}) = T_{ISP}(T_{sRGB \rightarrow RAW}(I_{day}) + I_{noise}) + I_{flare}, \quad (8)$$

which generates a degraded image I_{day} that closely resembles a dark nighttime scene. Furthermore, we use an online dynamic degradation process. It applies randomized parameter combinations, as defined in Equation 8, to simulate diverse nighttime driving conditions.

Fog Scene Simulator. Inspired by RIDCP [81], we design a foggy image degradation transform, denoted as T_{fog} , to synthesize realistic hazy images, as shown in Figure 8. Specifically, we simulate fog by introducing transmission maps $t(x)$ using depth estimation algorithms (e.g., Depth anything V2 [84]), combined with exponential attenuation $e^{\beta d(x)}$, where β controls haze density within the range [0.3, 1.5]. Additionally, poor lighting conditions are

modeled by applying a brightness adjustment factor $\gamma \in [1.5, 3.0]$, Gaussian noise \mathcal{N} , and atmospheric light variation $A + \Delta A$, where ΔA is sampled from $[-0.025, 0.025]$. To further enhance realism, JPEG compression artifacts are introduced by applying JPEG (\cdot) to the degraded image. The complete foggy image synthesis process is defined as:

$$T_{\text{fog}}(I_{\text{day}}) = \text{JPEG} \left(\mathcal{P} \left(I_{\text{day}}^\gamma + \mathcal{N}, e^{\beta d(x)}, A + \Delta A \right) \right), \quad (9)$$

where \mathcal{P} represents the hazy image formation process, I_{day} is the clean image, and $d(x)$ is the estimated depth map. The variable x refers to the spatial coordinates of the image. This dynamic degradation process is designed to operate online with randomized parameters, simulating diverse real-world foggy conditions.

Rain Scene Simulator. Inspired by PGDGN [54], we introduce a rain degradation transform, denoted as T_{rain} , to generate realistic rainy images (Figure 8). This transform synthesizes rainy images by combining a disentangled clean image with a physics-based rain rendering model. The degradation process is formulated as:

$$T_{\text{rain}}(I_{\text{day}}) = W_{\text{Mod}}(G(I_{\text{day}}), \hat{w}, z), \quad (10)$$

where I_{day} is the clean image, $G(I_{\text{day}})$ represents the disentangled base image, and W_{Mod} is the rain rendering model. W_{Mod} incorporates parameters $\hat{w} = \{\hat{w}_d, \hat{w}_{nd}\}$, with \hat{w}_d controlling differentiable aspects such as raindrop size and streak density, and \hat{w}_{nd} addressing nondifferentiable properties. The term z introduces stochastic noise for variability in rain effects. This process applies W_{Mod} to add realistic raindrop occlusions, rain streaks, and scene wetness to the disentangled image. $G(I_{\text{day}})$, generating a visually plausible rainy image $T_{\text{rain}}(I_{\text{day}})$ with controlled and diverse effects.

Snow Scene Simulator. Building on the img2img-turbo model [52], we introduce a snow transformation, denoted as T_{snow} , to generate realistic snowy images from daytime inputs. This process uses the SD-Turbo model with textual conditioning. It synthesizes snowy scenes by combining the input image with a latent diffusion-based generator and a textual prompt. The snow transformation is formulated as:

$$T_{\text{snow}}(I_{\text{day}}, C_{\text{snow}}) = G_{\text{snow}}(I_{\text{day}}, C_{\text{snow}}), \quad (11)$$

where I_{day} is the daytime input image, C_{snow} is the textual condition (e.g., “driving in the heavy snow”), and G_{snow} represents the generator. By employing LoRA adapters and skip connections, the generator enables precise control over scene characteristics while maintaining the structural integrity of the input image. This process applies G_{snow} to infuse the daytime image I_{day} with snowy features, guided by the contextual information in C_{snow} . The resulting synthetic image aligns closely with the visual expectations of a snowy environment while maintaining consistency with the original scene’s structure.

10. More ablation

To thoroughly investigate the proposed JarvisIR, we conducted an extensive array of ablation studies on the CleanBench-Real dataset. Four non-reference metrics are used for assessment: MUSIQ [35], MANIQA [86], CLIP-IQA+ [69], LIQE [95]. The specific elements of these studies are further expounded in the sections that follow.

10.1. MRRHF vs. vanilla RRHF

We evaluate the effectiveness of our proposed MRRHF by comparing it with vanilla RRHF [93]. The reward and diversity metrics over training iterations are illustrated in Table 5. Fine-tuning JarvisIR with MRRHF significantly improves the average values of both reward and diversity by 0.19 and 3.43, respectively, compared to using RRHF. The degradation in diversity and reward when using vanilla RRHF results from its offline sample generation strategy. As discussed in Sec. 5 in the manuscript, this strategy confines its generated samples to the finite sample space created by the SFT model using diverse beam search [68]. In contrast, our MRRHF employs a hybrid sample generation strategy and entropy regularization, providing sufficient sample exploration space to achieve globally optimal results.

10.2. Sample generation strategy and entropy regularization

In our manuscript, we examine the effects of the sample generation strategy and entropy regularization on the MRRHF tuning process, focusing on reward scores and response diversity. This section provides further evidence of the effectiveness of our hybrid sample generation strategy and entropy regularization. Specifically, as shown in Table 5, we assess their impact on performance using the CleanBench-Real validation set. The results demonstrate that our hybrid sampling approach and entropy regularization not only enhance training stability and facilitate high-quality exploration of the optimization space but also significantly improve testing performance.

10.3. Effectiveness of differentiated contrast weights

In Equation 4 of our manuscript, we refine the original ranking loss [93] by introducing differentiated contrast weights, expressed as $L_{\text{rank}} = \sum_{s_i < s_j} (s_j - s_i) \max(0, p_i - p_j)$. The term $(s_j - s_i)$ represents the differentiated contrast weights. We compare this with the original ranking loss $\hat{L}_{\text{rank}} = \sum_{s_i < s_j} \max(0, p_i - p_j)$. Table 5 presents the reward and diversity metrics over training iterations. When the ranking loss is applied without differentiated contrast weights \hat{L}_{rank} the average values of both reward and diversity decrease by 0.14 and 3.93, respectively, compared

Table 5. Ablation studies on tuning paradigm, differentiated contrast weights, different sample generation strategies, and entropy regularization. The ‘‘Reward’’ represents the average reward scores obtained during alignment tuning, spanning from -1 to 1. A negative score indicates a penalty, while a positive score represents a reward. The ‘‘Diversity’’ reflects the average number of unique responses produced during the training process. Additionally, we evaluate performance on the CleanBench-Real validation set using four non-reference metrics: MUSIQ, MANIQA, CLIP-IQA+, and LIQE. The reported values represent the average performance across all tested scenes.

Strategy	Reward	Diversity	MUSIQ \uparrow	MANIQA \uparrow	CLIP-IQA+ \uparrow	LIQE \uparrow
Vanilla RRHF	0.40	3.12	63.89	0.5090	0.5388	3.589
MRRHF (Ours)	0.67	6.55	71.43	0.7099	0.7296	4.411
w/o. differentiated contrast weights	0.53	2.62	63.22	0.5871	0.6130	3.597
w. differentiated contrast weights (Ours)	0.67	6.55	71.43	0.7099	0.7296	4.411
offline sample generation	0.43	3.63	64.12	0.5323	0.6012	3.620
online sample generation	-0.87	1.27	-	-	-	-
hybrid sample generation (Ours)	0.67	6.55	71.43	0.7099	0.7296	4.411
w/o. entropy regularization	0.50	4.56	65.06	0.6207	0.6915	3.867
w. entropy regularization (Ours)	0.67	6.55	71.43	0.7099	0.7296	4.411

to using L_{rank} . We attribute this to the differentiated contrast weights enabling the VLM to recognize that some negative examples are neutral (with reward scores close to positive examples) and thus should not be excessively penalized, which helps prevent confusion during VLM training. Specifically, assuming the system uses diverse beam search to obtain multiple responses $r_1, \dots, r_i, r_k, r_n$ the original RRHF algorithm treats the best response r_k as positive and the remaining responses $r_i < r_k$ as negative examples of r_k and applies the same penalty to them. However, this approach may not be reasonable, especially when the preference scores of different r_i are similar. For instance, when the preference of r_{k+1} is only slightly worse than r_k , while r_n is significantly worse than r_k , the model should differentiate and apply different penalty strengths, slightly penalizing r_{k+1} and heavily penalizing r_n compared to r_k . To address this, we propose using the score $\mathcal{S}(r_i)$ from a reward model $\mathcal{S}(\cdot)$ to indicate the numerical preference of r_i , i.e., the differentiated contrast weights ($s_j - s_i$).

10.4. Impact of reasoning for decision-making

As the pioneering work [76] points out, Chain-of-Thought (CoT) is ‘‘a series of intermediate reasoning steps’’ that has proven effective in complex reasoning tasks [36, 76, 96]. The main idea of CoT is to prompt large language models (LLMs) to output not only the final answer but also the reasoning process leading to it, resembling human cognitive processes. Inspired by this approach, we enable JarvisIR to provide detailed degradation and reasoning insights about the degraded image before making decisions, specifically before producing the task sequence with model selection. To assess the impact of reasoning on final decision-making, we perform ablation experiments on the CleanBench-Real validation set by comparing two variants: (1) directly requesting JarvisIR to output the task sequences, and (2) providing detailed degradation and reasoning insights before

Table 6. Ablation studies on the impact of reasoning for decision-making. We evaluate performance on the CleanBench-Real validation set using four non-reference metrics: MUSIQ, MANIQA, CLIP-IQA+, and LIQE. The reported values represent the average performance across all tested scenes.

Configurations	MUSIQ \uparrow	MANIQA \uparrow	CLIP-IQA+ \uparrow	LIQE \uparrow
w/o. reasoning	71.17	0.6942	0.7156	4.394
(Ours) w reasoning	71.43	0.7099	0.7296	4.411

outputting the task sequences. As shown in Table 6, providing detailed degradation and reasoning insights significantly enhances JarvisIR’s decision-making, leading to notable improvements in the four non-reference metrics. By explicitly describing degradations and reasoning insights, the model can use in-context learning to align selected tasks and restoration experts with the specific degradations present. This strategy not only enhances interpretability but also introduces constraints that make the model’s decisions more reliable in real-world scenarios.

10.5. Impact of reward model

To analyze how various reward model configurations affect model optimization, we conducted an ablation experiment exploring three distinct settings: (I) multiple VLM-based IQA models as a unified reward model (e.g., Q-instruct [80] and Q-align [79]). (II) using a single VLM-based IQA model (e.g., Q-Instruct [80] or Q-align [79]) or a traditional IQA model (e.g., MUSIQ [35] or MANIQA [86]). (III) multiple traditional IQA models as a unified model (e.g., MUSIQ [35] and MANIQA [86]). The results of JarvisIR-MRRHF trained with different reward models are summarized in Table 7. Based on the results, we make the following observations: (1) Using multiple VLM-based IQA models as the reward model significantly improves perception metrics, although it increases resource consumption during training. (2) Training with a

Table 7. Ablation studies on the impact of different reward model configurations. We evaluate performance on the CleanBench-Real validation set using four non-reference metrics: MUSIQ, MANIQA, CLIP-IQA+, and LIQE. The reported values represent the average performance across all tested scenes.

Configurations	MUSIQ \uparrow	MANIQA \uparrow	CLIP-IQA+ \uparrow	LIQE \uparrow
(I) Q-align [79] + Q-Instruct [80]	71.41	0.7094	0.7308	4.419
(II) Q-align	71.35	0.7086	0.7288	4.409
(II) Q-Instruct [80]	71.37	0.7093	0.7257	4.402
(II) MUSIQ [35]	71.64	0.6932	0.6977	3.955
(II) MANIQA [86]	68.49	0.7126	0.6805	3.981
(III) MUSIQ [35] + MANIQA [86]	71.52	0.7118	0.7068	4.127
(Ours) Q-Instruct [80]+ MUSIQ [35] + MANIQA [86]	71.43	0.7099	0.7296	4.411

single IQA model improves the corresponding metric significantly, but other metrics may experience some degradation. (3) Combining multiple traditional IQA models as the reward model enhances performance on certain metrics, but the improvements are asymmetrical—some traditional metrics exhibit very high performance while perception metrics are relatively low. Consequently, we opt to create the unified reward model by combining both VLM-based and non-VLM-based IQA models, such as Q-instruct [80], MUSIQ [35], and MANIQA [86]. This combination allows for a comprehensive evaluation of system responses while preserving training efficiency.

11. More visual results.

11.1. Perception restoration

Additional visual comparisons highlight the effectiveness of the proposed JarvisIR framework in real-world adverse weather conditions. Figure 9 illustrates the comprehensive workflow of JarvisIR, which begins by receiving user commands and degraded images. JarvisIR evaluates the image quality, identifies degradation factors, and formulates task sequences. It then selects appropriate models for tasks such as denoising, dehazing, and super-resolution. The outputs include evaluated inference insights, detailed restoration plans, and enhanced images, effectively bridging user instructions with image restoration plans.

Figure 10 illustrates the decision-making processes of both JarvisIR-MRRHF and JarvisIR-SFT. Experimental results indicate that the decision-making capability of JarvisIR-MRRHF surpasses that of JarvisIR-SFT. Specifically, JarvisIR-MRRHF makes correct decisions in cases where JarvisIR-SFT previously failed. For example, in coupled degraded real rain scenarios (the first row), JarvisIR-SFT yields a mediocre decision—“Enhancement (Img2img-turbo) \rightarrow Dehaze (RIDCP) \rightarrow DeRaindrop (IDT)” —which does not remove raindrops and blur the background. However, JarvisIR-MRRHF accurately identifies the appropriate restoration tasks and selects the opti-

mal models to solve them: “Denoise (SCUNet) \rightarrow DeRaindrop (IDT) \rightarrow Deblur (StableSR-turbo)”. This improvement confirms that MRRHF fine-tuning significantly enhances JarvisIR’s decision-making ability under real-world conditions, reduces hallucination errors, and improves generalization performance.

Figures 12, 13, 14, and 15 illustrate visual comparisons of our method and the baseline methods across four different scenes on the CleanBench-Real test set. Our results demonstrate that JarvisIR outperforms the comparative methods in terms of color enhancement, detail preservation, and the elimination of degradations, achieving a superior balance among these aspects. Conversely, the baseline methods perform poorly in real-world environments. They struggle to handle coupled degradations that occur simultaneously in natural settings, such as low light combined with fog or a mixture of rain and fog. These limitations may arise from their heavy dependence on specific degradation priors and significant domain gaps due to mismatches between synthetic training data distributions and real-world data. Consequently, they often produce subpar recovery results featuring artifacts, overexposure, underexposure, and amplified noise.

12. Limitations, broader impacts and future work

The primary limitation of our research is that JarvisIR is unable to address all real-world restoration scenarios. While it demonstrates effectiveness in handling most degradation scenarios relevant to autonomous driving, it does not extend to tasks such as underwater image restoration, old photo enhancement, or blind face restoration. By incorporating appropriate data and tools, rapid adaptation could be achieved through the proposed training paradigm. Furthermore, the tools currently employed are limited in scope and capability. In our future work, we will incorporate more advanced and robust restoration tools that might further enhance JarvisIR’s ability to address real-world coupled degradation challenges.

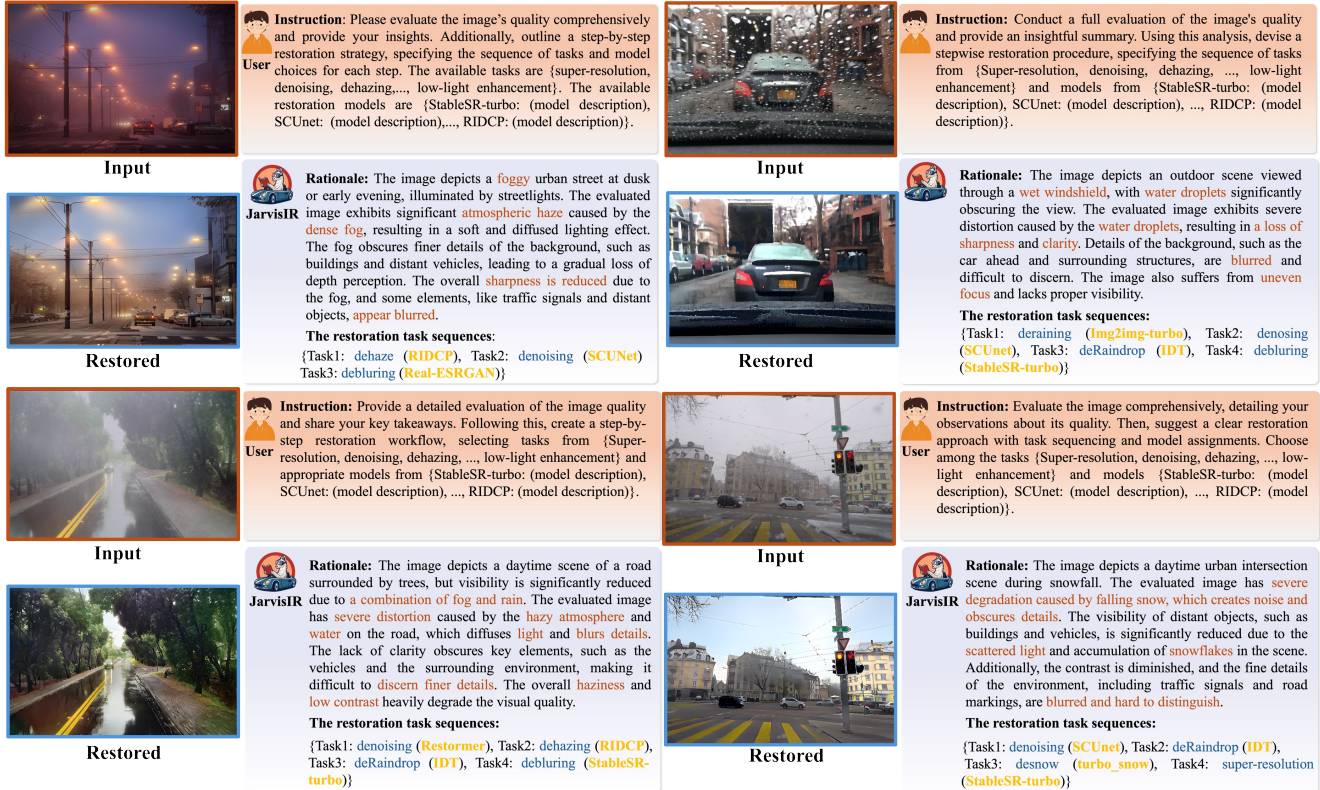


Figure 9. More examples of JarvisIR’s perception restoration are presented. Initially, JarvisIR assesses the degradation of the input images and parses user instructions to formulate a task plan, selecting appropriate expert models for each subtask. The selected experts perform their designated tasks and return the results to JarvisIR, which integrates the outcomes and provides the final answer to the user.



Figure 10. Comparison of the decision-making processes of JarvisIR-MRRHF and JarvisIR-SFT. The results indicate that the MRRHF version accurately predicts the correct task sequence and selects appropriate restoration models. Conversely, the SFT version often fails to make suitable decisions in real-world scenarios due to the domain gap between training and real data distributions.

Another future work could focus on retaining the original image resolution during training. Most current vision-language models (VLMs) resize input images to a fixed resolution, such as 336×336 , which may degrade performance, as resolution variation may affect the model’s perception of degradation. To mitigate this, future research could explore techniques to maintain original image resolutions. One approach involves adapting the position embeddings in CLIP [58] using bicubic interpolation to accommodate varying image dimensions.

This work focuses on building an autonomous, robust,

intelligent restoration system tailored for real-world challenges. To enhance system robustness, reduce hallucinations, and improve generalizability, we introduce a novel two-stage framework that integrates supervised fine-tuning with human feedback alignment. By utilizing human feedback and large-scale real unlabeled data, our method allows the VLM to be fine-tuned in an unsupervised manner. We believe that this paradigm can inspire future work to build more powerful and versatile intelligent systems.

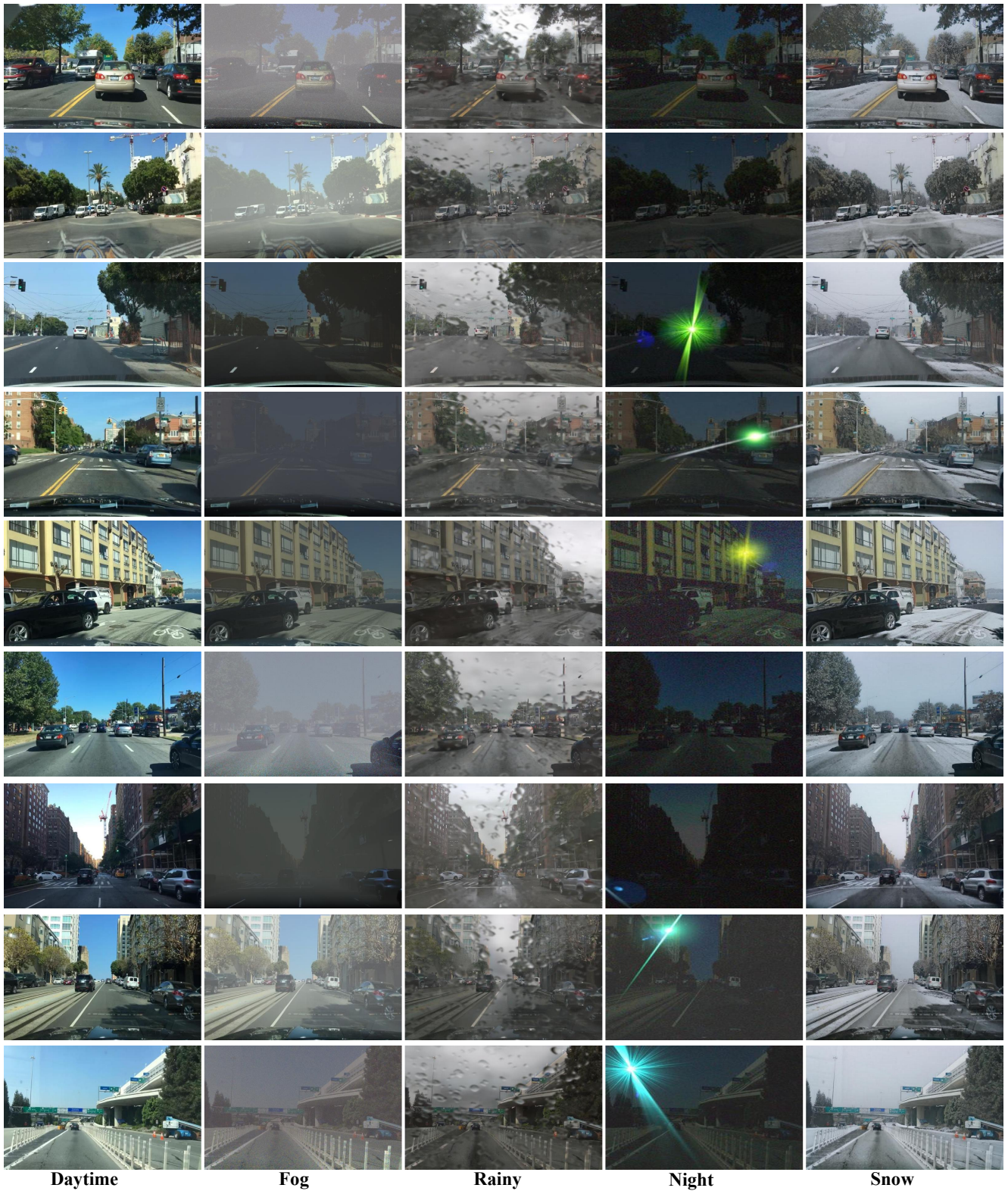


Figure 11. Examples of synthetic adverse weather scenarios in autonomous driving from the CleanBench dataset.

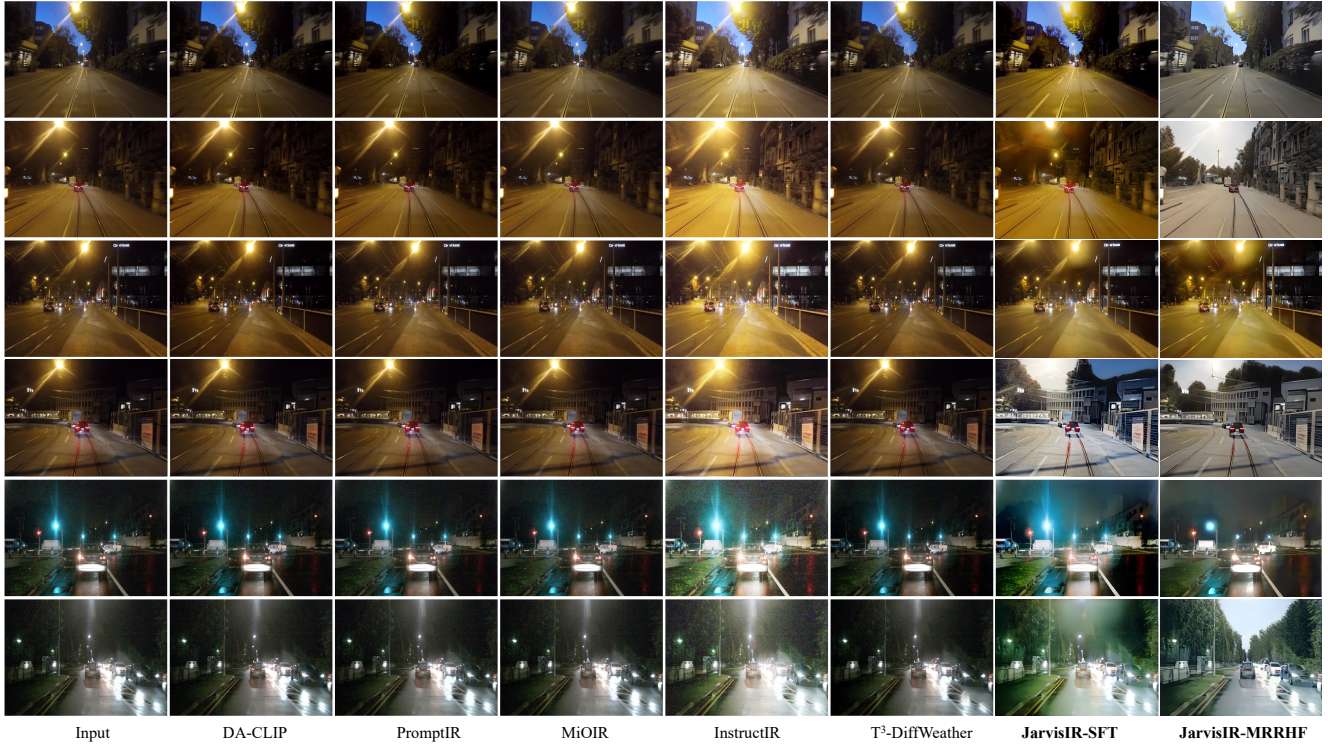


Figure 12. Visual comparisons among various methods on CleanBench-Real’s night scene validation set.

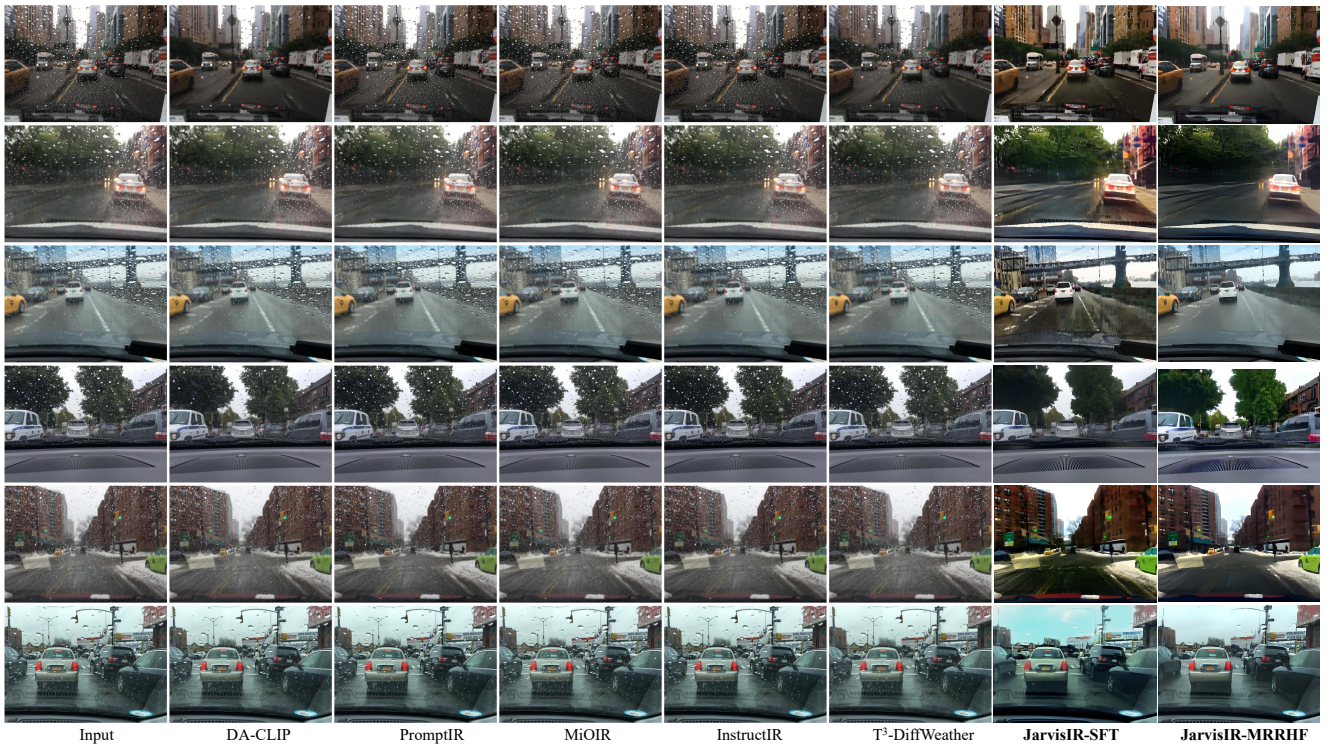


Figure 13. Visual comparisons among various methods on CleanBench-Real’s rain scene validation set.

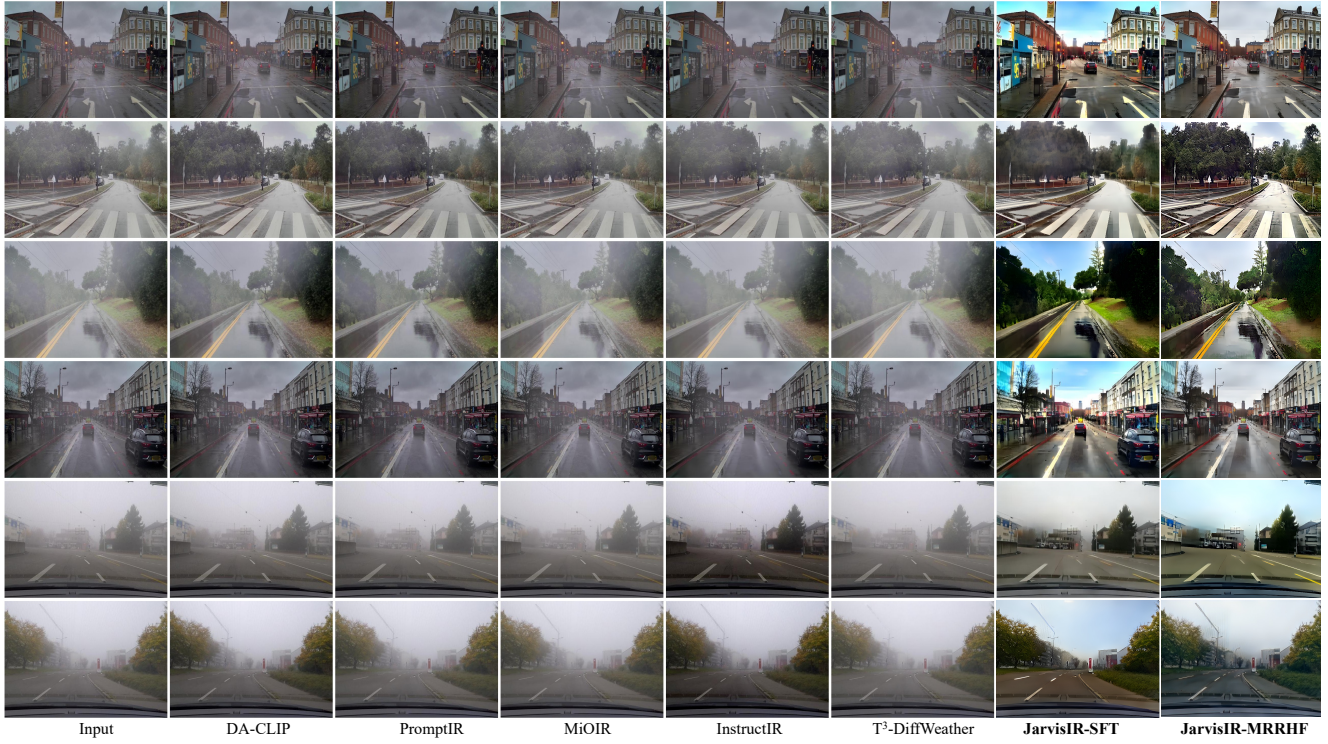


Figure 14. Visual comparisons among various methods on CleanBench-Real’s fog scene validation set.

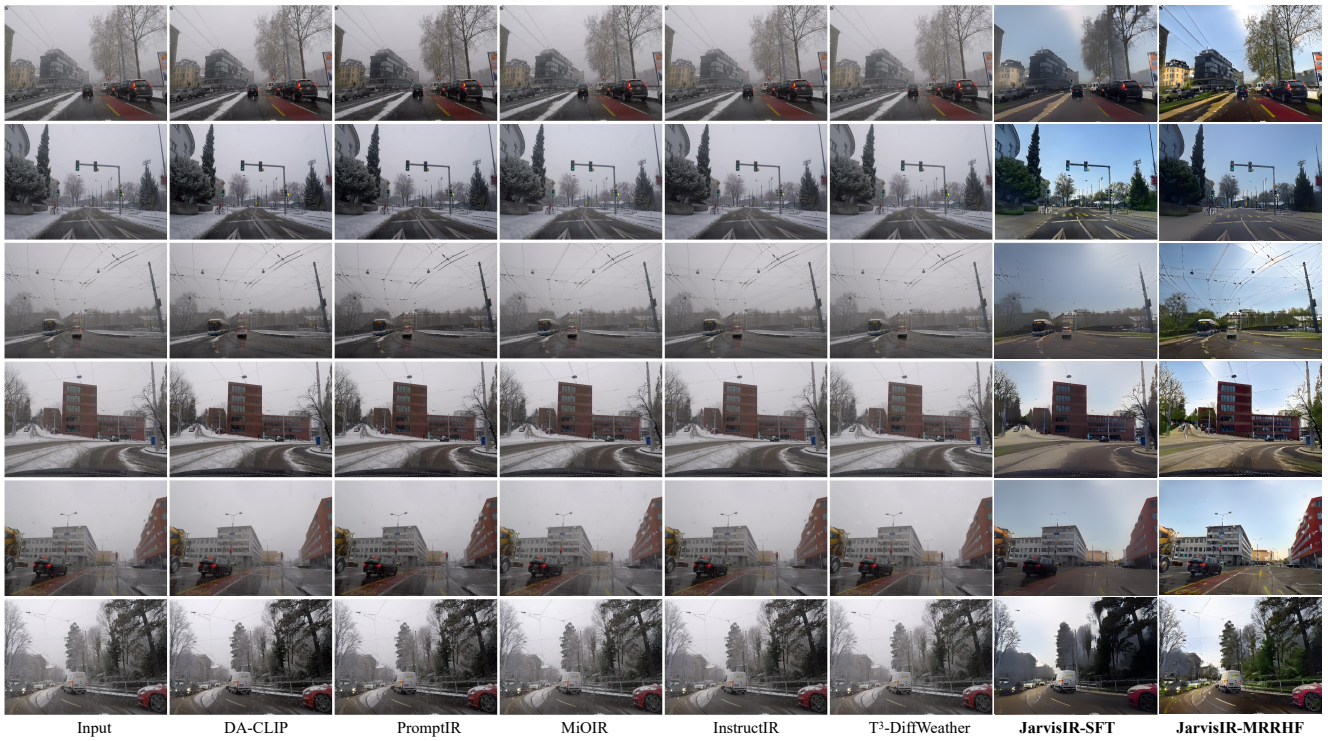


Figure 15. Visual comparisons among various methods on CleanBench-Real’s snow scene validation set.

Table 8. Instruction generated by GPT-4V using the self-instruct strategy [73]

#	Instruction
1	Please evaluate the image’s quality comprehensively and provide your insights. Additionally, outline a step-by-step restoration strategy, specifying the sequence of tasks and model choices for each step. The available tasks are super-resolution, denoising, dehazing,...., low-light enhancement. The available restoration models are StableSR-turbo: (model description), SCUnet: (model description),...., RIDCP: (model description).
2	Analyze the quality of the image comprehensively and provide your insights. Furthermore, propose a restoration strategy by detailing each task and model choice sequentially. The available tasks include super-resolution, denoising, dehazing,...., low-light enhancement. The restoration models are StableSR-turbo: (model description), SCUnet: (model description),...., RIDCP: (model description).
3	Assess the overall quality of the image and provide a detailed evaluation. Then, design a step-by-step restoration process, specifying tasks and model choices. Tasks available are super-resolution, denoising, dehazing,...., low-light enhancement, and models include StableSR-turbo: (model description), SCUnet: (model description),...., RIDCP: (model description).
4	Perform a comprehensive evaluation of the image quality and explain your observations. Additionally, develop a step-by-step restoration plan, identifying tasks and model choices. Available tasks are super-resolution, denoising, dehazing,...., low-light enhancement, and models include StableSR-turbo: (model description), SCUnet: (model description),...., RIDCP: (model description).
5	Conduct a thorough analysis of the image’s quality and provide your insights. Subsequently, create a restoration strategy step by step, specifying the tasks and model choices. The tasks available are super-resolution, denoising, dehazing,...., low-light enhancement, and models are StableSR-turbo: (model description), SCUnet: (model description),...., RIDCP: (model description).
6	Evaluate the quality of the image comprehensively and outline your findings. Moreover, formulate a sequential restoration plan, detailing tasks and model selections. Available tasks include super-resolution, denoising, dehazing,...., low-light enhancement, and models are StableSR-turbo: (model description), SCUnet: (model description),...., RIDCP: (model description).
7	Provide a detailed assessment of the image’s quality and share your observations. Then, create a restoration strategy in a step-by-step manner, specifying tasks and models. Available tasks are super-resolution, denoising, dehazing,...., low-light enhancement, and models include StableSR-turbo: (model description), SCUnet: (model description),...., RIDCP: (model description).
8	Examine the image comprehensively and describe its quality in detail. Additionally, outline a step-wise restoration process, specifying the tasks and models for each stage. Tasks available are super-resolution, denoising, dehazing,...., low-light enhancement, and models are StableSR-turbo: (model description), SCUnet: (model description),...., RIDCP: (model description).
9	Review the image quality comprehensively and present your findings. Furthermore, propose a restoration plan step by step, detailing the tasks and models used. Available tasks are super-resolution, denoising, dehazing,...., low-light enhancement, and restoration models are StableSR-turbo: (model description), SCUnet: (model description),...., RIDCP: (model description).
10	Perform a detailed evaluation of the image and provide your insights. Then, construct a step-by-step restoration strategy, specifying the sequence of tasks and the corresponding models. The available tasks include super-resolution, denoising, dehazing,...., low-light enhancement, and models are StableSR-turbo: (model description), SCUnet: (model description),...., RIDCP: (model description).

Table 9. Responses generated by GPT-4V using the self-instruct strategy [73]

#	Response
1	The image depicts an outdoor scene viewed through a wet windshield, with water droplets significantly obscuring the view. The evaluated image exhibits severe distortion caused by the water droplets, resulting in a loss of sharpness and clarity. Details of the background, such as the car ahead and surrounding structures, are blurred and difficult to discern. The image also suffers from uneven focus and lacks proper visibility.
2	This image portrays an outdoor scene as seen through a wet windshield, where water droplets heavily obstruct the view. The distortion caused by these droplets leads to a noticeable reduction in sharpness and clarity. Background elements, like the car in front and nearby structures, appear heavily blurred and challenging to distinguish. The image further suffers from uneven focus, reducing overall visibility.
3	An outdoor setting is captured in this image, viewed through a wet windshield. The water droplets significantly obscure the scene, causing severe distortion. Sharpness and clarity are greatly reduced, and background features, such as the car ahead and surrounding buildings, are blurred and hard to recognize. Uneven focus further detracts from the image's overall visibility.
4	This image showcases an outdoor scene as observed through a wet windshield, with water droplets causing significant obstruction. The distortion leads to substantial clarity and sharpness loss. Details in the background, including the car ahead and nearby structures, are blurred and hard to identify. Additionally, the image has uneven focus, which impairs its overall visibility.
5	The image illustrates an outdoor view observed through a wet windshield, where water droplets obscure much of the scene. The distortion caused by the droplets results in a notable reduction in sharpness and clarity. Background details, such as the car ahead and surrounding structures, appear blurry and indistinct. Uneven focus further degrades the visibility of the image.
6	In this image, an outdoor setting is depicted as seen through a wet windshield, with water droplets heavily obscuring the view. The distortion from the droplets results in a considerable loss of sharpness and clarity. Background elements, like the car ahead and nearby buildings, are blurred and hard to discern. The uneven focus of the image adds to its overall lack of visibility.
7	This image shows an outdoor scene viewed through a wet windshield, with water droplets creating significant visual obstruction. The distortion caused by the droplets reduces the sharpness and clarity of the image. Background features, including the car ahead and the surrounding structures, are blurred and difficult to distinguish. The image also suffers from uneven focus, limiting its visibility.
8	An outdoor view is displayed in this image, captured through a wet windshield with water droplets blocking much of the scene. The distortion caused by the droplets results in reduced clarity and sharpness. Background details, such as the car in front and nearby structures, are blurred and indistinct. Uneven focus further compromises the visibility of the image.
9	The image depicts an outdoor scene seen through a wet windshield, where water droplets obscure much of the view. This distortion leads to a significant loss of sharpness and clarity. Background elements like the car ahead and nearby structures are blurred and challenging to distinguish. Uneven focus further contributes to the poor visibility of the image.
10	This image represents an outdoor scene viewed through a wet windshield, with water droplets obscuring the visual field. The distortion caused by the droplets significantly affects sharpness and clarity, making background features like the car ahead and surrounding structures appear blurred and indistinct. The uneven focus further reduces the overall visibility of the image.