# InstantRetouch: Efficient and High-Fidelity Instruction-Guided Image Retouching with Bilateral Space

Jiarui Wu[1,2], Yujin Wang[1], Ruikang Li[1,2], Fan Zhang[1], Mingde Yao[2], Tianfan Xue[2,1,3]

[1]Shanghai AI Laboratory,[2]CUHK MMLab,[3]CPII under InnoHK
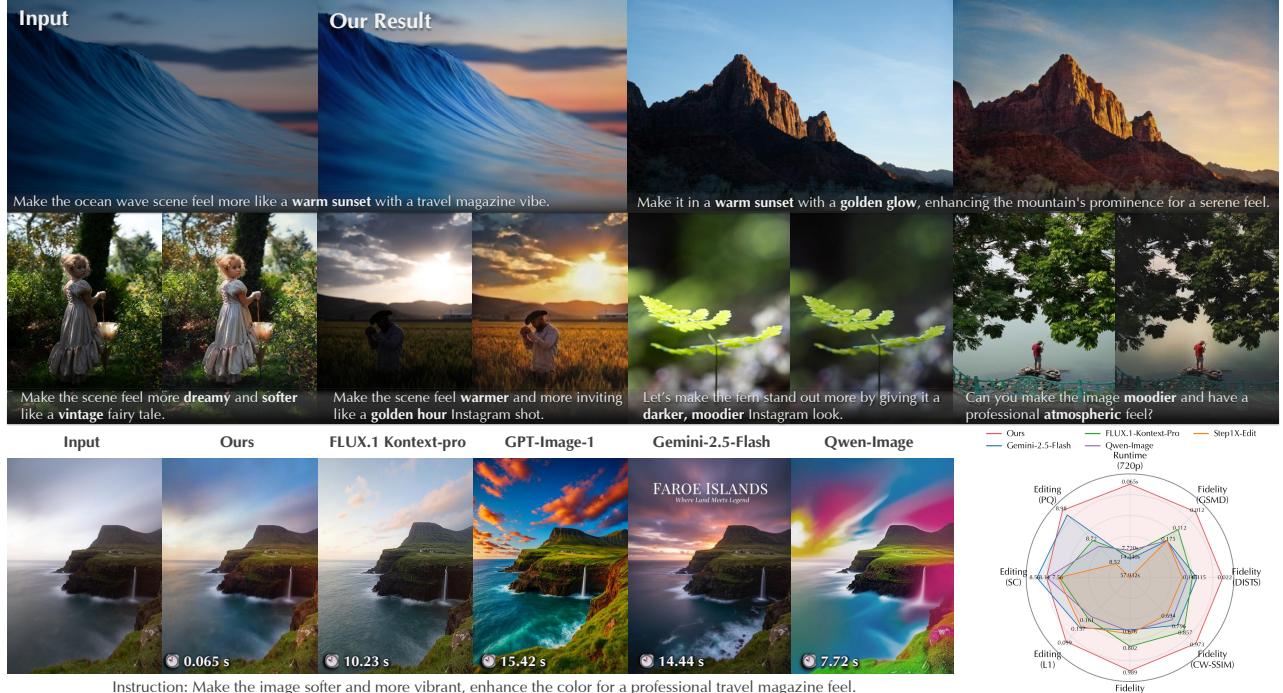
{wj024, tfxue}@ie.cuhk.edu.hk

Figure 1. Comparing our method with state-of-the-art image editing methods [6, 22, 28, 43]. As shown in the upper part, our method follows text instructions to generate visually pleasing retouching results while preserving high fidelity, whether for natural landscapes or portraits. In contrast, as demonstrated in the lower section, other state-of-the-art methods modify the original content. Finally, as depicted in the multi-dimensional comparison chart, our method outperforms others in terms of *fidelity*, *quality*, and *speed*.

## Abstract

*Language-guided photo retouching aims to adjust color and tone while preserving geometry and texture. Recently, diffusion-based retouching shows a superior visual quality, but often struggles with both fidelity issues due to its generative nature and efficiency because of its iterative sampling process. In this work, we propose an efficient and fidelity-preserving retouching method using bilateral space manipulation, which is both compact and content-decoupled. Specifically, instead of directly editing pixels or image latents, our model predicts a low-resolution bilateral grid of affine transforms, which are sliced using a learned guid-ance map and then applied to the full-resolution image. This approach yields both high fidelity and improved efficiency. To retain strong priors of a pretrained generative model, we distill a multi-step diffusion model into our bilateral grid framework using Variational Score Distillation, complemented by a prompt alignment loss to guide instruction-following behavior. Additionally, we introduce a new benchmark and evaluate our method across multiple dimensions: fidelity, instruction following, and efficiency. Compared to latest editing methods, like Gemini-2.5-Flash (Nano-Banana), our method can avoid content drift, significantly improve latency, and generate visually pleasing*

*edits, while maintaining a high level of fidelity.*

## 1. Introduction

The ability to automatically retouch photos using natural language instructions represents a significant advancement over traditional image enhancement algorithms [31, 48], which often lack expressive, fine-grained control. This paradigm shift has been driven by large-scale diffusion models [2, 34], capable of producing expressive and visually pleasing results guided by user instructions. Recent work continues to scale up these generative models for general-purpose image editing, as seen in Step-1X-Edit [28], FLUX.1-Kontext [22], Qwen-Image [43], or Gemini-2.5-Flash [6]. These models exhibit a remarkable ability in general-purpose editing, such as adding or removing objects, often producing results that are indistinguishable from real images.

Still, these generative models exhibit limitations in fidelity and efficiency when applied to image retouching. First, for photo retouching, changes must be restricted to photometric adjustments without affecting geometry or texture. Existing generative editing models, however, may not adequately disentangle these edits, leading to unwanted content drift, as shown in Fig. 1. Second, these models, often based on iterative diffusion processes, are computationally expensive and slow, limiting their application for high-resolution image retouching.

These limitations arise because generative editing directly modifies the variational latent of the input image in the diffusion process [34]. The latent representation consists of both actual image content and photometric information (brightness, color, etc.), which is unnecessarily large for retouching, slowing down the process. Manipulating in latent representation may also introduce the risk of changing actual visual content, texture, or geometric structure. Instead, retouch editing should only operate on a smaller representation that only focuses on visual appearance, without content information.

Therefore, we propose to only predict the parameters of a transformation in a compact and content-decoupled *bilateral space*, instead of directly touch original image content. The bilateral manipulation space [4, 5, 12] is instantiated as a low-resolution 3D bilateral grid of affine transforms. A learned guidance map slices this grid to produce per-pixel affine coefficients that are applied to the full-resolution image, enabling complex tonal adjustments. This representation is exceptionally efficient even at 4K resolution and ensures high fidelity by design. As shown in Fig. 1, our solution based on bilateral grids is 70-800 times faster and better preserves fidelity compared with baselines.

While the bilateral space offers an ideal representation for retouching, generating visually pleasing results from instructions still requires the rich semantic priors from dif-

fusion models. However, their slow, iterative inference is fundamentally at odds with our desired efficiency. We resolve this conflict by distilling a multi-step diffusion model into a fast, one-step generator that directly predicts the bilateral grid. In this way, we can leverage the *rich diffusion priors* for visually pleasing results guided by instruction, along with the *fidelity and efficiency* advantages of the bilateral space.

To enable this distillation, we first curate a large-scale, high-quality instruction-retouching dataset to fine-tune a multi-step teacher diffusion model. We then transfer its knowledge to an efficient student network which outputs the bilateral grid in a single forward pass, using a one-step bilateral distillation framework. In this one-step distillation, we employ Variational Score Distillation (VSD) [42, 47] as the core objective, which we augment with a CLIP-based [32] prompt alignment loss. This provides crucial semantic supervision to improve instruction following, particularly for ambiguous or stylistic prompts where pixel-level signals are weak. In addition, we design a bilateral loss to better regularize bilateral grid prediction. At last, we design a progressive distillation strategy to ensure training stability.

To evaluate performance on the instruction-guided retouching task, we introduce a new benchmark, iRetouch, composed of diverse real-world instruction-guided retouching scenarios. We assess models along three key axes: content fidelity, measured by the preservation of original texture and geometry; instruction following, evaluated via text-image alignment metrics and human preference studies; and efficiency, quantified by latency at various resolutions. As demonstrated in Fig. 1, our method is 70-800 times faster than large editing models [2, 6, 17, 22, 28, 43] and achieves superior content fidelity, all while maintaining comparable instruction-following performance.

## 2. Related Works

**Instruction-based Image editing.** Image editing enables intuitive image modifications driven by language. Early works, such as InstructPix2Pix [2], fine-tuned diffusion models by creating paired instruction-image datasets. Subsequent research [9, 13, 23, 29, 52] primarily focused on architectural optimizations to improve control granularity and consistency, while others [3, 11, 38, 50] concentrated on data-driven enhancements, expanding the range of instructions and diversifying editing examples. Additionally, some approaches [15, 16, 24, 51] integrated large language model reasoning with diffusion-based image synthesis, while others leveraged chain-of-thought (CoT) reasoning [10] to improve the model's reasoning ability for handling more complex editing tasks. Flow-edit [21] constructs an ordinary differential equation to map the source and target distributions, reducing transport costs in text-driven editing. JarvisArt [26], on the other hand, combines a multi-modal large language model (MLLM)-driven

agent that understands user intent and intelligently coordinates over 200 retouching tools. Recently, image editing has increasingly shifted toward large models with multi-modal fusion [1, 6, 7, 22, 28, 40, 49]. For example, FLUX.1 Kontext [22], as a generative flow matching model, integrates both image generation and editing tasks into a unified architecture, handling both local editing and generative in-context tasks.

**Image retouching.** Automating the complex task of image style adjustment has seen varied approaches. Early methods like 3D LUTs [31, 48] were fast but confined to fixed styles, while generative models [18] often lack sufficient interpretability and usually alter the original content of the image. More recent works utilized reinforcement learning to automate editing [14, 20, 44]. Tseng *et al.* [39] used neural networks to proxy different image processing modules and optimized the image processing pipeline parameters using a style loss function. However, those methods mentioned above typically handle a single style during training and cannot offer flexible control based on instructions.

## 3. Method Overview

Our goal is to leverage the *rich diffusion priors* for instruction-guided editing while retaining the *fidelity and efficiency* of the bilateral space. To this end, our method distills a multi-step diffusion model into a fast, one-step generator, $G_\theta$, that directly predicts a bilateral grid. The process unfolds in two main stages. First, we curate a large-scale, high-quality instruction-retouching dataset (Sec. 3.1) to fine-tune a multi-step diffusion teacher, $\epsilon_\phi$ (Sec. 3.2). Second, we distill the knowledge from this teacher into our one-step bilateral grid generator (Sec. 3.3) using a novel distillation framework (Sec. 3.4).

As illustrated in Fig. 2, our generator $G_\theta$ consists of two synergistic branches: a low-resolution one-step diffusion branch for semantic understanding and retaining rich diffusion priors, and a full-resolution bilateral processing branch that applies the learned edit to deliver high-fidelity retouching on high-resolution input. However, directly training the proposed bilateral processing network may introduce instability in training; we instead adopt a progressive training strategy. We first train the low-resolution branch by minimizing the Variational Score Distillation loss (Sec. 3.4.1), a data loss (Sec. 3.4.3), and our prompt alignment loss (Sec. 3.4.2). We then jointly optimize both branches, adding a bilateral loss (Sec. 3.4.4) to optimize the full-res bilateral branch.

### 3.1. Training Dataset

Training diffusion models relies on large-scale data. Existing instruction-editing datasets primarily focus on object-level or geometric edits and lack the fine-grained, high-fidelity examples needed for photo retouching. We therefore construct a new dataset of $\sim$ 200K triplets $(x, x^\star, c_T)$,

where $x$ is the input image, $x^\star$ is a high-quality retouched target, and $c_T$ is a textual instruction. Our dataset is built via a controlled degradation process.

**High-quality targets.** We curate visually pleasing images from public datasets and the web, filtered by no-reference image quality metrics MUSIQ [19] and LAION aesthetic score [37] with conservative thresholds, yielding targets $x^\star$.

**Input image generation.** For each target $x^\star$, we synthesize a degraded input $x$ by applying random photometric adjustments via a photo-finishing pipeline [44]. This includes perturbations to exposure, gamma, white balance, contrast, tone curves, saturation, shadows/highlights, and HSL. To simulate local retouching, we generate region masks using a Grounding-SAM procedure [27, 33] and additional soft masks from simple priors, applying different degradation parameters within each mask to induce spatially varying edits while preserving fidelity.

**Instruction generation.** Given $(x, x^\star)$, we prompt a multi-modal LLM (Qwen2.5-VL-72B [1]) in a role-playing template to generate concise, diverse photo-finishing instructions $c_T$ that describe the transformation from $x$ to $x^\star$. A small rule-based checker enforces diversity and filters content-edit verbs. Further details on dataset construction are in the Appendix.

### 3.2. Pretrained Multi-step Diffusion

Following InstructPix2Pix [2], our teacher model $\epsilon_\phi$ is a UNet trained to predict the noise added to a target image's latent representation. Let $x$ be the input image, $x^\star$ the target, and $c_T$ the text instruction. We operate in the VAE latent space of a pre-trained Stable Diffusion model [34], with encoder $\mathcal{E}_\phi$ and decoder $\mathcal{D}_\phi$. During training, noise $\epsilon$ is added to the target latent $z_0 = \mathcal{E}_\phi(x^\star)$ to create a noisy latent $z_t = \alpha_t z_0 + \beta_t \epsilon$. The teacher $\epsilon_\phi$ is trained with an MSE loss to predict this noise, conditioned on the input image latent $c_I = \mathcal{E}_\phi(x)$ and the text prompt $c_T$:

$$\mathcal{L}_{\text{teacher}}(\phi) = \mathbb{E}_{x,x^\star,c_T,t,\epsilon}\left[\left\|\epsilon - \epsilon_\phi(z_t,\ t,\ c_I,\ c_T)\right\|_2^2\right].$$
(1)

However, applying this multi-step diffusion model for retouching is slow and prone to content drift. We therefore distill from this heavy pretrained editor into a lightweight, one-step retouching model designed to guarantee content fidelity, discussed below.

### 3.3. One-step Bilateral Generator

As shown in Fig. 2, our lightweight one-step bilateral grid generator is composed of two branches: a low-resolution diffusion branch and a full-resolution bilateral processing branch. The low-resolution branch contains a frozen VAE encoder $\mathcal{E}_\theta$ and a one-step U-Net denoiser $\epsilon_\theta$, tasked with semantic understanding and preserving the rich diffusion priors. During training, a VAE decoder is temporarily employed to generate a low-resolution image, which helps to stabilize the distillation process.
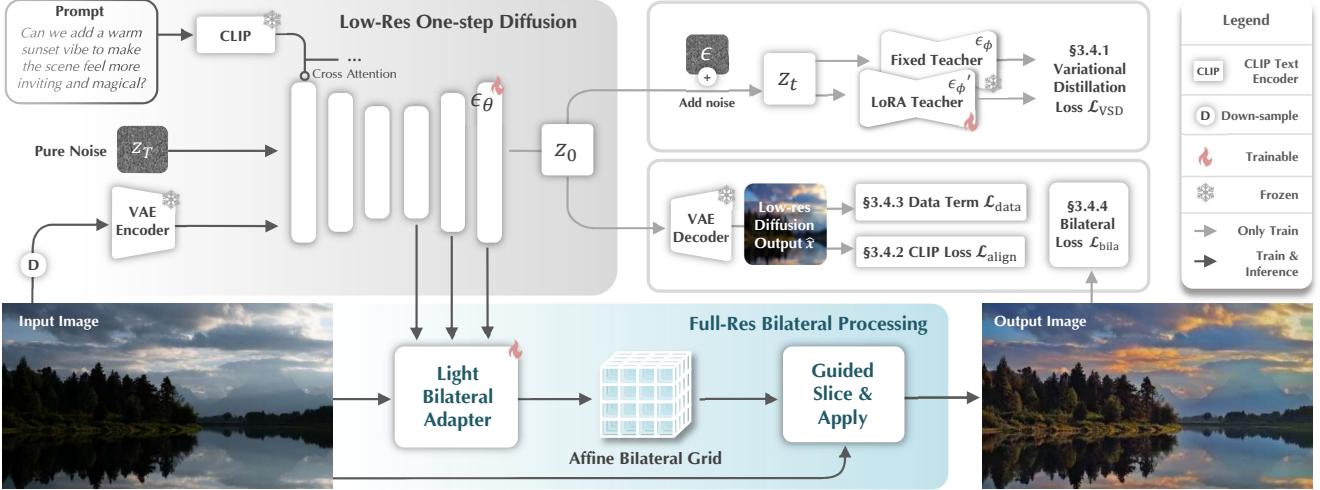
Figure 2. Our framework distills a multi-step diffusion teacher into a fast, one-step generator composed of two synergistic branches. The low-resolution diffusion branch processes the input image and text instruction to understand the edit, and then uses a light bilateral adapter to predict the parameters of a bilateral grid. The full-resolution branch then applies this grid to the original high-res image, producing the final high-fidelity result. We use Variational Score Distillation (VSD) to transfer the teacher's knowledge and a CLIP-based language alignment loss to ensure instruction alignment.

The full-resolution branch contains a lightweight bilateral adapter. In a single forward step, it generates a bilateral grid [4] $\Gamma \in \mathbb{R}^{H_g \times W_g \times D \times 12}$ that stores local affine transformation parameters in 3D space. This grid is then processed by a fully differentiable "slice-and-apply" operator that acts on the full-resolution input image of size $(H, W)$. For each input pixel with coordinates $(x', y')$ and color $(r, g, b)$, the operator first computes a grayscale guidance value $z = g(r, g, b)$ via a learned lookup table. It then uses the pixel's spatial coordinates and this guidance value to retrieve a specific affine matrix $A$ by slicing the grid with trilinear interpolation: $A = \Gamma(x'W_g/W, y'H_g/H, z/d)$. Finally, this matrix is applied to the original pixel color, $O = A \cdot (r, g, b, 1)^T$, to produce the final output. This entire mechanism delivers efficient and high-fidelity retouching directly on the high-resolution image.

Our model is highly efficient due to its design. The full-resolution operators have negligible runtime, even at 4K resolution. And the low-resolution branch has constant latency at different resolutions. This enables 4K image processing in just 68ms, vastly outperforming diffusion methods requiring over 10s for 720p inputs.

### 3.4. One-step Bilateral Distillation
Although the proposed one-step generator $G_\theta$ is super-efficient and guarantees no content drift by design, it has a very different structure compared with the pretrained teacher network $\epsilon_\phi$ (diffusion model), posing a challenge in distillation. Therefore, we proposed a novel progressive distillation strategy, described below.

#### 3.4.1. Variational Score Distillation in Latent Space
In the low-res one-step diffusion branch, the frozen VAE encoder and decoder are initialized from the weights of pre-

trained VAE $\mathcal{E}_\phi$ and $\mathcal{D}_\phi$, and the denoising network $\epsilon_\theta$ is initialized from the weights of pretrained denoiser $\epsilon_\phi$. Recall that diffusion models utilize a UNet to predict the noise $\hat{\epsilon}$ in noisy latent $z_t$, and the denoised latent can be obtained as $\hat{z}_0 = \frac{z_t - \beta_t \hat{\epsilon}}{\alpha_t}$. We directly conducting one-step denoising on the white noise $z_{t_{max}} \sim \mathcal{N}(0, I)$, conditioned on $c_I = \mathcal{E}_\theta(x)$ and $c_T$, to predict the clean latent $\hat{z}_0$ is calculated as:

$$\hat{z}_0 = \frac{z_{t_{max}} - \beta_t \epsilon_\theta(z_{t_{max}}, \ t_{max}, \ c_I, \ c_T)}{\alpha_t}, \quad (2)$$

and the corresponding low-resolution image is $\hat{x} = \mathcal{D}_\theta(\hat{z}_0)$. Note that during inference, the VAE decoder is not used, as $\hat{x}$ only helps to stabilize the distillation process, and is not needed in full-res bilateral processing. We regularize $\epsilon_\theta$ with a latent-space Variational Score Distillation (VSD) loss, $\mathcal{L}_{VSD}$, following the design in DMD [47].

VSD loss introduces a trainable regularizer $\epsilon_{\phi'}$ finetuned on the distribution of generated images $\hat{x}$ of the one-step generator $\epsilon_\theta$ to replicate its behaviour. Given the clean latent predicted by the one-step generator via Eqn. 2, we add noise to it to construct the noisy latent $\hat{z}_t = \alpha_t \hat{z}_0 + \beta_t \epsilon$. This $\hat{z}_t$ then serves as a common input to the teacher and regularizer to compute a stable gradient that steers the student towards the teacher. We adopt the latent form of VSD used in DMD [45, 47]. The gradient of the VSD loss w.r.t. $\theta \nabla_\theta \mathcal{L}_{VSD}$ is

$$\mathbb{E}_{t,\epsilon,\hat{z}_t}\left[\omega(t)\big(\epsilon_\phi(\hat{z}_t, t, c_I, c_T) - \epsilon_{\phi'}(\hat{z}_t, t, c_I, c_T)\big)\frac{\partial \hat{z}_0}{\partial \theta}\right], \quad (3)$$

To ensure the regularizer $\epsilon_{\phi'}$ remains a faithful proxy for the generator's current state, it is trained concurrently on noisy samples $\hat{z}_t$ derived from the generator's own outputs $\hat{z}_0$:

$$\mathcal{L}_{\text{diff}}(\phi') = \mathbb{E}_{t,\epsilon,c_I,c_T,\hat{z}_t} \left[ \|\epsilon_{\phi'}(\hat{z}_t, t, c_I, c_T) - \epsilon\|_2^2 \right]. \quad (4)$$

To further stabilize this process, we adopt a progressive schedule. Training begins with high noise levels ($t \in [t_{\text{hi}}, t_{max}]$) to learn coarse attributes like tone and exposure, before we gradually lower $t_{\text{hi}}$ to distill fine-grained color details.

### 3.4.2. Prompt Alignment Loss

Distilling a multi-step editor into one step often weaken the coupling between the instruction $c_T$ and often yields "plausible but misdirected" retouches under weak, aesthetic instructions. Thus, we need to add further supervision to ensure the output image follows the users' instructions.

Specifically, unlike object edits, retouching intents are mostly *directional and compositional* (e.g., warmer, dreamy, cinematic). While the VSD loss and data loss ensure feasibility, they do not guarantee that the change follows the intended semantic direction. We therefore convert user instruction $c_T$ into a small set of atomic *retouching attributes* $\mathcal{A}(c_T) = \{a\}$ using a rule-based matcher. Each attribute is an explicit edit direction tailored to photo retouching (e.g., `brightness:up`, `contrast:down`, `mood:cozy`, `temperature:warm`, `style:vintage`) and is paired with two short text prompts describing positive and negative directions. (e.g., "Bright Image" vs. "Dark Image"). This *attribute bank* turns a long, weak instruction into several stable, additive supervision signals. Let $\mathbf{e}^{\text{img}}(\cdot)$ and $\mathbf{e}^{\text{text}}(\cdot)$ be frozen CLIP image and text encoders. The cosine similarity of the two $\ell_2$-normalized image and text embeddings is used and its scalar value is denoted by $s$. For an attribute $a$ (e.g., `mood:cozy`) with prompts $(p_a^+, p_a^-)$, define $s_a^+ = \langle \mathbf{e}^{\text{img}}(\hat{x}), \mathbf{e}^{\text{text}}(p_a^+) \rangle$ and $s_a^- = \langle \mathbf{e}^{\text{img}}(\hat{x}), \mathbf{e}^{\text{text}}(p_a^-) \rangle$, where $\langle \cdot, \cdot \rangle$ denotes cosine similarity. The per-attribute InfoNCE loss [30] (viewed as a function of $a$) is

$$\ell_{\text{nce}}(a) = -\log \frac{\exp(s_a^+/\tau)}{\exp(s_a^+/\tau) + \exp(s_a^-/\tau)}. \quad (5)$$

Finally, with confidences $w_a$ from the matcher, the language alignment loss is applied to the one-step branch during distillation:

$$\mathcal{L}_{\text{align}} = \frac{1}{|\mathcal{A}(c_T)|} \sum_{a \in \mathcal{A}(c_T)} \left[ w_a \, \ell_{\text{nce}}(a) \right]. \quad (6)$$

This supervision restores directional alignment lost by step compression, and resolves ambiguity among many color transforms that could otherwise minimize the data term and VSD term yet deviate from $c_T$.

### 3.4.3. Data Supervision Loss

To stabilize distillation, we also add a data term that supervises the low-resolution output $\hat{x}$ with the ground truth target $x^\star$, following [45, 47]:

$$\mathcal{L}_{\text{data}} = \|\hat{x} - x^\star\|_2^2 + \lambda_{\text{LPIPS}}\mathcal{L}_{\text{LPIPS}}(\hat{x}, x^\star). \quad (7)$$

### 3.4.4. Bilateral Loss

The losses above focus on training the low-resolution branch. To also train the full-resolution bilateral processing branch, we introduce $\mathcal{L}_{\text{bila}}$. Let $\hat{x}_B$ be the final high-resolution output. This loss includes: (i) $\ell_1$ and LPIPS losses against the ground truth $x^\star$, (ii) a perceptual agreement term encouraging $\hat{x}_B$ to match the low-res prediction $\hat{x}$, and (iii) a laplacian regularizers on the bilateral grid $\Gamma$ for smoothness and a penalty that prevents RGB overflow:

$$\begin{aligned} \mathcal{L}_{\text{bila}} = &\; \lambda_1 \| \hat{x}_B - x^\star \|_1 + \lambda_2 \cdot \mathcal{L}_{\text{LPIPS}}(\hat{x}_B, x^\star) \\ &+ \lambda_3 \cdot \mathcal{L}_{\text{LPIPS}}(\hat{x}_B, \hat{x}) \\ &+ \lambda_4 \cdot \|\Delta^3\Gamma\|_2^2 + \lambda_5 \cdot \Psi(\hat{x}_B), \quad (8) \end{aligned}$$

where $\Delta^3$ is a 3D Laplacian regularizer to penalize differences between adjacent cells over the bilateral grid for smoothness, and $\Psi$ is a soft penalty discouraging out-of-gamut RGB.

### 3.4.5. Overall Objective and Distillation Strategy

Combining all training losses above, we finally design a novel two-stage progressive distillation strategy.

**Stage 1: Low-Resolution one-step diffusion branch training.** In this stage, we only train the low-resolution one-step diffusion branch. Note that the low-resolution branch shares the same network structure as the pretrained diffusion and thus distillation training is easier compared with the bilateral processing network. During training, we optimize the U-Net $\epsilon_\theta$ and the VSD regularizer $\epsilon_{\phi'}$. The objective combines VSD, the data term, and our prompt alignment loss:

$$\mathcal{L}_{\text{stage1}} = \mathcal{L}_{\text{data}} + \lambda_{\text{VSD}}\mathcal{L}_{\text{VSD}} + \lambda_{\text{align}}\mathcal{L}_{\text{align}}. \quad (9)$$

**Stage 2: Joint bilateral distillation.** After the first stage converges, we unfreeze the bilateral adapter and train the entire generator end-to-end. Since stage 1 already trains the relative heavy low-resolution network, finetuning the lightweight full-resolution bilateral processing is also trackable. To train the bilateral processing, a bilateral loss is added:

$$\mathcal{L}_{\text{stage2}} = \mathcal{L}_{\text{stage1}} + \lambda_{\text{bila}}\mathcal{L}_{\text{bila}}. \quad (10)$$

This complete *one-step bilateral distillation* framework yields an efficient model that guarantees high-fidelity, content-preserving retouching while retaining strong instruction-following capabilities.

## 4. Experiments

### 4.1. Experiment Setup

**New benchmark iRetouch.** For evaluation, we have created a new benchmark, iRetouch, consisting of 500 real-world before-and-after retouching pairs from the Adobe Lightroom community. Instructions for these pairs are generated using our method from Sec. 3.1, followed by manual refinement for clarity and diversity. The benchmark spans a wide variety of scenes (e.g., portraits, landscapes) and includes a rich vocabulary of retouching edits, such as global

Table 1. Comparison on iRetouch benchmark. Our method achieves state-of-the-art efficiency and content fidelity while remaining highly competitive in editing quality. Blank entries indicate models that cannot process high resolutions or are not instruction-driven.

| Method | Runtime(s) | | | | Content Fidelity | | | | Editing Quality | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 720p↓ | 1K↓ | 2K↓ | 4K↓ | SSIM↑ | CW-SSIM↑ | GSMD↓ | DISTS↓ | L1↓ | L2↓ | SC↑ | PQ↑ | O↑ |
| 3DLUT [48] | 0.066 | 0.079 | 0.112 | 0.201 | 0.982 | **0.981** | 0.013 | 0.024 | 0.136 | 0.034 | - | - | - |
| RSFNet [31] | **0.029** | **0.047** | 0.086 | 0.189 | 0.975 | 0.976 | **0.012** | 0.038 | 0.137 | 0.034 | - | - | - |
| InstructPix2Pix [2] | 4.632 | - | - | - | 0.742 | 0.768 | 0.149 | 0.177 | 0.164 | 0.050 | 7.11 | 7.58 | 7.34 |
| Step1X-Edit [28] | 57.932 | - | - | - | 0.706 | 0.694 | 0.174 | 0.167 | 0.140 | 0.036 | 7.63 | 8.52 | 8.06 |
| GPT-Image-1 [17] | 15.427 | 21.889 | - | - | 0.505 | 0.397 | 0.242 | 0.216 | 0.215 | 0.082 | 8.09 | 8.56 | 8.32 |
| Qwen-Image [43] | 7.720 | - | - | - | 0.689 | 0.744 | 0.174 | 0.147 | 0.168 | 0.054 | 8.12 | 8.67 | 8.39 |
| FLUX.1-Kontext-Pro [22] | 10.235 | - | - | - | 0.802 | 0.857 | 0.112 | 0.132 | 0.161 | 0.050 | 7.56 | 8.72 | 8.12 |
| Gemini-2.5-Flash [6] | 14.440 | - | - | - | 0.676 | 0.796 | 0.175 | 0.115 | 0.137 | 0.036 | **8.56** | 8.94 | **8.74** |
| **Ours** | 0.065 | 0.065 | **0.066** | **0.068** | **0.989** | 0.973 | **0.012** | 0.022 | **0.099** | **0.018** | 8.14 | **8.98** | 8.54 |

adjustments, specific styles (cinematic, dreamy), moods, and local effects (see Appendix for a detailed breakdown).

**Content fidelity metrics.** Retouching is non-destructive, so edits must preserve structure and texture without repaints. To factor out intentional tone changes, we convert outputs to grayscale and histogram-match them to the input, then compute SSIM [41] (structural similarity), CW-SSIM [35] (geometry and texture distortion), DISTS [8] (textural similarity), and GMSD [46] (gradient-magnitude consistency).

**Editing quality metrics.** Following prior work [28, 50], we report L1/L2 distances, instruction–image alignment (SC, 0–10), perceptual quality (PQ, 0–10), and the overall score $O = \sqrt{SC \times PQ}$. SC and PQ are generated using GPT-4o, similar to [28]. Additional details are provided in the Appendix.

**Implementation.** Our one-step bilateral generator is built upon a pre-trained Stable Diffusion editor. We initialize our student U-Net from the teacher's weights and freeze the VAE. VSD distillation follows a three-stage curriculum over timesteps to learn from coarse structure and tone (high $t$), then instruction alignment (mid $t$), and finally fine-grained color details (low $t$). Training is at 512px using AdamW with EMA, mixed precision, and gradient clipping. Inference is a single pass: the model generates a bilateral grid and applies it to the native resolution input, yielding constant-time performance regardless of image size. We train on our instruction-retouching dataset in Sec. 3.1. See the Appendix for full implementation details.

**Runtime.** We measure end-to-end latency across resolutions from 720p to 4K. Open-source models are benchmarked on a server with 8 NVIDIA RTX 4090 GPUs. For proprietary models, we report the full end-to-end API latency, including data transfer.

## 4.2. Evaluation and Results

We compare our method with baselines across three categories: (1) traditional enhancement methods [31, 48], (2) open-source image editing models [2, 28, 43], and (3) proprietary large-scale editing models [6, 17, 22].

### 4.2.1. Evaluation on Our iRetouch Benchmark

As shown in Tab. 1, our method outperforms others in terms of runtime, fidelity, and editing quality.

**Efficiency.** Our model demonstrates exceptional efficiency, maintaining a near-constant inference time of 0.065–0.068s from 720p up to 4K resolutions. This represents a 70–900× speedup over generative baselines at 720p. The blank runtime entries for some baselines highlight a critical limitation: most diffusion-based models cannot natively process resolutions beyond 1K, a barrier our design overcomes.

**Fidelity.** Our approach achieves state-of-the-art content fidelity among all instruction-guided models. This confirms our bilateral branch successfully prevents the textural distortions common in pure diffusion editors.

**Quality.** For editing quality, our model's overall score (O) of 8.54 is highly competitive with the top proprietary system (Gemini-2.5-Flash at 8.74) and significantly surpasses other open-source editors. The blank quality scores for traditional methods like 3D LUT exist because they are not instruction-driven and thus cannot be evaluated for semantic alignment. In summary, our method delivers near–state-of-the-art editing quality with state-of-the-art fidelity and 4K-constant runtime. These results support our design goal: instruction-guided retouching that is high-fidelity, fast, and stable across resolutions. We also provide a more detailed analysis of the relationship between quality and fidelity in the Appendix.

**Visual comparison.** Fig. 3 provides a qualitative comparison across a range of instructions. The results highlight a common failure in competing methods: a trade-off between editing quality and content fidelity. Generative editors like InstructPix2Pix and GPT-Image-1 often introduce severe artifacts, hallucinations, or unwanted text overlays, fundamentally altering the source image. Even capable models like Gemini-2.5-Flash can subtly change key features. Our method, however, successfully follows both global and local instructions while maintaining high fidelity, applying the desired stylistic edits without distorting content or compromising the original photograph's integrity.

Make the scene more dramatic and moody, emphasizing the pier's colors and the water's movement for a cinematic feel.

Can we make it feel like a warm sunset with a cozy travel vibe?

Make the image moodier and more atmospheric, brightening the person's face and highlighting the vintage airplane backdrop.

Make the image brighter and more vibrant, let the elephant pops more.

Can we make the background warmer and more autumnal for a cozy fall Instagram look?

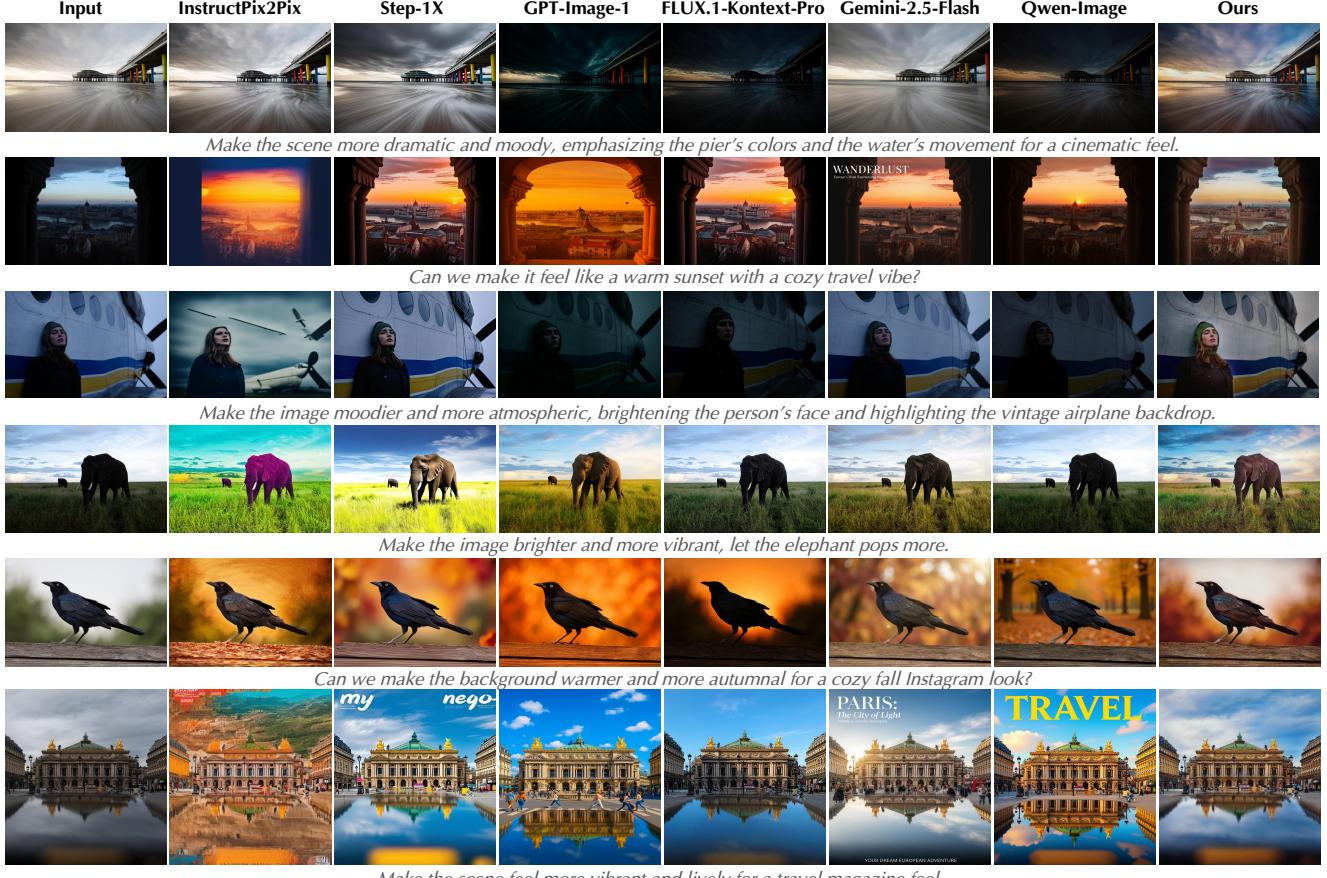Make the scene feel more vibrant and lively for a travel magazine feel.

Figure 3. Visual comparisons of different image editing methods on our iRetouch benchmark.
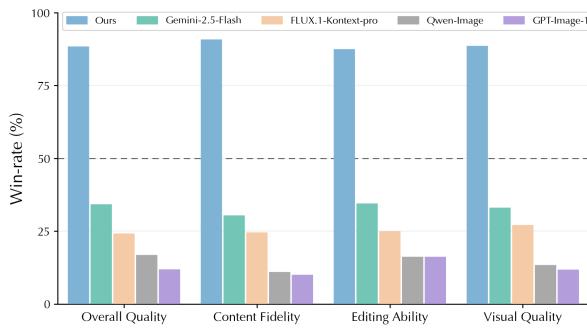
### 4.2.2. User Study



Figure 4. User preference study results on iRetouch benchmark.

To assess subjective user preference, we conducted a user study with 30 participants, who evaluated 20 retouching examples from our iRetouch benchmark. They compared our method against four leading baselines (FLUX.1-Kontext-pro, Gemini-2.5-Flash, Qwen-Image, GPT-Image-1) across four dimensions: content fidelity, editing ability, visual quality, and overall preference. As shown in Fig. 4, the results reveal a clear and consistent preference for our method. Our approach achieved the highest ratings in all categories, confirming that users favor its artifact-free, high-fidelity edits that accurately reflect their intent.

### 4.2.3. Evaluation of Identity Preservation on PPR10K



Make the image feel more intimate and balanced, bringing out the soft, cinematic glow.
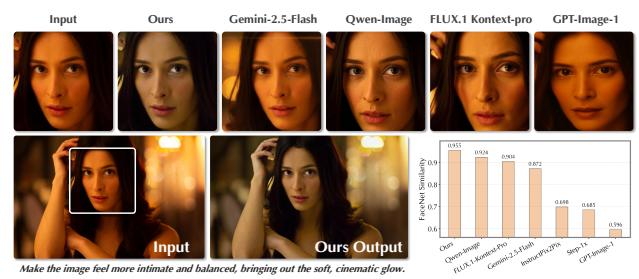
Figure 5. Results of identity preservation comparison on the PPR10K dataset. Our model scores highest in facial similarity and avoids the identity-altering artifacts.

In portrait editing, content fidelity is crucial as it requires strict identity preservation. To evaluate identity preservation on this task, we test on 100 images from the PPR10K dataset [25] with MLLM generated instructions. We quantify identity preservation by extracting facial embeddings from the input and output images using FaceNet [36] and then computing their cosine similarity. As shown quantitatively in Fig. 5, our method achieves the highest face similarity score. We also include qualitative comparison in the figure; our model retouches the portrait while strictly preserving fidelity, whereas competing methods introduce no-

Figure 6. Visualization of ablation study on the loss configuration of one-step bilateral distillation.

ticeable repainting that distorts the subject's identity.

## 4.3. Ablation

We conduct a series of ablation studies to validate our key design choices, focusing on our framework and the components of our distillation algorithms.

Table 2. Ablation study on our framework. We evaluate content fidelity, editing quality, and runtime. Our full model effectively combines the strengths of diffusion priors and bilateral processing, achieving high scores across all criteria.

| Method | Runtime(s)↓ | Content Fidelity | | | Editing Quality | | |
|---|---|---|---|---|---|---|---|
| | | SSIM↑ | GSMD↓ | DISTS↓ | SC↑ | PQ↑ | O↑ |
| Bilateral Grid Prediction | 0.001 | 0.996 | 0.003 | 0.005 | 4.48 | 8.28 | 6.09 |
| Teacher (Multi-step Diffusion) | 4.602 | 0.833 | 0.095 | 0.121 | 7.96 | 8.71 | 8.33 |
| Hybrid (Teacher Features + Bilateral) | 0.065 | 0.904 | 0.073 | 0.107 | 5.65 | 7.62 | 6.56 |
| Student (Diffusion-Only) | 0.319 | 0.788 | 0.130 | 0.152 | 8.43 | 8.85 | 8.64 |
| **Ours (Full Model)** | **0.065** | **0.989** | **0.012** | **0.022** | **8.14** | **8.98** | **8.54** |

**Ablation on framework.** We first analyze the contribution of our framework components in Tab. 2. We compare our full model against four key baselines: (1) *Bilateral Grid Prediction*, a model that directly predicts a bilateral grid from the input image without diffusion priors, trained on our dataset; (2) our *Teacher (Multi-step Diffusion)* model; (3) a *Hybrid* model that uses features from the multi-step teacher to predict a bilateral grid; and (4) our *Student (Diffusion-Only)*, which corresponds to the low-resolution RGB output from our distilled U-Net without the bilateral branch.

The results in Tab. 2 reveal a clear trade-off. Purely diffusion-based models (Teacher, Student-Only) achieve high editing quality but low fidelity. In contrast, a simple Bilateral Grid Prediction model preserves content perfectly (0.996 SSIM) but fails to follow instructions (6.09 O-score). Our full model uniquely resolves this conflict by merging the semantic strength of diffusion (8.54 O-score) with the structural preservation of bilateral processing (0.989 SSIM), all while maintaining high efficiency. This validates our dual-branch design for balancing fidelity, quality, and speed.

**Ablation on one-step bilateral distillation.** Next, we validate the effectiveness of the core loss components in our one-step bilateral distillation process. As shown in Tab. 3, we start with a base objective, $\mathcal{L}_{\text{base}}$, which includes only the data term and bilateral losses ($\mathcal{L}_{\text{data}} + \mathcal{L}_{\text{bila}}$). We then

Table 3. Ablation on our distillation loss components. Both VSD and our prompt alignment loss ($\mathcal{L}_{\text{align}}$) are critical for achieving high editing quality.

| Loss Configuration | Editing Quality | | |
|---|---|---|---|
| | SC↑ | PQ↑ | O↑ |
| $\mathcal{L}_{\text{base}}$ | 5.978 | 8.280 | 7.036 |
| $\mathcal{L}_{\text{base}} + \mathcal{L}_{\text{VSD}}$ | 7.257 | 9.013 | 8.087 |
| $\mathcal{L}_{\text{base}} + \mathcal{L}_{\text{VSD}} + \mathcal{L}_{\text{align}}$ | 8.14 | 8.984 | 8.553 |

progressively add our main distillation loss, $\mathcal{L}_{\text{VSD}}$, and our prompt alignment loss, $\mathcal{L}_{\text{align}}$.

As shown in Tab. 3, the base model ($\mathcal{L}_{\text{base}}$) alone is insufficient for quality editing. Adding $\mathcal{L}_{\text{VSD}}$ is critical, dramatically boosting the score by transferring the teacher's generative priors. Incorporating our prompt alignment loss ($\mathcal{L}_{\text{align}}$) provides a final, significant gain. This confirms its role in providing essential directional supervision for interpreting stylistic prompts where VSD alone falls short. We also visualize this ablation in Fig. 6.

## 4.4. Fine-grained Control



Figure 7. Visualization of continuous control on editing strength.

Our framework's control extends beyond language prompts to include more fine-grained control over the retouching effect. As shown in Fig. 7, users can continuously adjust the retouching intensity by applying a scalar $s$ to the per-pixel affine transforms. Thanks to the linearity of the affine transform in bilateral space, setting $s = 0$ yields the input, while $s > 1$ enhances the effect. This transforms our model into a smart, language-guided filter, offering precise control where language can be ambiguous. Moreover, we support fine-grained regional control using a soft bilateral blending strategy, which is further detailed in the Appendix.

## 5. Conclusion

In this work, we present an efficient and fidelity-preserving approach to image retouching that addresses both fidelity degradation and computational inefficiency. Instead of manipulating pixels or latent features, our method operates in a compact, content-decoupled bilateral space, enabling high fidelity with significantly improved efficiency. To preserve strong generative priors, we distill a multi-step diffusion model into our bilateral grid framework via variational score distillation, enhanced with a CLIP-based contrastive loss for instruction following. We further introduce a new benchmark dataset for instruction-guided

retouching and evaluate fidelity, instruction alignment, and efficiency. Compared to recent image editing methods such as Gemini-2.5-Flash (Nano Banana), our approach runs orders of magnitude faster while achieving superior content fidelity and comparable instruction-following performance.

# References

[1] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 3

[2] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023. 2, 3, 6

[3] Tuhin Chakrabarty, Kanishk Singh, Arkadiy Saakyan, and Smaranda Muresan. Learning to follow object-centric image editing instructions faithfully. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9630–9646, Singapore, 2023. Association for Computational Linguistics. 2

[4] Jiawen Chen, Sylvain Paris, and Frédo Durand. Real-time edge-aware image processing with the bilateral grid. *ACM Transactions on Graphics (TOG)*, 26(3):103–es, 2007. 2, 4

[5] Jiawen Chen, Andrew Adams, Neal Wadhwa, and Samuel W Hasinoff. Bilateral guided upsampling. *ACM Transactions on Graphics (TOG)*, 35(6):1–8, 2016. 2

[6] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025. 1, 2, 3, 6

[7] Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao Yu, Xiaonan Nie, Ziang Song, et al. Emerging properties in unified multimodal pretraining. *arXiv preprint arXiv:2505.14683*, 2025. 3

[8] Keyan Ding, Kede Ma, Shiqi Wang, and Eero P Simoncelli. Image quality assessment: Unifying structure and texture similarity. *IEEE transactions on pattern analysis and machine intelligence*, 44(5):2567–2581, 2020. 6

[9] Zheng-Peng Duan, Jiawei Zhang, Zheng Lin, Xin Jin, Xun-Dong Wang, Dongqing Zou, Chun-Le Guo, and Chongyi Li. Diffretouch: Using diffusion to retouch on the shoulder of experts. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2825–2833, 2025. 2

[10] Tsu-Jui Fu, Wenze Hu, Xianzhi Du, William Yang Wang, Yinfei Yang, and Zhe Gan. Guiding instruction-based image editing via multimodal large language models. *arXiv preprint arXiv:2309.17102*, 2023. 2

[11] Zigang Geng, Binxin Yang, Tiankai Hang, Chen Li, Shuyang Gu, Ting Zhang, Jianmin Bao, Zheng Zhang, Houqiang Li, Han Hu, et al. Instructdiffusion: A generalist modeling interface for vision tasks. In *Proceedings of the IEEE/CVF Con-*

*ference on computer vision and pattern recognition*, pages 12709–12720, 2024. 2

[12] Michaël Gharbi, Jiawen Chen, Jonathan T Barron, Samuel W Hasinoff, and Frédo Durand. Deep bilateral learning for real-time image enhancement. *ACM Transactions on Graphics (TOG)*, 36(4):1–12, 2017. 2

[13] Qin Guo and Tianwei Lin. Focus on your instruction: Fine-grained and multi-instruction image editing by attention modulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6986–6996, 2024. 2

[14] Yuanming Hu, Hao He, Chenxi Xu, Baoyuan Wang, and Stephen Lin. Exposure: A white-box photo post-processing framework. *ACM Transactions on Graphics (TOG)*, 37(2): 1–17, 2018. 3

[15] Yujia Hu, Songhua Liu, Zhenxiong Tan, Xingyi Yang, and Xinchao Wang. Image editing as programs with diffusion models. *arXiv preprint arXiv:2506.04158*, 2025. 2

[16] Yuzhou Huang, Liangbin Xie, Xintao Wang, Ziyang Yuan, Xiaodong Cun, Yixiao Ge, Jiantao Zhou, Chao Dong, Rui Huang, Ruimao Zhang, et al. Smartedit: Exploring complex instruction-based image editing with multimodal large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8362–8371, 2024. 2

[17] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 2, 6

[18] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 3

[19] 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. 3

[20] Satoshi Kosugi and Toshihiko Yamasaki. Unpaired image enhancement featuring reinforcement-learning-controlled image editing software. In *Proceedings of the AAAI conference on artificial intelligence*, pages 11296–11303, 2020. 3

[21] Vladimir Kulikov, Matan Kleiner, Inbar Huberman-Spiegelglas, and Tomer Michaeli. Flowedit: Inversion-free text-based editing using pre-trained flow models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19721–19730, 2025. 2

[22] Black Forest Labs, Stephen Batifol, Andreas Blattmann, Frederic Boesel, Saksham Consul, Cyril Diagne, Tim Dockhorn, Jack English, Zion English, Patrick Esser, et al. Flux. 1 kontext: Flow matching for in-context image generation and editing in latent space. *arXiv preprint arXiv:2506.15742*, 2025. 1, 2, 3, 6

[23] Sijia Li, Chen Chen, and Haonan Lu. Moecontroller: Instruction-based arbitrary image manipulation with mixture-of-expert controllers. *arXiv preprint arXiv:2309.04372*, 2023. 2

[24] Shufan Li, Harkanwar Singh, and Aditya Grover. Instruc-tany2pix: Flexible visual editing via multimodal instruction following. *arXiv preprint arXiv:2312.06738*, 2023. 2

[25] Jie Liang, Hui Zeng, Miaomiao Cui, Xuansong Xie, and Lei Zhang. Ppr10k: A large-scale portrait photo retouch-ing dataset with human-region mask and group-level consis-tency. In *Proceedings of the IEEE/CVF Conference on Com-puter Vision and Pattern Recognition*, pages 653–661, 2021. 7

[26] Yunlong Lin, Zixu Lin, Kunjie Lin, Jinbin Bai, Panwang Pan, Chenxin Li, Haoyu Chen, Zhongdao Wang, Xinghao Ding, Wenbo Li, et al. Jarvisart: Liberating human artistic creativ-ity via an intelligent photo retouching agent. *arXiv preprint arXiv:2506.17612*, 2025. 2

[27] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *Euro-pean conference on computer vision*, pages 38–55. Springer, 2024. 3

[28] Shiyu Liu, Yucheng Han, Peng Xing, Fukun Yin, Rui Wang, Wei Cheng, Jiaqi Liao, Yingming Wang, Honghao Fu, Chun-rui Han, et al. Step1x-edit: A practical framework for general image editing. *arXiv preprint arXiv:2504.17761*, 2025. 1, 2, 3, 6

[29] Chaojie Mao, Jingfeng Zhang, Yulin Pan, Zeyinzi Jiang, Zhen Han, Yu Liu, and Jingren Zhou. Ace++: Instruction-based image creation and editing via context-aware content filling. *arXiv preprint arXiv:2501.02487*, 2025. 2

[30] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Repre-sentation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 5

[31] Wenqi Ouyang, Yi Dong, Xiaoyang Kang, Peiran Ren, Xin Xu, and Xuansong Xie. Rsfnet: A white-box image retouch-ing approach using region-specific color filters. In *Proceed-ings of the IEEE/CVF International Conference on Com-puter Vision*, pages 12160–12169, 2023. 2, 3, 6

[32] 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 supervi-sion. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2

[33] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 3

[34] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2, 3

[35] Mehul P Sampat, Zhou Wang, Shalini Gupta, Alan Conrad Bovik, and Mia K Markey. Complex wavelet structural sim-ilarity: A new image similarity index. *IEEE transactions on image processing*, 18(11):2385–2401, 2009. 6

[36] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clus-tering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015. 7

[37] Christoph Schuhmann and Romain Beaumont. Laion-aesthetics. *LAION. AI*, 2022. 3

[38] Shelly Sheynin, Adam Polyak, Uriel Singer, Yuval Kirstain, Amit Zohar, Oron Ashual, Devi Parikh, and Yaniv Taigman. Emu edit: Precise image editing via recognition and gen-eration tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8871–8879, 2024. 2

[39] Ethan Tseng, Yuxuan Zhang, Lars Jebe, Xuaner Zhang, Zhi-hao Xia, Yifei Fan, Felix Heide, and Jiawen Chen. Neu-ral photo-finishing. *ACM Transactions on Graphics*, 41(6): 3555526, 2022. 3

[40] Peng Wang, Yichun Shi, Xiaochen Lian, Zhonghua Zhai, Xin Xia, Xuefeng Xiao, Weilin Huang, and Jianchao Yang. Seededit 3.0: Fast and high-quality generative image editing. *arXiv preprint arXiv:2506.05083*, 2025. 3

[41] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Si-moncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 6

[42] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distilla-tion. *Advances in neural information processing systems*, 36: 8406–8441, 2023. 2

[43] Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng-ming Yin, Shuai Bai, Xiao Xu, Yilei Chen, et al. Qwen-image technical report. *arXiv preprint arXiv:2508.02324*, 2025. 1, 2, 6

[44] Jiarui Wu, Yujin Wang, Lingen Li, Fan Zhang, and Tianfan Xue. Goal conditioned reinforcement learning for photo fin-ishing tuning. In *Advances in Neural Information Processing Systems*, pages 46294–46318. Curran Associates, Inc., 2024. 3

[45] Rongyuan Wu, Lingchen Sun, Zhiyuan Ma, and Lei Zhang. One-step effective diffusion network for real-world image super-resolution. *Advances in Neural Information Process-ing Systems*, 37:92529–92553, 2024. 4, 5

[46] Wufeng Xue, Lei Zhang, Xuanqin Mou, and Alan C Bovik. Gradient magnitude similarity deviation: A highly efficient perceptual image quality index. *IEEE transactions on image processing*, 23(2):684–695, 2013. 6

[47] Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shecht-man, Fredo Durand, William T Freeman, and Taesung Park. One-step diffusion with distribution matching distillation. In *Proceedings of the IEEE/CVF conference on computer vi-sion and pattern recognition*, pages 6613–6623, 2024. 2, 4, 5

[48] Hui Zeng, Jianrui Cai, Lida Li, Zisheng Cao, and Lei Zhang. Learning image-adaptive 3d lookup tables for high perfor-mance photo enhancement in real-time. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(4):2058–2073, 2020. 2, 3, 6

[49] Hong Zhang, Zhongjie Duan, Xingjun Wang, Yuze Zhao, Weiyi Lu, Zhipeng Di, Yixuan Xu, Yingda Chen, and Yu Zhang. Nexus-gen: A unified model for image understanding, generation, and editing. *arXiv preprint arXiv:2504.21356*, 2025. 3

[50] Kai Zhang, Lingbo Mo, Wenhu Chen, Huan Sun, and Yu Su. Magicbrush: A manually annotated dataset for instruction-guided image editing. *Advances in Neural Information Processing Systems*, 36:31428–31449, 2023. 2, 6

[51] Zechuan Zhang, Ji Xie, Yu Lu, Zongxin Yang, and Yi Yang. In-context edit: Enabling instructional image editing with in-context generation in large scale diffusion transformer. *arXiv preprint arXiv:2504.20690*, 2025. 2

[52] Haozhe Zhao, Xiaojian Shawn Ma, Liang Chen, Shuzheng Si, Rujie Wu, Kaikai An, Peiyu Yu, Minjia Zhang, Qing Li, and Baobao Chang. Ultraedit: Instruction-based fine-grained image editing at scale. *Advances in Neural Information Processing Systems*, 37:3058–3093, 2024. 2