# Modular Neural Image Signal Processing

Mahmoud Afifi     Zhongling Wang     Ran Zhang     Michael S. Brown

AI Center-Toronto, Samsung Electronics

{m.afifi1, z.wang2, ran.zhang, michael.b1}@samsung.com

| Input raw image | Intermediate stages of our modular ISP | Our final rendered image | Re-rendering w/ diff. picture styles & edits |

Figure 1. We present a modular neural image signal processing (ISP) framework that offers full control over every stage of the pipeline and can handle unseen cameras without requiring re-training. On top of this framework, we built a user-interactive tool that supports post-editable re-rendering, allowing users to re-process saved outputs with different picture styles and manual adjustments. The shown image was captured in raw format using the iPhone 13 main camera, then denoised and processed by our modular ISP, with intermediate stages and multiple picture-style and manual-adjustment results displayed. None of our models were trained on data from iPhone cameras.

## Abstract

*This paper presents a modular neural image signal processing (ISP) framework that processes raw inputs and renders high-quality display-referred images. Unlike prior neural ISP designs, our method introduces a high degree of modularity, providing full control over multiple intermediate stages of the rendering process. This modular design not only achieves high rendering accuracy but also improves scalability, debuggability, generalization to unseen cameras, and flexibility to match different user-preference styles. To demonstrate the advantages of this design, we built a user-interactive photo-editing tool that leverages our neural ISP to support diverse editing operations and picture styles. The tool is carefully engineered to take advantage of the high-quality rendering of our neural ISP and to enable unlimited post-editable re-rendering. Our method is a fully learning-based framework with variants of different capacities, all of moderate size (ranging from ∼0.5 M to ∼3.9 M parameters for the entire pipeline), and consistently delivers competitive qualitative and quantitative results across multiple test sets.*

## 1. Introduction and Related Work

Image signal processing (ISP) is a set of computational operations, typically organized as a sequential pipeline, that transform linear raw sensor data into display-referred, high-quality images [34]. These operations include raw image enhancement [23, 32, 57, 59, 79], white balancing and color-space conversion [14, 18, 22, 24, 46], and various quality enhancement [2, 39, 51, 77, 78, 83, 90], each targeting a specific stage of the pipeline and collectively often requiring careful calibration and engineering to function as a unified ISP system.

Recent learning-based methods model the entire ISP pipeline as a single black-box neural network trained end-to-end to map raw images from a specific camera to their display-referred outputs (e.g., [44, 48, 49, 69, 70, 85, 87, 91]). However, such monolithic designs tend to generalize poorly to unseen cameras, as their learned mappings are tightly coupled to the characteristics of the training camera [6, 68]. Moreover, many of these networks (e.g., [44, 48, 70]) have high memory and computational requirements, limiting their practicality for deployment scenarios such as on-device processing or real-time photo ren-

dering. Beyond computational cost, these black-box ISPs are difficult to interpret, debug, or extend in real-world deployments, where continuous improvement and scalability are critical (e.g., supporting new picture styles or handling user-specific corner cases [3, 4, 17, 27, 41]).

A few recent attempts have explored multi-stage or modular ISP designs (e.g., [10, 33, 53, 58, 60, 76]); however, these approaches remain limited in several ways. Some adopt overly generic stage definitions (e.g., restoration vs. enhancement [58], or local vs. global stages [10, 53]), while others require post-training fine-tuning—providing no assurance that the stages retain their intended functionality and thereby reducing interpretability [60]. Another line of work relies on a commercial raw-processing pipeline to generate training data for each stage [76]. However, the pipeline is closed-source, and intermediate-stage outputs are difficult to access, with limited information available about its internal processing order. Others operate on non-linear display-referred images, without a clear scheme for adapting them to the linear raw domain [33].

This research gap motivates us to develop a learning-based modular ISP framework that achieves high-quality results while maintaining a high degree of modularity and interpretability to facilitate scalability, debugging, and flexibility. Unlike conventional neural ISPs, our method explicitly decomposes the pipeline into well-defined, learnable stages, each responsible for a specific part of the image rendering process. This structure not only preserves functional transparency but also allows individual modules to be improved, replaced, or reused across cameras and picture styles without retraining the entire system. With this modular design, we achieve competitive image quality while gaining the ability to analyze and isolate the causes of corner cases, replace camera-specific modules with generic ones for unseen cameras, and extend the framework to additional picture styles without duplicating the full pipeline. Moreover, the modularity provides a foundation for interactive and user-controllable image rendering. To showcase this flexibility, we integrate several image-editing operators directly into the learnable pipeline, allowing users to interactively adjust the final image appearance (see Fig. 1).

**Contributions.** We introduce a fine-grained modular ISP framework that provides explicit control over the entire raw rendering process through well-defined, interpretable components, while supporting diverse picture styles with lightweight to moderate model capacity. Our framework achieves state-of-the-art results across multiple picture styles and offers full user control over the rendering pipeline for further customization. To demonstrate its practicality, we develop an interactive photo-editing tool built upon our modular ISP, enabling users to process raw images from unseen cameras, select or interpolate between picture styles, and apply photo-editing adjustments—all within a fast and lightweight system. Furthermore, we describe how the tool supports unlimited re-rendering of saved images by embedding raw data within output JPEGs, and how it extends to editing standard sRGB images produced by unknown third-party cameras or software.

## 2. Method

Given a demosaiced raw image $\mathbf{I}_{\mathtt{raw}} \in \mathbb{R}^{H \times W \times 3}$, our objective is to render a high-quality, display-referred version $\mathbf{I}_{\mathtt{out}} \in \mathbb{R}^{H \times W \times 3}$ in a standard color space (assumed to be sRGB in this paper) that closely matches a ground-truth reference $\mathbf{I}_{\mathtt{GT}}$ rendered from the same raw image under a specific picture style. Beyond high-quality rendering, our goal is to maintain a high degree of modularity across the rendering process, with interpretable stages that enhance scalability, facilitate debugging, and provide fine-grained control over the entire pipeline. To this end, we propose a learning-based modular ISP framework, illustrated in Fig. 2. The framework consists of a raw enhancement stage (Sec. 2.1), followed by a color correction stage (Sec. 2.2) that outputs $\mathbf{I}_{\mathtt{LsRGB}}$ in linear sRGB color space. This image is then downsampled to $\mathbf{I}_{\mathtt{LsRGB}}^{\downarrow} \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times 3}$ and processed by a modular photofinishing module (Sec. 2.3) to render the image, after which a guided upsampling step (Sec. 2.4) produces the photofinished result at the original input resolution. Lastly, this image is refined through a detail-enhancement process (Sec. 2.5) to generate the final output $\mathbf{I}_{\mathtt{out}}$. Each learnable stage in the pipeline is trained independently to preserve stage-level modularity, allowing individual stages to be replaced or updated without retraining the entire framework. The remainder of this section describes our framework pipeline, whereas network architecture details are provided in Sec. C of the supplementary material.

### 2.1. Raw Enhancement

Raw images often exhibit noticeable noise and detail loss, particularly in dark or underexposed regions, due to low photon counts and sensor imperfections [28]. Among various raw degradations, we focus on denoising, as it is critical for preserving fine details in subsequent ISP stages [1]. While multi-frame burst denoising [35, 40, 42, 66] achieves high quality, it requires multiple captures of the same scene. To maintain broad applicability (supporting both on-device and software-based rendering), our framework adopts single-image raw denoising:

$$\mathbf{I}_{\mathtt{enh-raw}} = f_{\mathtt{enh-raw}}(\mathbf{I}_{\mathtt{raw}}), \qquad (1)$$

where $f_{\mathtt{enh-raw}}(\cdot)$ denotes our denoising function, implemented as a fully convolutional network $\mathcal{D}_{\mathtt{raw}}$.

We train $\mathcal{D}_{\mathtt{raw}}$ in a supervised manner using a pixel-wise $\ell_1$ loss between the predicted and ground-truth raw images.
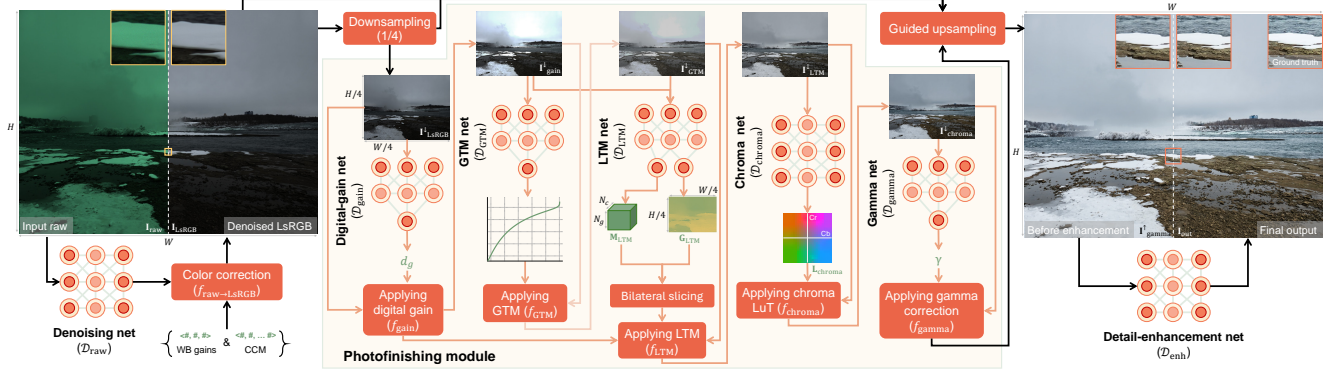
Figure 2. Overview of our modular framework. The pipeline begins with image denoising, followed by color correction to map the denoised raw image to the linear sRGB space. The photofinishing module then processes a downsampled version of the linear sRGB image through five parametric stages, where neural networks predict image-based parameters for each stage: digital gain map, global tone mapping (GTM), local tone mapping (LTM), chroma mapping, and gamma correction. A guided upsampling step, using the full-resolution linear sRGB image as guidance, reconstructs the full-resolution photofinishing output, which is then refined by a detail-enhancement stage to produce the final image. The shown example is from the S24 dataset [16].

Ground-truth images are expected to have minimal noise and can be generated using third-party denoisers (e.g., the AI-based Adobe Lightroom denoiser) rather than physically captured references [1], which are time consuming to acquire and often lack scene diversity. This practical choice balances accuracy and generalization, allowing training on a broader range of raw data. Although such pseudo ground truths are imperfect, they provide stable supervision and help preserve fine image structures. A workflow for generating pseudo ground truths when third-party denoisers are unavailable is described in Sec. B.2 of the supplementary material, with additional training details in Sec. G.1.

## 2.2. Color Correction

After raw denoising, we apply a color-correction function $f_{\texttt{raw}\rightarrow\texttt{LsRGB}}(\cdot)$ that maps the denoised raw image to the linear sRGB domain. This function places the image in a camera-agnostic color space, making subsequent operators less dependent on camera-specific characteristics—unlike denoising and color correction, which are inherently camera-dependent. The function $f_{\texttt{raw}\rightarrow\texttt{LsRGB}}$ is defined as:

$$\mathbf{I}_{\texttt{LsRGB}(x,y)} = \mathbf{M}_{\texttt{CCM}}\big(\mathbf{D}_{\texttt{WB}}\,\mathbf{I}_{\texttt{enh-raw}(x,y)}\big), \qquad (2)$$

where $\mathbf{D}_{\texttt{WB}}$ is a diagonal matrix encoding the red and blue white-balance (WB) gains, and $\mathbf{M}_{\texttt{CCM}}$ is the color correction matrix (CCM) interpolated based on these gains using pre-calibrated camera-specific matrices [52]. WB gains are typically stored in DNG metadata and estimated by the camera's on-board auto white balance (AWB) module, but can also be predicted using learning-based AWB methods (e.g., [20, 52, 62]). Further discussion of learning-based AWB integration is provided in Sec. I.3 of the supplementary material.

## 2.3. Photofinishing

The photofinishing module finalizes the image's overall 'look and feel', covering both perceptual enhancements and artistic picture styles. Rather than a single black-box network, we design it as a modular module, since this stage largely determines the final visual appearance and is key to maintaining flexibility.

As illustrated in Fig. 2, the module operates on the down-sampled image $\mathbf{I}_{\texttt{LsRGB}}^{\downarrow}$ for efficiency and comprises five steps: 1) digital gain to adjust brightness, 2) global tone mapping (GTM) to refine global contrast, maintain perceptual brightness, and preserve highlights, 3) local tone mapping (LTM) to enhance local contrast and details, 4) chroma mapping to adjust chromaticities, and 5) gamma correction to produce a display-referred output. These are implemented by the functions $f_{\texttt{gain}}$, $f_{\texttt{GTM}}$, $f_{\texttt{LTM}}$, $f_{\texttt{chroma}}$, and $f_{\texttt{gamma}}$, each parameterized by image-specific coefficient(s) predicted by lightweight neural networks ($\sim$200K parameters in total): $\mathcal{D}_{\texttt{gain}}$, $\mathcal{D}_{\texttt{GTM}}$, $\mathcal{D}_{\texttt{LTM}}$, $\mathcal{D}_{\texttt{chroma}}$, and $\mathcal{D}_{\texttt{gamma}}$.

A key challenge is the lack of ground-truth supervision for individual photofinishing functions, making independent optimization infeasible. We therefore train the entire module end-to-end, with the main difficulty being to ensure that each function performs its intended role (e.g., GTM adjusts global tone rather brightness). Our design encourages such separation and interpretability, as detailed below.

The process begins with a digital gain adjustment:

$$\mathbf{I}_{\texttt{gain}}^{\downarrow} = f_{\texttt{gain}}\left(\mathbf{I}_{\texttt{LsRGB}}^{\downarrow}; d_g\right) = d_g\,\mathbf{I}_{\texttt{LsRGB}}^{\downarrow}, \qquad (3)$$

where $d_g$ is a global gain factor predicted by $\mathcal{D}_{\texttt{gain}}$. Next, tone mapping refines image contrast while maintaining perceptual brightness. Although one could apply tone mapping only to the luminance channel (i.e., leaving chroma to be modified solely by $f_{\texttt{chroma}}$), we found this design to

Figure 3. User-interactive photo-editing tool built on our modular ISP, providing full control over the rendering process, picture styles, and editing options. The interface supports selecting or interpolating between styles and adjusting white balance, exposure, color, and overall appearance. See the supplementary material (Sec. I) and the video (click to view) for details.

yield suboptimal results (see Sec. K.1.4 of the supplementary material). Instead, tone mapping is applied directly to the scaled linear sRGB channels of $\mathbf{I}_{\text{gain}}^{\downarrow}$, treating all channels equally. The GTM function, $f_{\text{GTM}}$, is formulated as:

$$\mathbf{I}_{\text{GTM}(x,y)}^{\downarrow(\rho)} = f_{\text{TM}}\big(\mathbf{I}_{\text{gain}(x,y)}^{\downarrow(\rho)}; a_{\text{GTM}}, b_{\text{GTM}}, c_{\text{GTM}}\big), \qquad (4)$$

where $(x, y)$ indexes spatial locations and $\rho \in \{\text{R, G, B}\}$. Here, $f_{\text{GTM}}$ is simply an application of the shared tone-mapping function $f_{\text{TM}}$ using global parameters. The tone-mapping function $f_{\text{TM}}$ is defined as:

$$f_{\text{TM}}(x; a_{\text{TM}}, b_{\text{TM}}, c_{\text{TM}}) = \frac{x^{a_{\text{TM}}}}{x^{a_{\text{TM}}} + \big(c_{\text{TM}}(1-x)\big)^{b_{\text{TM}}}}, \qquad (5)$$

where $x \in [0, 1]$ is the normalized RGB intensity. The image-specific parameters $a_{\text{GTM}}$, $b_{\text{GTM}}$, and $c_{\text{GTM}}$, predicted by $\mathcal{D}_{\text{GTM}}$, primarily control: 1) midtone contrast through the exponent $a_{\text{GTM}}$, 2) shadow compression via the slope parameter $b_{\text{GTM}}$, and 3) highlight roll-off by scaling $c_{\text{GTM}}$, which modulates the curvature at the upper end of the tone-mapping curve.

The LTM stage complements the GTM by providing spatially adaptive tone control. To achieve this, $\mathcal{D}_{\text{LTM}}$ comprises two subnetworks: 1) a multi-scale guidance subnetwork that outputs a guidance map $\mathbf{G}_{\text{guide}} \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times 1}$, and 2) a grid-prediction subnetwork that processes both $\mathbf{I}_{\text{gain}}^{\downarrow}$ and $\mathbf{I}_{\text{GTM}}^{\downarrow}$ (concatenated along channels) to predict a coarse grid of tone-mapping parameters $\mathbf{M}_{\text{LTM}} \in \mathbb{R}^{N_g \times N_g \times N_c \times 5}$, where $N_g, N_c \ll H/4, W/4$. Including both inputs provides the grid-prediction subnetwork with cues about the global tone-mapping behavior (from $\mathbf{I}_{\text{GTM}}^{\downarrow}$) and the pre-tonemapped image content (from $\mathbf{I}_{\text{gain}}^{\downarrow}$), helping it predict the coefficients applied to both inputs by the LTM function $f_{\text{LTM}}$, defined as:

$$\begin{aligned}
\mathbf{I}_{\text{LTM}(x,y)}^{\downarrow(\rho)} = {} & (1 - \mathbf{W}_{\text{LTM}(x,y)})\, \mathbf{I}_{\text{GTM}(x,y)}^{\downarrow(\rho)} \\
& + \mathbf{W}_{\text{LTM}(x,y)}\, f_{\text{TM}}\big(\mathbf{X}_{\text{LTM}(x,y)}^{(\rho)}; \mathbf{A}_{\text{LTM}(x,y)}, \\
& \mathbf{B}_{\text{LTM}(x,y)}, \mathbf{C}_{\text{LTM}(x,y)}\big),
\end{aligned} \qquad (6)$$

where $\mathbf{X}_{\text{LTM}(x,y)}$ is the locally scaled version of the gain-adjusted linear sRGB image:

$$\mathbf{X}_{\text{LTM}(x,y)}^{(\rho)} = \mathbf{I}_{\text{gain}(x,y)}^{\downarrow(\rho)} \mathbf{G}_{\text{LTM}(x,y)}, \qquad (7)$$

and the spatial coefficient maps:

$$\mathbf{A}_{\text{LTM}}, \mathbf{B}_{\text{LTM}}, \mathbf{C}_{\text{LTM}}, \mathbf{G}_{\text{LTM}}, \mathbf{W}_{\text{LTM}}' \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4}}$$

are generated by sampling the grid $\mathbf{M}_{\text{LTM}}$ via trilinear interpolation using the guidance map $\mathbf{G}_{\text{guide}}$. The preliminary map $\mathbf{W}_{\text{LTM}}'$ is then passed through a sigmoid activation to obtain $\mathbf{W}_{\text{LTM}}$. This formulation enables the LTM to build upon the globally tone-mapped image while flexibly re-tone-mapping the gain-adjusted image as needed. See Sec. K.1.4 in the supplementary material for ablations.

After tone mapping, chroma mapping refines the image's color components for enhancement or artistic stylization. We implement this step using $f_{\text{chroma}}$, an image-specific learnable 2D chroma LuT operating in the CbCr space. We opt for a 2D chroma LuT instead of predicting an image-specific 3D RGB LuT to reduce memory usage and ensure that this stage affects only chromaticity. Specifically, the tone-mapped image $\mathbf{I}_{\text{LTM}}^{\downarrow}$ is converted to YCbCr[1]. The chroma network $\mathcal{D}_{\text{chroma}}$ computes a differentiable 2D histogram of the CbCr channels (see Sec. D in the supplementary material for details), which is processed by a learnable encoder-decoder network. The encoder extracts chroma features, which are modulated by a latent vector derived from the Y channel. This vector is produced by a lightweight luminance-guidance subnetwork whose sigmoid-activated output scales the encoded chroma features for brightness-dependent modulation. The modulated features are decoded to predict a residual 2D LuT in CbCr, added to a learnable global base LuT to form the final $\mathbf{L}_{\text{chroma}}$. The LuT is applied to the chroma channels via bilinear interpolation, producing chroma-mapped CbCr values that are combined with Y and converted back to a pre-gamma, quasi-linear sRGB image, $\mathbf{I}_{\text{chroma}}^{\downarrow}$.

For artistic picture styles that involve stronger or more stylized color manipulations, we found that augmenting $f_{\text{chroma}}$ with an image-independent learnable 3D LuT improves color expressiveness. This optional step performs

---

[1] We use the YCbCr ↔ RGB conversion matrices defined in ITU-R BT.709, assuming a 2.2 gamma-encoded sRGB space. Although our input is not strictly 2.2 gamma, we adopt this approximation for simplicity and differentiability.

a trilinear lookup over an $11\times11\times11$ LuT ($\mathbf{L}_{\text{RGB}}$) on the tone-mapped RGB values before $\mathcal{D}_{\text{chroma}}$ and $f_{\text{chroma}}$. The 3D LuT can capture rich color transformations that complement the subsequent 2D chroma LuT ($\mathbf{L}_{\text{chroma}}$), enhancing expressiveness in artistic picture modes. However, since learning $\mathbf{L}_{\text{RGB}}$ can interfere with the intended role of other tone-mapping stages, we keep it optional.

At the final stage of photofinishing, we apply gamma correction using $f_{\text{gamma}}$, defined as:

$$\mathbf{I}^{\downarrow}_{\text{gamma}} = f_{\text{gamma}}\left(\mathbf{I}^{\downarrow}_{\text{chroma}}; \gamma\right) = \left(\mathbf{I}^{\downarrow}_{\text{chroma}}\right)^{(1/\gamma)}, \quad (8)$$

where $\gamma$ is an image-specific factor predicted by $\mathcal{D}_{\text{gamma}}$.

We train the five networks of our photofinishing module ($\mathcal{D}_{\text{gain}}$, $\mathcal{D}_{\text{GTM}}$, $\mathcal{D}_{\text{LTM}}$, $\mathcal{D}_{\text{chroma}}$, and $\mathcal{D}_{\text{gamma}}$), and optionally $\mathbf{L}_{\text{RGB}}$, jointly in an end-to-end manner by minimizing:

$$\begin{aligned}
\mathcal{L}_{\text{total}} = {} & \lambda_1\,\ell_1 + \lambda_{\text{SSIM}}\,\ell_{\text{SSIM}} + \lambda_{\Delta E}\,\ell_{\Delta E} + \lambda_{\text{perc}}\,\ell_{\text{perc}} \\
& + \lambda_{\text{CbCr}}\,\ell_{\text{CbCr}} + \lambda_{\text{LuT-s}}\,\ell_{\text{LuT-s}} + \lambda_{\text{TM}}\,\ell_{\text{TM}} \\
& + \lambda_{\text{LTM-s}}\,\ell_{\text{LTM-s}} + \lambda_{\text{luma}}\,\ell_{\text{luma}}.
\end{aligned} \quad (9)$$

The total loss integrates fidelity, perceptual, and regularization components. Low-level terms ($\ell_1$, $\ell_{\text{SSIM}}$, $\ell_{\text{CbCr}}$) ensure pixel- and structure-level accuracy; $\ell_{\text{CbCr}}$ measures chroma error between the predicted and de-gammaed (using the predicted $\gamma$) ground-truth CbCr channels. Perceptual terms ($\ell_{\Delta E}$ and $\ell_{\text{perc}}$) enforce perceptual color and feature similarity, where $\ell_{\Delta E}$ is a differentiable CIE $\Delta E$ metric and $\ell_{\text{perc}}$ is VGG-based. Regularization terms ($\ell_{\text{LuT-s}}$, $\ell_{\text{LTM-s}}$, $\ell_{\text{TM}}$, and $\ell_{\text{luma}}$) stabilize learning and improve interpretability: $\ell_{\text{LuT-s}}$ and $\ell_{\text{LTM-s}}$ are total-variation smoothness penalties for $\mathbf{L}_{\text{chroma}}$ and the LTM coefficient maps. $\ell_{\text{TM}}$ enforces luminance consistency between the downsampled global and full-resolution local tone-mapped Y channels and the Y channel of the de-gammaed ground truth, balancing global contrast and local detail refinement, while encouraging the GTM and LTM subnetworks to complement each other (avoiding dominance by either). Finally, $\ell_{\text{luma}}$ regularizes the GTM stage to preserve the average brightness of the gain-adjusted image, ensuring contrast refinement without altering global brightness.

The coefficients $\lambda_j$ control the relative strength of each term. Details of the loss definitions, weighting factors, training setup, and ablation studies are provided in the supplementary material (Secs. F, G.2, and K.1.3).

## 2.4. Upsampling

For efficiency, photofinishing is performed on a downscaled image, $\mathbf{I}^{\downarrow}_{\text{LsRGB}}$, and produces $\mathbf{I}^{\downarrow}_{\text{gamma}}$, which is then upsampled using the high-resolution linear sRGB image $\mathbf{I}_{\text{LsRGB}}$ as guidance. We adopt bilateral grid upsampling (BGU) [29], which computes an affine transform per grid cell by solving a regularized least-squares system. The original Halide

BGU constrains regularization to be achromatic (a single scalar gain across channels) and uses grid blurring to handle empty cells—both of which introduce limitations: 1) enforced achromaticity causes color crosstalk, and 2) grid blurring trades off detail for smoothness. We address these issues with per-channel gated regularization that removes the need for grid blurring. Each cell is regularized independently by channel, while empty cells fall back to global per-channel gains derived from pooled grid statistics. This design yields sharper, more faithful reconstructions (see Sec. E of the supplementary material for details). The guided upsampling produces the photofinished image $\mathbf{I}^{\uparrow}_{\text{gamma}} \in \mathbb{R}^{H\times W\times 3}$, which is then refined by the detail-enhancement stage.

## 2.5. Detail Enhancement

The final stage of our pipeline applies a detail-enhancement step to compensate for residual artifacts from denoising and guided upsampling. The output image $\mathbf{I}_{\text{out}}$ is obtained as:

$$\mathbf{I}_{\text{out}} = f_{\text{enh}}\left(\mathbf{I}^{\uparrow}_{\text{gamma}}\right), \quad (10)$$

where $f_{\text{enh}}(\cdot)$ is implemented as a compact fully convolutional network $\mathcal{D}_{\text{enh}}$. To train $\mathcal{D}_{\text{enh}}$, we first generate $\mathbf{I}^{\uparrow}_{\text{gamma}}$ for each training image using all preceding stages (with pretrained models), and use them as inputs. The network is optimized with a pixel-wise $\ell_1$ loss between the predicted and ground-truth sRGB images. Additional details are provided in Sec. G.3 of the supplementary material.

We explored merging the raw- and detail-enhancement stages by fine-tuning $\mathcal{D}_{\text{raw}}$ after training the photofinishing module, aiming to pre-enhance image details before resampling and thereby mitigate detail loss. However, this approach was unstable and often failed to converge. Using $\mathcal{D}_{\text{enh}}$ as a separate final stage (with only ~50K parameters) proved more robust and easier to train.

## 2.6. Photo-Editing Tool

To demonstrate the advantages of our modular design, we developed a user-interactive photo-editing tool built on top of our neural ISP. The tool provides full control over the *entire* rendering process, supporting multiple picture styles and additional editing adjustments (e.g., highlights, shadows, contrast, and exposure) directly within the ISP pipeline (see Fig. 3). To enhance camera-agnostic usability, we integrated a "generic" denoiser trained on a diverse set of synthetic and real noisy images, aimed at improving robustness to unseen noise patterns (see Sec. H of the supplementary material). The modular design also allows users to either apply the camera's AWB estimates or recompute WB gains using recent illuminant estimation models [11, 16, 92], including both camera-specific [16] and cross-camera variants [11].

Figure 4. Qualitative comparison between our method and recent neural ISP methods (ISPDiffuser [70], LiteISP [91], and ParamISP [53]) on an example from the S24 test set [16]. Results are shown for the default picture style (Style #0) and the remaining artistic styles (Styles #1–5). PSNR values with respect to the ground truth are shown in the lower-left corner of each image.

We draw inspiration from recent work on raw image compression [15, 56, 80] and use the learning-based method of [15] to embed the compressed raw data into the final JPEG image. This enables unlimited post-editable re-rendering under new settings with only a modest file size increase. Additionally, a lightweight linearization network [10] is included to synthesize a raw-like representation from external sRGB images, allowing the tool to process native DNG files, sRGB JPEGs saved by our tool with embedded raw, or standard sRGB inputs. Despite its versatility, the entire tool (including photofinishing networks for multiple picture styles) requires only ∼3.9 M parameters, far fewer than competing neural ISPs (e.g., ISPDiffuser [70], ∼20.9 M parameters for a single style; see Table 1 for parameter comparisons across methods). Comprehensive details and additional demonstrations are provided in the supplementary material (Sec. I) and in the accompanying video (click to view).

## 3. Experimental Results

We evaluated our modular ISP framework against recent neural ISP methods, including black-box end-to-end models (PyNet [48], LAN [69], LiteISP [91], Invertible-ISP [85], FourierISP [44], MicroISP [49], ISPDiffuser [70]) and multi-stage architectures (CIE XYZ Net [10], ParamISP [53], and FlexISP [60]). We used the S24 test

dataset [16], which provides all data needed to train and evaluate our framework. The dataset includes pseudo ground-truth denoised raw images generated by Adobe Lightroom's AI-based denoiser (for training our denoising module) and six ground-truth sRGB images per raw input, corresponding to one default style (Style #0) and five artistic styles (Styles #1–5), making it a suitable benchmark for evaluating the scalability of our method in handling multiple styles. The dataset comprises 2,619 training, 205 validation, and 400 test pairs. For results on the MIT-Adobe FiveK dataset [26], see Sec. K.2.2 in the supplementary material.

We trained three denoising variants ($\mathcal{D}_{raw}$): lite (0.25 M parameters), base (0.93 M), and large (3.6 M), described in Sec. C.1 of the supplementary material. These yield three configurations of our full pipeline. Each module (denoising, photofinishing, and detail enhancement) was trained separately. Owing to our modular design, only the photofinishing and detail-enhancement networks are style-specific, while the denoiser is shared across all styles. This greatly reduces training and memory requirements when supporting multiple picture styles, unlike prior methods that must retrain and load the entire ISP pipeline into memory to adapt to new styles. All baselines were trained per style using their official implementations (see Sec. L in the supplementary material for details).

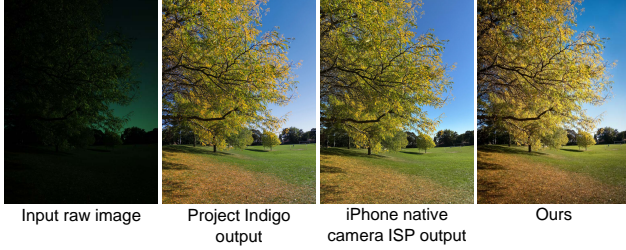| Input raw image | Project Indigo output | iPhone native camera ISP output | Ours |

Figure 5. Comparison among Project Indigo [55], the iPhone native camera ISP, and our method (using the generic denoiser and cross-camera auto white balance). The image was captured using the iPhone 13 Pro Max main camera.

## 3.1. Quantitative Results

We report the peak signal-to-noise ratio (PSNR), structural similarity index measure (SSIM) [81], LPIPS [89], and $\Delta E_{2000}$ [72] for quantitative evaluation. Table 1 summarizes results on the default S24 style along with each model's parameter count. We present results for our three variants (lite, base, and large) and ablations with/without detail enhancement.

Further ablations are detailed in Sec. K.1 of the supplementary material, covering: 1) denoising (Sec. K.1.1), 2) guided upsampling (Sec. K.1.2), 3) photofinishing loss design (Sec. K.1.3), and 4) photofinishing design (Sec. K.1.4). As shown in Table 1, our approach achieves state-of-the-art results across all variants (even the lite model with 0.5M parameters), while LiteISP [91], the closest competitor, uses 9M parameters and yields ∼2 dB lower PSNR.

Table 2 reports PSNR for all artistic styles (Styles #1–5), with full metrics (SSIM, LPIPS, and $\Delta E_{2000}$) provided in Sec. K.2.4 of the supplementary material. For completeness, we also compare with cmKAN [67], a lightweight color-matching method trained to transfer colors from the default style to the remaining styles. At test time, we rendered images using our pipeline in the default style and applied cmKAN as a post-processing color transfer (denoted 'PP (cmKAN)' in Table 2). This experiment aims to highlight the advantage of performing style rendering directly from raw data vs. applying color transfer as a post-processing step after rendering to the default style.

We further report the total parameters required to support all five artistic styles. Because our denoiser is shared, the added cost per style is minimal—a key benefit of our modular design. Table 2 lists results for our lite, base, and large variants with/without detail enhancement and the optional 3D LuT ($\mathbf{L}_{RGB}$).

As shown, our method achieves state-of-the-art results across all styles while requiring a moderate number of parameters to support multiple styles compared to other methods. Furthermore, learning $\mathbf{L}_{RGB}$ consistently improves the final results. However, we found that when training for the default style—which primarily targets natural color rendering and high image quality rather than aggressive artistic

Table 1. Results on the S24 test set [16]. We report PSNR, SSIM [81], LPIPS [89], and $\Delta E$ 2000 [72], along with the total number of parameters for each method. Our method is evaluated with different denoising model capacities (lite, base, and large), with and without the enhancement network. The best results are highlighted in <mark>yellow</mark>. Our method achieves state-of-the-art results, offers a modular design with full ISP control, requires a moderate number of parameters, and runs efficiently on a single GPU (∼0.7 sec with the lite denoising model, ∼0.95 sec with the base model, and ∼1.4 sec with the large model on an NVIDIA GeForce RTX 4080 SUPER).

| Method | S24 Test Set | | | | # params |
| | PSNR↑ | SSIM↑ | LPIPS↓ | ΔE 2000↓ | |
|---|---|---|---|---|---|
| ISPDiffuser [70] | 24.03 | 0.881 | 0.159 | 5.918 | 20,938,890 |
| PyNet [48] | 23.80 | 0.869 | 0.100 | 6.277 | 47,548,170 |
| CIE XYZ Net [10] | 23.32 | 0.860 | 0.124 | 7.024 | 1,348,789 |
| FlexISP [60] | 24.01 | 0.886 | 0.110 | 6.938 | 25,705,065 |
| Invertible-ISP [85] | 22.87 | 0.820 | 0.147 | 7.374 | 1,413,760 |
| LAN [69] | 23.11 | 0.812 | 0.116 | 6.765 | 46,847 |
| MicroISP [49] | 20.55 | 0.775 | 0.180 | 11.012 | 13,560 |
| ParamISP [53] | 24.32 | 0.841 | 0.115 | 6.135 | 1,420,000 |
| LiteISP [91] | 25.49 | 0.897 | 0.074 | 5.521 | 9,094,000 |
| FourierISP [44] | 24.50 | 0.913 | 0.096 | 5.928 | 7,589,736 |
| Ours (lite, w/o enhancement) | 26.36 | 0.878 | 0.071 | 4.413 | 452,447 |
| Ours (base, w/o enhancement) | 26.48 | 0.883 | 0.065 | 4.282 | 1,139,907 |
| Ours (large, w/o enhancement) | 26.51 | 0.884 | 0.064 | 4.253 | 3,841,547 |
| Ours (lite, w/ enhancement) | 27.37 | 0.916 | 0.060 | 4.059 | 503,082 |
| Ours (base, w/ enhancement) | 27.52 | 0.922 | 0.055 | 3.938 | 1,190,542 |
| Ours (large, w/ enhancement) | **27.57** | **0.923** | **0.054** | **3.913** | 3,892,182 |

Table 2. Results across the S24 dataset picture styles (Styles #1–5) [16]. We report the average PSNR for each target style, along with the total number of parameters required to support all five styles. The best results in each column are highlighted in <mark>yellow</mark>.

| Method | S24 Test Set | | | | | # params |
| | S #1 | S #2 | S #3 | S #4 | S #5 | (for all styles) |
|---|---|---|---|---|---|---|
| ISPDiffuser [70] | 25.60 | 27.30 | 25.02 | 25.93 | 26.83 | 104,694,450 |
| PyNet [48] | 24.36 | 25.94 | 24.70 | 24.34 | 26.32 | 237,740,850 |
| CIE XYZ Net [10] | 22.40 | 24.05 | 22.00 | 22.26 | 24.67 | 6,743,945 |
| PP (cmKAN) [67] | 20.93 | 23.04 | 21.85 | 20.91 | 21.3 | 384,535 |
| FlexISP [60] | 24.86 | 27.47 | 25.23 | 23.96 | 24.65 | 128,525,325 |
| Invertible-ISP [85] | 23.48 | 26.35 | 23.84 | 23.33 | 24.90 | 7,068,800 |
| LAN [69] | 22.98 | 23.74 | 23.47 | 22.80 | 25.38 | 234,235 |
| MicroISP [49] | 20.30 | 23.66 | 21.45 | 20.34 | 22.68 | 67,800 |
| ParamISP [53] | 24.97 | 27.11 | 24.77 | 24.18 | 25.43 | 7,100,000 |
| LiteISP [91] | 26.66 | 28.33 | 26.31 | 25.04 | 28.07 | 45,470,000 |
| FourierISP [44] | 25.19 | 28.03 | 25.38 | 24.74 | 27.41 | 37,948,680 |
| Ours (lite, w/o enh.) | 25.16 | 28.09 | 25.66 | 25.47 | 27.08 | 1,281,343 |
| Ours (base, w/o enh.) | 25.28 | 28.19 | 25.73 | 25.55 | 27.19 | 1,968,803 |
| Ours (large, w/o enh.) | 25.31 | 28.22 | 25.75 | 25.58 | 27.23 | 4,670,443 |
| Ours (lite, w/o enh., w/ 3D LuT) | 26.39 | 29.21 | 26.69 | 26.35 | 27.93 | 1,347,950 |
| Ours (base, w/o enh., w/ 3D LuT) | 26.52 | 29.29 | **26.79** | 26.47 | 28.13 | 2,035,410 |
| Ours (large, w/o enh., w/ 3D LuT) | 26.56 | **29.31** | 26.83 | 26.51 | 28.19 | 4,737,050 |
| Ours (lite, w/ enh., w/ 3D LuT) | 26.56 | 28.92 | 26.78 | 26.66 | 28.73 | 1,550,490 |
| Ours (base, w/ enh., w/ 3D LuT) | **26.71** | 28.99 | 26.78 | **26.79** | 28.95 | 2,237,950 |
| Ours (large, w/ enh., w/ 3D LuT) | **26.75** | 29.01 | **26.83** | **26.84** | **29.03** | 4,939,590 |

stylization—learning the 3D LuT may not provide benefits and can even lead to slight degradations in some metrics. See Sec. K.1.4 of the supplementary material for further ablation analysis.

## 3.2. Qualitative Results

Figure 4 shows a qualitative comparison between our method and recent neural ISP approaches [53, 70, 91] on an example from the S24 test set. Rendered images are shown for the default picture style (Style #0) and the artistic styles, along with their corresponding ground truths. Our method consistently delivers higher visual quality across all styles. Additional examples are provided in Sec. K.2.5 of the supplementary material.

**Cross-Camera Generalization.** As described in Sec. 2.6, our photo-editing tool integrates generic denoisers and cross-camera AWB models to extend applicability to unseen cameras. Figure 5 presents a qualitative comparison on an image captured with the iPhone 13 Pro Max main camera, comparing our result (using the generic denoiser and cross-camera AWB) with the native iPhone ISP and Project Indigo [55]. Our method delivers visual quality comparable to both, despite not using any iPhone data during training. This capability is further illustrated in Fig. 1, showing results on another camera unseen during training. Such generalization arises from our modular design, which allows switching between camera-specific models (e.g., trained on the S24 dataset) and generic ones for unseen devices, while maintaining visually pleasing output. Additional examples and discussion on cross-camera generalization are provided in the supplementary material (Secs. H and I.3).

**User Study.** To evaluate perceptual quality, we conducted a user study comparing our method against the Samsung S24 native camera ISP and Adobe Lightroom. We captured 45 scenes with the S24 main camera in Pro mode (saving DNG files) and re-captured the same scenes using the native app. The captured scenes covered indoor, outdoor daylight, sunset, and low-light conditions. DNG files were processed using Adobe Lightroom (auto enhancement) and our method. For each scene, participants viewed three versions (ours, native ISP, and Lightroom) in random order and selected their preferred image under four criteria: 'color quality', 'brightness & contrast', 'sharpness & detail', and 'overall preference'. Twenty participants took part in the study. Our method was consistently preferred, achieving 53.2% in 'color quality', 46.4% in 'brightness & contrast', 43.4% in 'sharpness & detail', and 51.4% in 'overall preference', out of three versions compared, outperforming the closest competitor by +13.9–27.0%. More details are provided in Sec. J of the supplementary material.

## 4. Conclusion and Discussion

In this paper, we have presented a learning-based modular ISP framework that offers fine-grained control across the raw-to-sRGB image rendering process. We evaluated three variants of our framework—lite (∼0.5 M parameters), base (∼1.2 M), and large (∼3.9 M)—all of which have lightweight to moderate capacity and can process a 12-megapixel image in about a second or less on a single GPU. Across all variants, our method achieves state-of-the-art results on nearly all picture styles in the S24 dataset [16], and competitive results on the MIT-Adobe FiveK dataset [26] (see Sec. K.2.2 in the supplementary material), while requiring significantly fewer parameters than the closest competing methods. Beyond high-quality
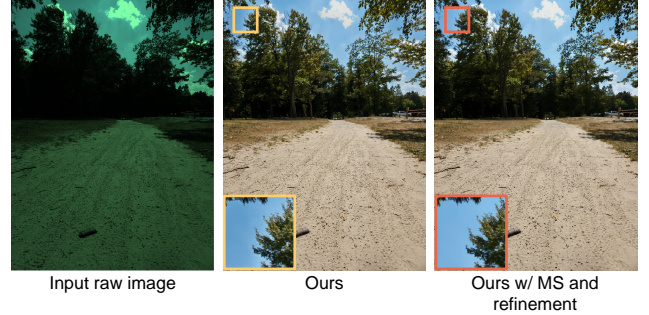


Figure 6. Our pipeline may produce halo artifacts or color inconsistencies in flat regions near edges. These artifacts can be mitigated via multi-scale (MS) processing and post-refinement of LTM maps. Image captured using the S24 main camera.

rendering, our framework consists of interpretable stages, making debugging and handling corner cases easier compared to previous black-box, non-interpretable end-to-end neural ISP designs. Furthermore, the modularity of our design makes scalability more practical and enables better generalization to unseen cameras. Our framework also offers strong flexibility throughout the rendering process. To highlight this flexibility, we implemented a user-interactive photo-editing tool built on top of our modular ISP that enables fine-grained control over the rendering pipeline, including picture styles and editing options.

Despite these advantages, our method faces a few challenges. First, in certain corner cases, halo artifacts may appear near edges in backlit scenes. Because of the modular structure of our framework, we were able to analyze this issue and identify its origin in the LTM process. These artifacts can be mitigated by applying multi-scale (MS) processing and post-processing refinement of the predicted LTM maps. We discuss this in detail in the supplementary material (Sec. B.1) and show an example in Fig. 6, comparing the original artifact and the mitigated result using MS and refinement.

Another challenge in training our framework is the need for reference denoised images to supervise the denoising modules, along with camera AWB and CCM data for color correction. While pseudo ground-truths for denoising can be obtained using AI-based third-party denoisers, the absence of DNG files makes it difficult to generate such pseudo ground-truth data or to obtain the necessary color-correction data for training. In Sec. B.2 of the supplementary material, we discuss this issue and present a practical strategy for obtaining the missing data when DNG files are unavailable. Using the Zurich Raw-to-sRGB dataset [48] as a case study, we demonstrate how our workflow can train the proposed method effectively, achieving results comparable to the top-performing methods despite the lack of DNG files and with roughly half or fewer parameters compared to existing alternatives—while maintaining a high degree of modularity not offered by other methods.

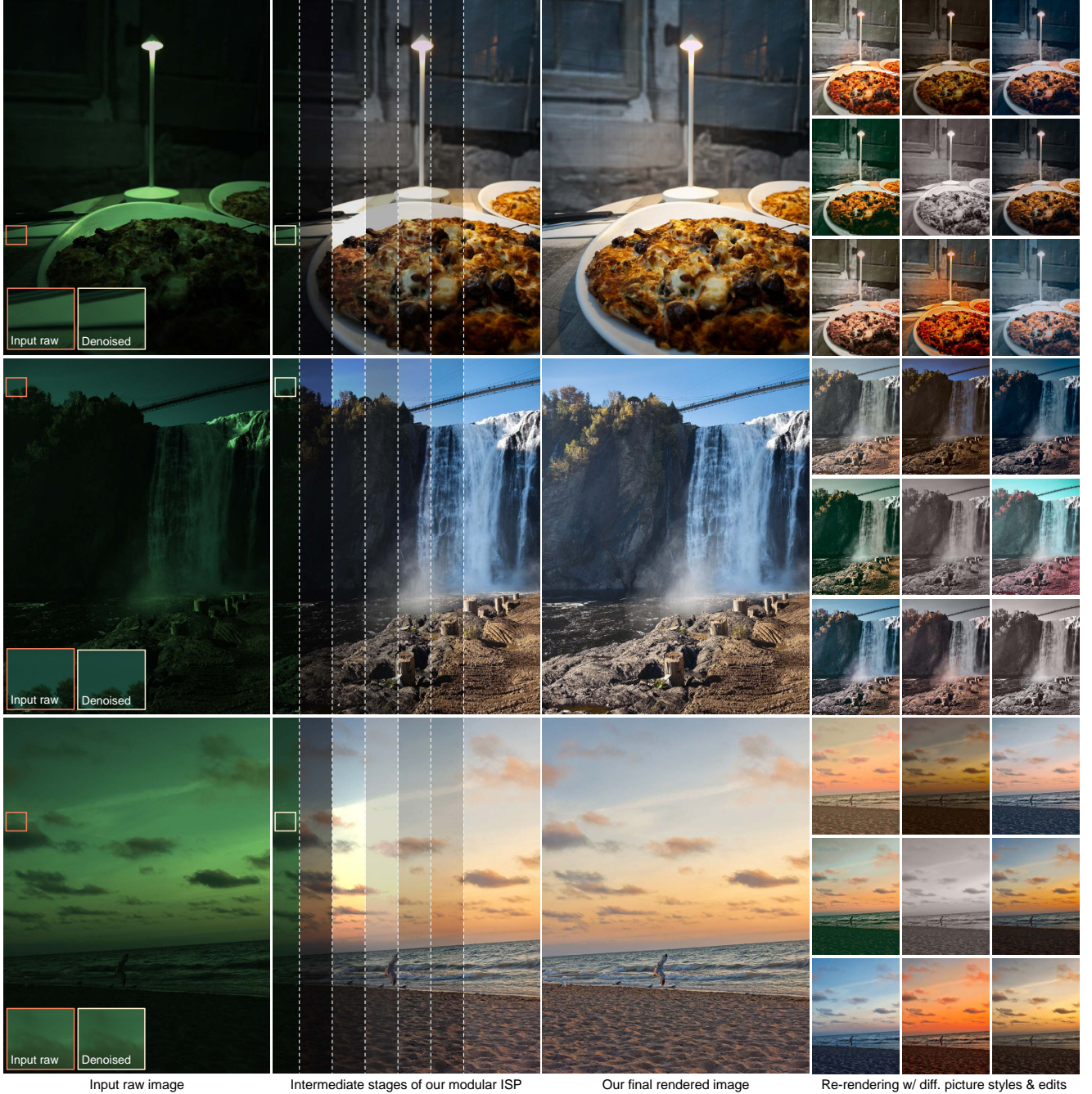# Modular Neural Image Signal Processing

## Supplementary Material



Figure 7. Intermediate outputs of our modular neural ISP, including denoising, color correction, digital gain, global tone mapping, local tone mapping, chroma mapping, gamma correction and detail enhancement. Different picture styles and edits are also shown (right). From top to bottom and left to right (for the styles and edits on the right), the first five images correspond to Styles #1–5 of the S24 dataset [16], while the remaining ones show results with additional edits and style adjustments. The first example was captured using an iPhone 15 main camera, the second using an iPhone 13 main camera, and the third using a Samsung S24 main camera. Note that none of the training images were captured with, or include any synthetic data derived from, iPhone devices, underscoring the robustness and generalization capability of our method to unseen cameras.

1

In the main paper, we presented our modular neural image signal processing (ISP) framework, which enables accurate raw-to-sRGB rendering through a modular design that provides full control over the rendering process and supports different picture styles. This modularity allows the framework to be flexibly tuned to produce desired capture-time outputs or to function as an interactive photo editor (see Fig. 7). We provide this supplementary material to further clarify how our method is developed, how it can be configured for capture-time deployment, and how it can alternatively be used as an interactive image-processing tool. Specifically, in Sec. A, we describe our GPU-accelerated iterative bilateral solver, which we use to mitigate halo artifacts that may occur in some corner cases. Handling challenging artifacts (with the help of the GPU-accelerated bilateral solver) and training with incomplete datasets are discussed in detail in Sec. B.

We then elaborate on the design of the deep networks used in this work in Sec. C and provide additional information about the differentiable histogram used in the chroma mapping network in Sec. D. Afterwards, we describe the guided upsampling regularization used in our implementation in Sec. E. Details of the photofinishing loss functions are provided in Sec. F, and the training details of our networks are presented in Sec. G.

In Sec. H, we discuss our efforts to improve the generalization of our method across cameras. Sec. I describes our graphical user interface tool, built on top of our method, which includes additional editing capabilities. In Sec. J, we present further details of the user study conducted as part of our evaluation. Lastly, in Sec. K, we provide extensive ablation studies of the pipeline components, along with additional results and comparisons. Further evaluation details for both the main paper and this supplementary material are provided in Sec. L.

## A. GPU-Accelerated Iterative Bilateral Solver

To mitigate residual artifacts that occasionally appear in corner cases of our local tone-mapping (LTM) stage, we add an optional guided refinement step. Edge-aware smoothing is a natural choice for this purpose. The Fast Bilateral Solver (FBS) [19] is particularly effective: it minimizes a quadratic objective with bilateral affinities, thereby enforcing edge adherence while propagating low-frequency information. However, the original solver is not GPU-friendly. FBS constructs a sparse bilateral grid and solves a large symmetric positive definite (SPD) system using preconditioned conjugate gradient (PCG). This involves repeated gather-scatter operations between pixel and grid spaces, dominated by memory access rather than arithmetic, and converges slowly under the Jacobi preconditioner. These factors make direct GPU acceleration of FBS inefficient.

### A.1. Quadratic Objective in Image Space

We propose a lightweight image-space variant that preserves the FBS energy while avoiding the bilateral grid. Given an initial tensor $\mathbf{M}$ and a guidance image $\mathcal{Z}$, we solve for a refined output $\mathbf{Y}$ by minimizing the following energy:

$$\min_{\mathbf{Y}} \ \lambda \sum_p \|\mathbf{Y}_p - \mathbf{M}_p\|_2^2$$
$$+ \sum_p \sum_{q \in \mathcal{N}_k(p)} \mathbf{W}_{pq}(\mathcal{Z}) \|\mathbf{Y}_p - \mathbf{Y}_q\|_2^2, \quad (11)$$

where $\mathcal{N}_k(p)$ denotes a local $k{\times}k$ neighborhood around pixel $p$, and $\mathbf{W}_{pq}(\mathcal{Z})$ are bilateral affinities defined based on the guidance image $\mathcal{Z}$:

$$\mathbf{W}_{pq}(\mathcal{Z}) = \exp\Big( - \frac{\|p - q\|_2^2}{2\sigma_s^2} - \frac{(\mathcal{Z}_p - \mathcal{Z}_q)^2}{2\sigma_r^2} \Big), \quad (12)$$

with $(\sigma_s, \sigma_r)$ controlling the spatial and range scales. The guidance image, $\mathcal{Z}$, is computed as the luminance of the input RGB image using a weighted combination of the red, green, and blue channels, with respective weights of 0.2989, 0.5870, and 0.1140.

### A.2. Iterative Solver

Instead of solving Eq. 11 with PCG, we pre-compute the bilateral weights $\mathbf{W}_{pq}(\mathcal{Z})$ once (per image) and apply a fixed number of successive over-relaxation (SOR) updates [86]. Initializing $\mathbf{Y}^{(0)} = \mathbf{M}$, each pixel $p$ is updated as:

$$\tilde{\mathbf{Y}}_p^{(t+1)} = \frac{\lambda \mathbf{M}_p + \sum_{q \in \mathcal{N}_k(p)} \mathbf{W}_{pq}(\mathcal{Z}) \mathbf{Y}_q^{(t)}}{\lambda + \sum_{q \in \mathcal{N}_k(p)} \mathbf{W}_{pq}(\mathcal{Z})}, \quad (13)$$

$$\mathbf{Y}^{(t+1)} = \mathbf{Y}^{(t)} + \omega\Big( \tilde{\mathbf{Y}}^{(t+1)} - \mathbf{Y}^{(t)} \Big), \quad (14)$$

with relaxation $\omega \in [1, 2)$. In our implementation, $\mathbf{W}_{pq}(\mathcal{Z})$ are normalized such that $\sum_{q \in \mathcal{N}_k(p)} \mathbf{W}_{pq}(\mathcal{Z}) = 1$, making the denominator $\lambda + 1$. The same normalized weights are reused across channels and iterations.

### A.3. GPU Implementation and Efficiency

All steps in Eq. 14 are implemented using dense tensor primitives (`unfold`, pointwise operations, and reductions) that are highly optimized on GPUs (see Algorithm 1). We use reflective padding to avoid border bias and pre-compute $\mathbf{W}_{pq}(\mathcal{Z})$ once per image, reusing it across channels and iterations. Unless otherwise stated, our implementation uses the following default values: $k{=}7$, $\sigma_s{=}3.0$ px, $\sigma_r{=}0.01$, $\lambda{=}10^{-3}$, $n_{\text{iter}}{=}80$, and $\omega{=}1.6$.

With $n_{\text{iter}}{=}80$, our GPU-based iterative solver achieves an $\approx 11\times$ wall-clock speedup over the original CPU-based FBS while producing visually comparable refinements. On CPU, however, it is slower than FBS due to the overhead

**Algorithm 1** GPU-Accelerated Iterative Bilateral Solver

**Require:** Guidance $\mathcal{Z}$, input tensor $\mathbf{M}$, kernel size $k$, scales $(\sigma_s, \sigma_r)$, smoothness $\lambda$, iterations $n_{\text{iter}}$, relaxation $\omega$

1: Pad $\mathcal{Z}$ with `reflect`, extract $k \times k$ neighborhoods using `unfold`
2: Compute range weights $\exp(-(\Delta \mathcal{Z})^2 / 2\sigma_r^2)$
3: Compute spatial weights $\exp(-\|\Delta p\|^2 / 2\sigma_s^2)$
4: Form bilateral affinities $\mathbf{W}_{pq}$ and normalize so $\sum_q \mathbf{W}_{pq} = 1$
5: Initialize $\mathbf{Y}^{(0)} \leftarrow \mathbf{M}$
6: **for** $t = 0$ to $n_{\text{iter}} - 1$ **do**
7: $\quad$ Extract $k \times k$ neighborhoods of $\mathbf{Y}^{(t)}$ with `unfold`
8: $\quad$ Compute smooth term $\sum_q \mathbf{W}_{pq} \mathbf{Y}_q^{(t)}$
9: $\quad$ Compute target $(\lambda \mathbf{M}_p + \text{smooth})/(\lambda + 1)$
10: $\quad$ Update $\mathbf{Y}^{(t+1)} \leftarrow \mathbf{Y}^{(t)} + \omega(\text{target} - \mathbf{Y}^{(t)})$
11: **end for**
12: **return** $\mathbf{Y}^{(n_{\text{iter}})}$



Figure 8. Runtime of the guided refinement process on a $750 \times 1000$ guide image for different iteration counts $n_{\text{iter}}$. We compare the Fast Bilateral Solver (FBS) [19] on CPU with our modified bilateral refinement on both CPU and GPU. Since our pipeline is primarily intended for GPU deployment, the modified bilateral refinement provides an efficient and practical replacement for FBS. Runtimes were measured on an Intel Core i7-14700K CPU and an NVIDIA GeForce RTX 4080 SUPER GPU (16 GB VRAM).

of repeated dense tensor operations (see Fig. 8 for runtime vs. $n_{\text{iter}}$ on a $750 \times 1000$ guidance image). Since our framework assumes GPU availability in most deployment scenarios, the GPU-accelerated solver offers a clear practical advantage.

Our GPU-accelerated solver retains the FBS energy [19] but replaces the bilateral-grid PCG solver with a local SOR scheme. While this sacrifices exact convergence guarantees, it enables efficient and straightforward GPU implementation with sufficient quality in practice. Beyond mitigating artifacts in our LTM module (see Fig. 9), the proposed GPU-accelerated refinement can also serve as a *general* edge-aware post-processing method (see Fig. 10). Providing a thorough evaluation of our modified solver across different datasets and tasks is beyond the scope of this paper, but we consider this an important direction for future work.

## B. Challenges and Potential Solutions

In this section, we discuss key challenges of the proposed framework along with potential solutions. We begin with artifacts in the LTM stage, where in rare cases (particularly under strong backlighting), halo artifacts may appear. To address this, we introduce two optional lightweight steps that can be user-enabled, since applying them indiscriminately may reduce accuracy (as explained in the next subsection). We next consider the case of training on an incomplete dataset that lacks some of the essential information required by our method. For this purpose, we use the Zurich Raw-to-sRGB dataset [48], which, in addition to the aforementioned challenges, introduces further difficulties due to unaligned raw and ground-truth sRGB training pairs.
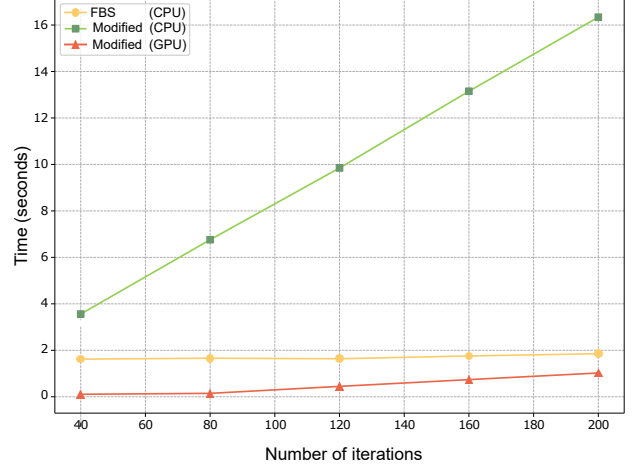
Specifically, we describe how our framework can still be trained when key data are missing, such as ground-truth denoised raw images or DNG metadata (e.g., illuminant color and color correction matrices).

### B.1. Mitigating Artifacts

We optionally apply two lightweight steps that suppress halos while preserving edges: 1) a multi-scale aggregation of LTM coefficient predictions, followed by 2) the edge-aware bilateral refinement discussed in Sec. A.

We adopt a multi-scale strategy to improve robustness and suppress halo artifacts in challenging cases (e.g., strong backlighting). Specifically, the input image after digital gain $\mathbf{I}_{\text{gain}}^{\downarrow}$ and its global tone-mapped version $\mathbf{I}_{\text{GTM}}^{\downarrow}$ are processed at progressively downsampled resolutions using scale factors: $\mathcal{S} = \{1.0, 0.5, 0.25, 0.125, 0.0625\}$.

At each scale $s \in \mathcal{S}$, if $s \neq 1.0$, both $\mathbf{I}_{\text{gain}}^{\downarrow}$ and $\mathbf{I}_{\text{GTM}}^{\downarrow}$ are bilinearly downsampled; otherwise the original resolution is used. From the input $\mathbf{I}_{\text{gain}(s)}^{\downarrow}$ at each scale $s$, we then generate a guidance map using our multi-scale guidance subnetwork (see Sec. C.5), followed by a smoothing step. In particular, the guidance map is first extended using reflection padding. We then apply average pooling over a $5 \times 5$ spatial window, producing a gently varying guidance signal that stabilizes the slicing coefficients while preserving the overall luminance structure needed for high-quality upsampling.

The pair $(\mathbf{I}_{\text{gain}(s)}^{\downarrow}, \mathbf{I}_{\text{GTM}(s)}^{\downarrow})$ is concatenated and fed into
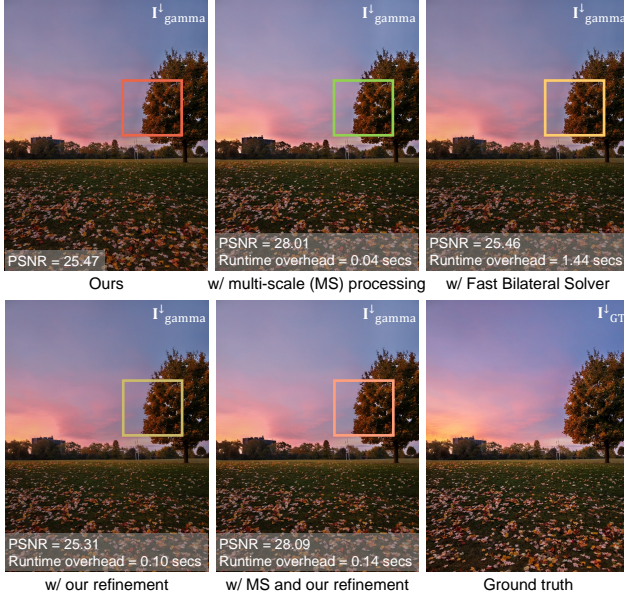
Figure 9. Comparison of 1) multi-scale processing for generating LTM coefficient maps, 2) Fast Bilateral Solver (FBS) [19] as a post-processing step, 3) our modified bilateral refinement in place of FBS, and 4) our refinement applied after multi-scale processing. The input is a linear sRGB image generated from a pseudo ground-truth denoised image in the S24 validation set [16]. We report PSNR between the photofinishing output (downsampled to one-quarter resolution) and the corresponding ground truth, along with the runtime overhead for each approach, measured on an Intel Core i7-14700K CPU and an NVIDIA GeForce RTX 4080 SUPER GPU.

the grid-prediction subnetwork (see Sec. C.5) to produce a coefficient grid. Bilateral slicing using the guidance map $\mathbf{G}_{\text{guide}(s)}$, after the aforementioned smoothing step, yields scale-specific coefficient maps, which are upsampled to full resolution when $s \neq 1.0$. The upsampled coefficient maps from all scales are then averaged to obtain the final output.

This strategy integrates information from coarse-to-fine representations, improving stability and reducing artifacts. After multi-scale processing, we apply our GPU-accelerated bilateral refinement (Sec. A) to further refine the predicted coefficient maps, using $\mathbf{I}^{\downarrow}_{\text{gain}}$ as the guide. These two steps together significantly mitigate halo artifacts in difficult scenes (see Fig. 11).

Applying these steps blindly was found to reduce accuracy, as they also affect the strength of the LTM. To quantify this effect, we report results in Table 3 using the photofinishing module. We use pseudo ground-truth denoised images at one-quarter resolution as input (after white balancing and color correction) and evaluate against the corresponding sRGB ground-truth images at the same resolution from the S24 test set [16]. We compare the accuracy of our trained photofinishing module (Style #0, the default style)

Table 3. Results of our method with and without multi-scale processing and refinement of local tone mapping in the photofinishing module. Input images are linear sRGB generated from pseudo ground-truth denoised images in the S24 test set [16], downsampled to one-quarter resolution. Outputs are compared against the ground truth at the same resolution. The best results are highlighted in **yellow**.

| Method | S24 Test Set 1/4 (LsRGB/sRGB) | |
|---|---|---|
| | PSNR ↑ | SSIM ↑ |
| Default (w/o multi-scale and refinement) | **27.49** | **0.939** |
| w/ multi-scale and refinement | 26.28 | 0.929 |

with and without the multi-scale processing and refinement designed to mitigate the rare halo artifacts observed in challenging scenes.

As shown, applying these steps indiscriminately degrades accuracy, so they remain user-controlled and can be enabled or disabled as needed.

## B.2. Misaligned and Incomplete Datasets

Training our method requires access to scene illuminant vectors and color correction matrices (CCMs) for the training images (information typically available in DNG files) along with denoised "ground-truth" images. The latter can be generated from DNG files using AI-based blind denoisers such as the one in Adobe Lightroom, as in the S24 dataset [16]. Missing DNG files (or the extracted data they contain) pose a challenge for training our method. In addition, if the ground-truth sRGB images are misaligned with the corresponding raw inputs, this introduces another source of difficulty.

The Zurich Raw-to-sRGB dataset [48] provides an ideal test case for studying these issues, since it lacks both DNG files and the required metadata, and its raw-sRGB pairs are not spatially aligned. We therefore expect a degradation in accuracy, as our design does not explicitly handle misaligned training pairs. Nevertheless, this scenario is valuable for illustrating how our framework can be adapted to operate with missing data and for analyzing the resulting impact on accuracy.

The Zurich dataset consists of raw images captured with a smartphone camera and their corresponding JPEGs from a DSLR camera, but the pairs are not strictly pixel-aligned due to lens distortion, hand-held capture, and other factors. Although such a dataset is useful for evaluating robustness, it does not reflect the typical scenario, as obtaining aligned paired datasets is feasible in practice, as demonstrated by the S24 dataset [16] and the MIT-Adobe FiveK dataset [26].

The remainder of this section describes how we address these missing elements to enable training of our framework and analyzes the effect of misalignment in the Zurich dataset on the resulting performance.

Figure 10. Comparison of our modified bilateral refinement with the original Fast Bilateral Solver (FBS) [19]. The proposed modified bilateral refinement achieves comparable or improved results while running efficiently on GPU. The reference image is from the S24 test set [16]; the high-resolution source was generated in Adobe Photoshop and downsampled to 100×75 to obtain the low-resolution input. Runtimes were measured on an Intel Core i7-14700K CPU and an NVIDIA GeForce RTX 4080 SUPER GPU (16 GB VRAM).



Figure 11. Halo artifact suppression using the proposed multi-scale (MS) processing and guided refinement of the predicted LTM coefficients. Images were captured with the S24 main camera.

**Pseudo Ground-Truth Denoised Images.** To generate pseudo ground-truth denoised images, we apply an inverse 2.2 gamma to the ground-truth sRGB image, producing a "linearized" estimate $\mathbf{I}_{\mathtt{lin}}$. The use of 2.2 gamma is a practical approximation; in reality, camera ISPs typically apply a sequence of non-linear operators that cannot be fully inverted by this simple correction [9]. We then compute a global non-linear color mapping (CM) function $f_{\mathtt{CM}} : \mathbb{R}^3 \to \mathbb{R}^3$ that maps $\mathbf{I}_{\mathtt{lin}}$ to the input raw domain by solving the following optimization problem:

$$\arg \min_{f_{\mathtt{CM}}} \left\| f_{\mathtt{CM}}(\mathbf{I}_{\mathtt{lin}}) - \mathbf{I}_{\mathtt{raw}} \right\|_2^2, \qquad (15)$$

where $f_{\mathtt{CM}}$ is computed via linear regression with a polynomial kernel expansion, after saturated pixels are discarded. The pseudo ground-truth denoised raw is then obtained as:

$$\mathbf{I}_{\mathtt{pseudo}} = f_{\mathtt{CM}}(\mathbf{I}_{\mathtt{lin}}). \qquad (16)$$

**Illuminant Estimation.** Since no illuminant metadata is provided in the Zurich dataset, we estimate the camera illuminant as the mean RGB of the demosaiced raw image, equivalent to applying the gray-world assumption [25]. While this provides only a rough estimate, the resulting error can be compensated for by the computed CCM described below.

**Color Correction.** The Zurich dataset also lacks CCMs, so we solve for a 3×3 matrix $\mathbf{C}$ that maps white-balanced raw colors $\mathbf{r}'$ (with the estimated gray-world illuminant) to the corresponding sRGB values $\mathbf{s}$:

$$\min_{\mathbf{C}} \left\| \mathbf{r}' \mathbf{C}^\top - \mathbf{s} \right\|_2^2, \qquad (17)$$

subject to each row of $\mathbf{C}$ summing to 1 (to approximately preserve intensity). We solve this constrained least-squares optimization with non-negativity bounds using sequential least squares programming (SLSQP). In practice, the computed CCM compensates for errors in the estimated illuminant to produce colors close to the target images.

After these steps, we obtain the main components required to train our framework: input raw images, pseudo ground-truth denoised images, illuminant vectors, and CCMs to map the denoised raw images into linear sRGB space before training the photofinishing module. This allows us to train the framework as described in the main paper (more details of the training are available in Sec. G). Specifically, we train the denoising network and
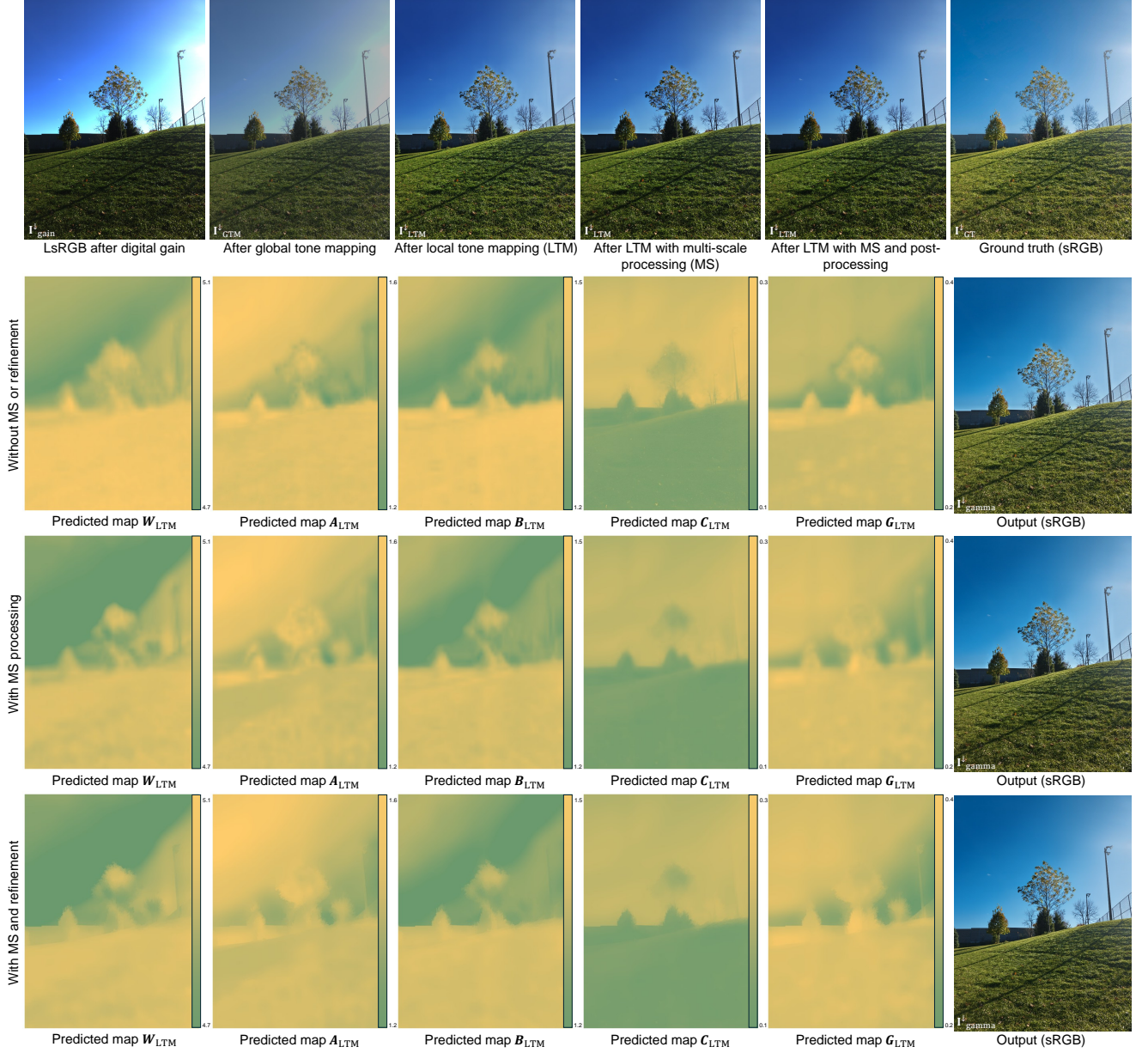
5

Figure 12. Comparison of predicted local tone mapping coefficients under different configurations: 1) without multi-scale processing, 2) with multi-scale processing, and 3) with multi-scale processing and guided refinement. For each configuration, we show intermediate outputs after LTM (first row) and the final sRGB outputs. The input images to the LTM module (linear sRGB after digital gain and after global tone mapping) and the ground truth are shown in the top row. The shown example is from the S24 test set [16].

the photofinishing module separately, then generate semi-final images before training the detail-enhancement network. Since the Zurich dataset provides $448 \times 448$ semi-aligned patches, we train the detail-enhancement network on patches of the same size, rather than the $512 \times 512$ patches used in the main paper experiments. Moreover, we disable the downsampling step before the photofinishing module and the guided upsampling after it when generating paired examples for the enhancement network and dur-

ing evaluation, since testing is also performed on $448 \times 448$ patches.

We report the results of our method in Table 4. As shown, our method does not achieve state-of-the-art results on the Zurich dataset, unlike on the S24 dataset [16]. Nevertheless, our workflow enables training on incomplete datasets such as Zurich and still yields competitive results, outperforming methods with a comparable number of parameters that lack the modularity of our framework (e.g., [10, 69]). Compared

6

Table 4. Results on the Zurich Raw-to-sRGB dataset [48]. We report PSNR, SSIM [81], LPIPS [89], and $\Delta$E 2000 [72], along with the total number of parameters for each method. While the Zurich dataset poses challenges for our framework due to misalignment and missing metadata, our method still achieves competitive results, with about a 1 dB drop in PSNR, while requiring significantly fewer parameters. The best results are highlighted in yellow.

| Method | Zurich Raw-to-sRGB Test Set | | | | |
| | PSNR↑ | SSIM↑ | LPIPS↓ | $\Delta$E 2000↓ | # params |
|---|---|---|---|---|---|
| PyNet [48] | 21.19 | 0.747 | 0.193 | NA | 47,548,170 |
| CIE XYZ Net [10] | 19.75 | 0.697 | 0.408 | 9.283 | 1,348,789 |
| LiteISP [91] | 21.55 | 0.749 | 0.187 | NA | 9,094,000 |
| FourierISP [44] | 21.65 | 0.755 | 0.182 | NA | 7,589,736 |
| LAN [69] | 19.46 | 0.730 | NA | NA | 46,847 |
| Ours (lite, w/o enhancement) | 19.57 | 0.717 | 0.401 | 9.448 | 452,447 |
| Ours (base, w/o enhancement) | 19.57 | 0.718 | 0.405 | 9.415 | 1,139,907 |
| Ours (large, w/o enhancement) | 19.73 | 0.721 | 0.397 | 9.257 | 3,841,547 |
| Ours (lite, w/ enhancement) | 20.68 | 0.726 | 0.390 | 8.306 | 503,082 |
| Ours (base, w/ enhancement) | 20.71 | 0.727 | 0.393 | 8.246 | 1,190,542 |
| Ours (large, w/ enhancement) | 20.76 | 0.729 | 0.386 | 8.221 | 3,892,182 |

to the best-performing method [44] in As shown in Table 4, our model requires substantially fewer parameters ($\sim$500 K with the lite denoising network and $\sim$3.9 M with the large variant; see Sec. C.1 for configurational details of the denoising network variants), while achieving less than a 1 dB drop in PSNR (compared to $\sim$7.6 M parameters for the reference method [44]). Moreover, the modular design of our framework provides greater flexibility, making the rendering pipeline easier to control, scale, debug, and customize.

## C. Design of Networks

This section details the architectures of the networks that form our proposed framework.

### C.1. Raw-Denoising Network

For the raw-denoising network ($\mathcal{D}_{\text{raw}}$), we employ three variants of the NAFNet architecture [30], which differ only in their base channel width and thus span different points on the accuracy-efficiency trade-off. The lite, base, and large variants differ in their initial channel width (4, 8, and 16 channels, respectively), which is doubled after each encoder stage. This scaling leads to progressively larger bottleneck widths (64, 128, and 256 channels, respectively) and total capacities of approximately 0.25 M, 0.93 M, and 3.6 M parameters. All three networks share the same overall design: a four-stage encoder–decoder with skip connections and a middle stage of four residual blocks. Each network operates in a residual-learning manner, predicting a correction that is added back to the input raw image. The encoder and decoder are composed of [2, 2, 4, 8] and [2, 2, 2, 2] NAF blocks per stage, respectively. Each block consists of lightweight NAF modules with depthwise convolutions, channel attention, and simple gating to realize an activation-free nonlinearity. Layer normalization and learnable scale parameters are applied to stabilize training. The networks

are fully convolutional and maintain compatibility with raw images of arbitrary resolution.

### C.2. Attention and Multi-Branch Blocks

In our photofinishing module, the networks employ a lightweight attention mechanism and multi-branch convolutional (MBConv) blocks. Specifically, we adopt the coordinate attention (CA) mechanism [45] with some modifications. Specifically, we replace Batch Normalization [50] with Group Normalization [82] to ensure consistent behavior under small batch sizes. We employ the LeakyReLU activation in all CA blocks. To prevent overly narrow bottlenecks in low-channel layers, we impose a lower bound of eight channels on the reduced dimensionality, computed as $\max(8, C_{\text{in-CA}}/r_{\text{CA}})$, where $C_{\text{in-CA}}$ is the input channel count and $r_{\text{CA}}$ is the reduction ratio. The reduction ratio is set to 4 by default, except for the luminance-guidance subnetwork within the chroma-mapping network, where we use $r_{\text{CA}}=2$.

The MBConv block is designed to capture features at multiple receptive-field scales while maintaining low computational cost. As shown in Fig. 13, the block consists of three depthwise convolution branches that operate in parallel. The first branch uses a 3×3 kernel without dilation to capture fine local details. The second branch also uses a 3×3 kernel but with dilation set to 2, allowing the convolution to gather broader contextual information without increasing parameters. The third branch employs a larger 5×5 kernel to further expand the receptive field and capture smoother structural variations. Each branch begins with reflection padding, followed by a depthwise convolution and a LeakyReLU activation function. The outputs from all branches are summed and then fused by a 1×1 convolution to produce the final aggregated feature map. This design effectively enhances receptive-field diversity without increasing the number of parameters or channels.

We provide ablation studies on the impact of using the MBConv blocks and the CA mechanism on the final photofinishing results in Sec. K.1.4.

### C.3. Gain and Gamma-Correction Networks

Both the digital-gain ($\mathcal{D}_{\text{gain}}$) and gamma-correction ($\mathcal{D}_{\text{gamma}}$) networks share the same lightweight convolutional architecture, designed to predict a single global scalar factor that adjusts exposure or that is used for gamma correction. The two networks differ only in their output ranges.

As illustrated in Fig. 14, the input image is first resized to a fixed resolution of 128×128 using bilinear interpolation. The network then begins with a 3×3 convolution layer with reflection padding, followed by Group Normalization (two groups) and a LeakyReLU activation. The features are processed by an MBConv block to capture multi-scale spatial context, followed by a CA block to encode long-range
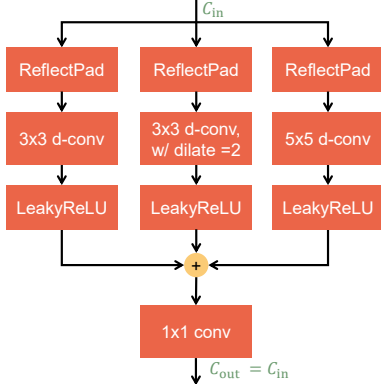
Figure 13. Structure of the multi-branch convolutional block (MB-Conv). The input feature map is processed in three parallel depth-wise convolution branches. Each branch applies reflection padding before the convolution and a LeakyReLU activation afterward. The outputs from all branches are summed and then fused through a convolution layer to produce the final aggregated feature map. This block preserves the number of channels (i.e., $C_{\text{in}} = C_{\text{out}}$).

dependencies along both spatial directions. Another $3\times3$ convolution and LeakyReLU activation refine the features before global pooling.

Two stages of average pooling are used to progressively compress spatial information: first to one-quarter of the input resolution (i.e., $128\times128\rightarrow32\times32$), and then to a single $1\times1$ feature vector. The resulting vector is passed through a fully connected layer with a Sigmoid activation to produce a normalized scalar value in the range $[0, 1]$. Reflection padding is applied to all convolutions to prevent border artifacts, and the total number of trainable parameters for each network (digital-gain and gamma-correction) is 6,587.

For the digital-gain stage, the predicted scalar is linearly mapped to the range $[0.25, 4.0]$ (equivalent to approximately exposure value [EV] $[-2, +2]$), producing a gain factor $d_g$ that scales image intensities before tone and color processing. For gamma correction, the normalized output is mapped to the range $[1.2, 3.0]$, yielding a gamma factor $\gamma$.

## C.4. Global Tone Mapping Network

The global tone-mapping (GTM) network ($\mathcal{D}_{\text{GTM}}$), consisting of 28,369 learnable parameters, predicts three positive parameters that define a parametric tone curve applied to the linear sRGB image after digital gain ($\mathbf{I}^{\downarrow}_{\text{gain}}$). As shown in Fig. 15, the input image is first resized to $128\times128$ using bilinear interpolation. The backbone begins with a $3\times3$ convolution (10 channels) with reflection padding, followed by Group Normalization (two groups) and a LeakyReLU activation. The resulting features are processed by an MB-Conv block and a CA block to enhance spatial context modeling. Two subsequent $3\times3$ convolutions (20 channels each, reflection padding) with LeakyReLU activations further refine the features. Spatial information is then pro-
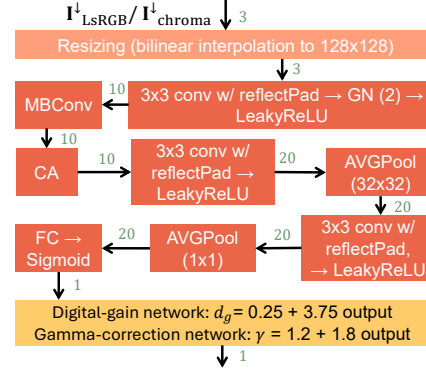


Figure 14. Architecture of the network used for both the digital gain and gamma correction (6,587 parameters). The input image is first resized to a fixed resolution $128\times128$, then processed by a convolutional layers with Group Normalization (2 groups) and LeakyReLU activation, followed by multi-branch convolutional and coordinate-attention blocks. After two stages of adaptive pooling, a fully connected layer with Sigmoid activation predicts a scalar in $[0, 1]$, which is linearly mapped to the target range for either digital gain ($[0.25, 4.0]$) or gamma ($[1.2, 3.0]$). The number of channels at each stage is indicated in green.

gressively aggregated through adaptive average pooling to $16\times16$, followed by a $3\times3$ convolution (40 channels, reflection padding) and LeakyReLU activation, another adaptive pooling to $4\times4$, and a final $3\times3$ convolution (40 channels, reflection padding) with LeakyReLU activation. A final global average pooling reduces the feature map to $1\times1$. The flattened feature vector is passed through a fully connected layer that produces three outputs, followed by a Softplus activation to ensure $(a_{\text{GTM}}, b_{\text{GTM}}, c_{\text{GTM}}) > 0$.

## C.5. Local Tone-Mapping Network

The LTM network, $\mathcal{D}_{\text{LTM}}$ (120,215 learnable parameters), predicts spatially varying coefficients that locally modulate the tone-mapping behavior applied to the linear sRGB image after digital gain ($\mathbf{I}^{\downarrow}_{\text{gain}}$) and global tone mapping ($\mathbf{I}^{\downarrow}_{\text{GTM}}$). Unlike the global tone-mapping network (Sec. C.4), which predicts a single set of global parameters ($a_{\text{GTM}}, b_{\text{GTM}},$ and $c_{\text{GTM}}$), the LTM network produces per-pixel coefficient maps ($\mathbf{A}_{\text{LTM}}, \mathbf{B}_{\text{LTM}}, \mathbf{C}_{\text{LTM}}, \mathbf{G}_{\text{LTM}},$ and $\mathbf{W}_{\text{LTM}}$) that enable locally adaptive tone mapping and exposure adjustment.

The LTM network consists of two main components: 1) a multi-scale guidance subnetwork and 2) a grid-prediction subnetwork, as shown in Fig. 16. The multi-scale guidance subnetwork derives a guidance map from one of the color channels of $\mathbf{I}^{\downarrow}_{\text{gain}}$. It processes three progressively downsampled versions of the luminance channel ($\times1$, $\times1/2$, and $\times1/4$ scales) using parallel convolutional branches, each composed of a $3\times3$ convolution (reflection padding), Group Normalization (2 groups), LeakyReLU activation, an MB-Conv block, and a CA block, followed by four additional
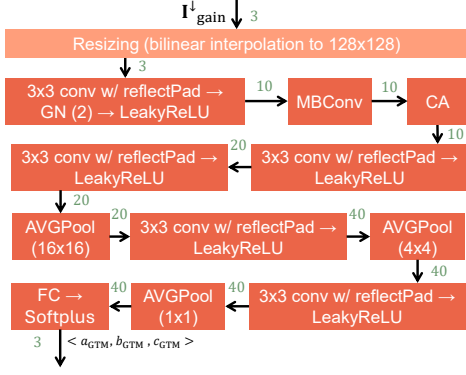
8

Figure 15. Architecture of the global tone-mapping (GTM) network (28,369 parameters). The input linear sRGB image after digital gain ($\mathbf{I}_{\texttt{gain}}^{\downarrow}$) is first resized to $128 \times 128$ and processed by convolutional layers with Group Normalization (2 groups) and LeakyReLU activations, followed by multi-branch convolutional and coordinate-attention blocks. Two additional convolutional layers refine the features before a series of adaptive pooling operations ($16 \times 16 \to 4 \times 4 \to 1 \times 1$). A fully connected layer with Softplus activation outputs the three positive parameters ($a_{\texttt{GTM}}$, $b_{\texttt{GTM}}$, $c_{\texttt{GTM}}$) that define the GTM curve. The number of channels at each stage is shown in green.

convolutional layers. The outputs from the three scales are upsampled (bilinear), concatenated, and fused using convolutional layers followed by a Tanh activation to produce the final guidance map $\mathbf{G}_{\texttt{guide}} \in [-1, 1]$. This multi-scale guidance design is intended to improve robustness against scale and contrast variations (see Sec. K.1.4 for ablation study).

In parallel, the grid-prediction subnetwork estimates a coarse grid of tone-mapping coefficients conditioned on both $\mathbf{I}_{\texttt{gain}}^{\downarrow}$ and $\mathbf{I}_{\texttt{GTM}}^{\downarrow}$. The concatenated image pair is downsampled to $384 \times 384$ and processed by a series of convolutional layers (with Group Normalization and LeakyReLU activations), MBConv and CA blocks, and average pooling to form a latent grid representation of size $N_g \times N_g$ with depth $N_c$. A final $1 \times 1$ convolution outputs $N_c \times 5$ feature channels corresponding to five coefficient volumes, activated by Softplus to ensure non-negativity.

The predicted $N_g \times N_g \times N_c \times 5$ coefficient grid, $\mathbf{M}_{\texttt{LTM}}$, is then sampled via trilinear interpolation using the guidance map $\mathbf{G}_{\texttt{guide}}$ as the depth coordinate, producing the spatial coefficient maps ($\mathbf{A}_{\texttt{LTM}}$, $\mathbf{B}_{\texttt{LTM}}$, $\mathbf{C}_{\texttt{LTM}}$, $\mathbf{G}_{\texttt{LTM}}$, and $\mathbf{W}_{\texttt{LTM}}$). In our implementation, we set $N_c=18$ and $N_g=64$, corresponding to a $64 \times 64$ spatial grid with 18 depth slices. See Sec. K.1.4 for ablation results on the impact of grid size.

During training, the spatial coefficient maps are regularized using the LTM smoothness loss, $\ell_{\texttt{LTM}-\texttt{s}}$ (see Sec. F). At inference time, an optional multi-scale processing and bilateral refinement step can be applied to smooth the coefficient maps and mitigate artifacts (see Sec. B.1 for details).
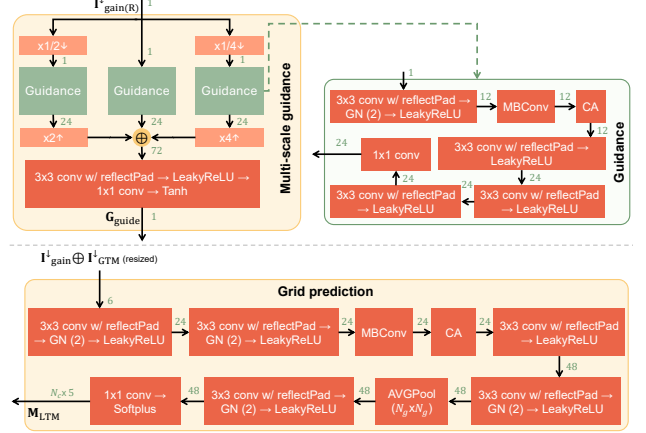


Figure 16. Architecture of the local tone-mapping (LTM) network (120,215 parameters). It consists of two main components: a multi-scale guidance subnetwork that processes a single channel of the linear sRGB image after digital gain to predict the guidance map $\mathbf{G}_{\texttt{guide}}$, and a grid-prediction subnetwork that processes the concatenated ($\oplus$) linear sRGB image after digital gain and the globally tone-mapped image to estimate the coefficient grid $\mathbf{M}_{\texttt{LTM}}$. The predicted grid is trilinearly sampled using $\mathbf{G}_{\texttt{guide}}$ to produce five spatial coefficient maps that control local tone-mapping behavior. The number of channels at each stage is indicated in green.

## C.6. Chroma-Mapping Network

The chroma-mapping network ($\mathcal{D}_{\texttt{chroma}}$) predicts an image-specific 2D lookup table (LuT) that transforms the chrominance components (CbCr) of the tone-mapped image produced by the preceding stages. The LuT is learned in an end-to-end fashion and applied via differentiable grid sampling to enable backpropagation through the photofinishing module. The network operates in the YCbCr color space, using the luminance channel of the tone-mapped image ($\mathbf{Y}_{\texttt{LTM}}$) as an auxiliary guidance signal.

Given the tone-mapped image (processed by both GTM and LTM operators) in YCbCr format, the image is first resized to a fixed spatial resolution of $128 \times 128$. A differentiable 2D histogram representation of the CbCr channels, $\hat{\mathbf{H}}^{CbCr}$, is then computed (see Sec. D for details). This histogram encodes the joint distribution of chroma values across $N_h \times N_h$ bins (with a value range of $[-0.5, 0.5]$) and provides a compact, differentiable summary of the chrominance content.

The histogram $\hat{\mathbf{H}}^{CbCr}$ is processed by a shallow convolutional subnetwork (hereafter referred to as the hist subnetwork) to extract chroma features (four channels). These features are concatenated with an identity 2D meshgrid, denoted as $\mathbf{H}_{\texttt{pos}} \in \mathbb{R}^{N_h \times N_h \times 2}$, which encodes the normalized CbCr coordinates and provides the convolutional layers with explicit knowledge of the histogram bin locations [11]. The resulting six-channel tensor is then used as input to an

9

encoder–decoder network with luminance-guided attention, as illustrated in Fig. 17.

The encoder consists of three convolutional stages. The first encoder stage applies a $3\times3$ convolution with reflection padding, Group Normalization (four groups), and a LeakyReLU activation. The second encoder stage includes a $3\times3$ convolution with reflection padding followed by a LeakyReLU activation and an MBConv block, while the third stage applies only a $3\times3$ convolution with reflection padding followed by a LeakyReLU. All stages preserve the channel dimensionality (28) and spatial resolution. The final encoder output is passed through a CA block to encode long-range dependencies along both chroma dimensions.

In parallel, the luminance channel ($\mathbf{Y}_{\text{LTM}}$) is processed by a lightweight auxiliary subnetwork that produces a 28-D attention vector. This luminance-guidance subnetwork begins with a $3\times3$ convolution (8 channels) followed by Group Normalization (two groups) with LeakyReLU activation, a CA block (reduction ratio 2), and an MBConv block. After adaptive average pooling to $8\times8$, the features are refined by another $3\times3$ convolution and LeakyReLU activation, followed by global pooling and a fully connected layer with Sigmoid activation. The resulting normalized vector ($[0, 1]^{28}$) modulates the encoder bottleneck features by channel-wise multiplication, allowing luminance-dependent adaptation of the predicted chroma mapping.

The decoder mirrors the encoder structure, consisting of three convolutional stages. The final layer uses a $1\times1$ convolution with two output channels and a Tanh activation to generate the residual chroma LuT ($N_h\times N_h\times2$). The predicted LuT ($\mathbf{L}_{\text{chroma-res}}$) is added to a learnable, image-independent base LuT ($\mathbf{L}_{\text{chroma-base}}$) to produce the final 2D LuT. The base LuT, $\mathbf{L}_{\text{chroma-base}}$, is initialized as an identity mapping in the CbCr space before training. During training, the final constructed LuT is regularized using the chroma LuT smoothness loss, $\ell_{\text{LuT-s}}$ (Sec. F). The total number of learnable parameters in the chroma-mapping network, including the base LuT, is 45,466.

## C.7. Detail-Enhancement Network

For the detail-enhancement network ($\mathcal{D}_{\text{enh}}$), we employ a lightweight variant of the NAFNet architecture [30], which predicts a residual correction added back to the input image. The detail-enhancement network follows the same encoder-decoder design with skip connections as our denoising models, but with a significantly reduced depth and width to minimize complexity. Specifically, it is configured with an initial channel width of 8, two encoder stages, two decoder stages, and four blocks in the middle stage. This compact setup results in approximately 50 K trainable parameters while providing sufficient capacity for fine-detail enhancement within our pipeline.
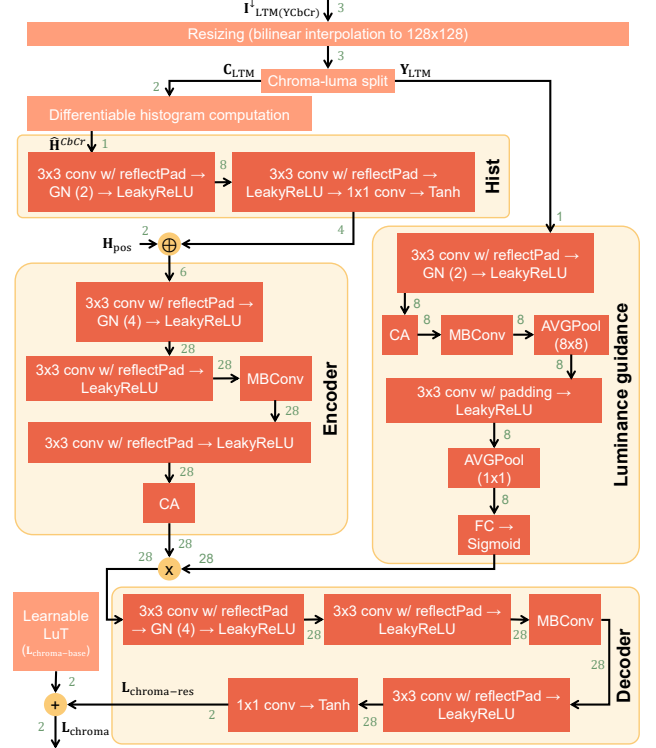


Figure 17. Architecture of the chroma-mapping network (45,466 parameters). The CbCr channels are first converted into a differentiable 2D histogram ($\hat{\mathbf{H}}^{CbCr}$; see Sec. D), which is processed by a shallow subnetwork (hist). The resulting chroma features are concatenated ($\oplus$) with an identity meshgrid ($\mathbf{H}_{\text{pos}}$) to provide subsequent layers with explicit awareness of the histogram bin positions. The encoder-decoder backbone (with MBConv and CA blocks) predicts a residual $N_h\times N_h\times2$ LuT that is then added to a learnable base LuT. A luminance-guided subnetwork processes the $\mathbf{Y}_{\text{LTM}}$ channel to produce an attention vector that modulates the encoder bottleneck features. The number of channels at each stage is indicated in green.

## D. Differentiable Histogram Computation

As described in the main paper and in this supplemental material (Sec. C.6), our chroma mapping network ($\mathcal{D}_{\text{chroma}}$) relies on a differentiable histogram representation of the chrominance channels to enable end-to-end training of the photofinishing module. Given an input image with Cb and Cr channels, we construct a 2D histogram by softly assigning each pixel to its corresponding histogram bins [7, 12].

Let $N_h$ denote the number of bins in the 2D histogram (and, accordingly, in the chroma 2D LuT constructed by the chroma-mapping network), and let $[v_{\min}, v_{\max}]$ denote the value range. We place $N_h$ uniformly spaced bin centers for both Cb and Cr channels as follows:

$$c_i = v_{\min} + \frac{i}{N_h-1}(v_{\max} - v_{\min}),$$
$$i = 0, \ldots, N_h - 1. \tag{18}$$

Each pixel $(cb, cr)$ contributes to all bins with a Gaussian weighting:

$$\mathcal{W}_{ij}(cb, cr) = \exp\left(-\frac{(cb - c_i)^2 + (cr - c_j)^2}{2\sigma_{\text{hist}}^2}\right), \quad (19)$$

where $cb$ and $cr$ denote the chrominance values of a pixel, and $c_i, c_j$ are the uniformly spaced bin centers along the Cb and Cr axes, respectively. The term $\sigma_{\text{hist}}$ controls the softness of the bin assignment. This formulation follows the differentiable histogram approach in [7, 12]. The histogram is computed by summing over all pixels:

$$\mathbf{H}_{ij}^{CbCr} = \sum_p \mathcal{W}_{ij}(cb_p, cr_p), \quad (20)$$

then normalized to unit sum and square-rooted:

$$\hat{\mathbf{H}}^{CbCr} = \sqrt{\frac{\mathbf{H}^{CbCr}}{\sum_{i,j} \mathbf{H}_{ij}^{CbCr} + \epsilon}}. \quad (21)$$

In our implementation, we set the default values to $N_h{=}24$ bins (see Sec. K.1.4 for ablation results with $N_h{=}12$), $v_{\min}{=}-0.5$, $v_{\max}{=}0.5$, and $\sigma_{\text{hist}}{=}0.075$.

## E. Guided Upsampling Regularization

In the main paper, we explained that our method employs bilateral guided upsampling (BGU) [29] with a modified regularization that yields better results. Here, we provide additional details of our regularization. Our photofinishing module predicts a low-resolution gamma-corrected output, $\mathbf{I}_{\text{gamma}}^{\downarrow} \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times 3}$. To recover full resolution, we apply BGU using the full-resolution linear sRGB guide image $\mathbf{I}_{\text{LsRGB}} \in \mathbb{R}^{H \times W \times 3}$. The high-resolution photofinishing result (after upsampling) is denoted as $\mathbf{I}_{\text{gamma}}^{\uparrow} \in \mathbb{R}^{H \times W \times 3}$. The fitting grid is constructed in the bilateral space using the low-resolution guide $\mathbf{I}_{\text{LsRGB}}^{\downarrow} \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times 3}$, while the local affine fitting statistics are accumulated from the corresponding low-resolution input-output pair $(\mathbf{I}_{\text{LsRGB}}^{\downarrow}, \mathbf{I}_{\text{gamma}}^{\downarrow})$, with guidance derived from the high-resolution image $\mathbf{I}_{\text{LsRGB}}$ during upsampling. The grid size is set to $\frac{H}{64} \times \frac{W}{64} \times 16$ to achieve a good balance between speed and accuracy.

Let $\mathbf{S}_m \in \mathbb{R}^{(C+1) \times (C+1)}$ and $\mathbf{T}_m \in \mathbb{R}^{C \times (C+1)}$ denote the accumulated statistics in grid cell $m$, with $c_m$ samples. For RGB image, $C = 3$. The affine transform for cell $m$, $\mathbf{A}_m \in \mathbb{R}^{C \times (C+1)}$, is estimated by solving:

$$\begin{aligned}\mathbf{T}_m \mathbf{S}_m^{\top} + \lambda_m \mathbf{R}_m = \mathbf{A}_m(\mathbf{S}_m \mathbf{S}_m^{\top} + \lambda_m \mathbf{I}_{C+1}), \\ \lambda_m = \lambda(c_m + 1),\end{aligned} \quad (22)$$

where $\mathbf{I}_{C+1}$ is the $(C+1) \times (C+1)$ identity matrix and $\lambda$ is a scalar scaling factor. Intuitively, $\mathbf{S}_m$ and $\mathbf{T}_m$ encode the

relationship between the low-resolution input $\mathbf{I}_{\text{gamma}}^{\downarrow}$ and the guide $\mathbf{I}_{\text{LsRGB}}$ within grid cell $m$, such that $\mathbf{A}_m$ models a local affine mapping that transfers structure and edges from the guide to the upsampled output. The regularization matrix $\mathbf{R}_m$ controls the per-channel stability of this fitting. Thus, the design of $\mathbf{R}_m$ directly affects how well color and fine details are preserved in the final $\mathbf{I}_{\text{gamma}}^{\uparrow}$ result.

**Original Halide Regularization.** The Halide implementation of BGU [29] tries to solve:

$$\mathbf{T}_m^{\text{blur}} \mathbf{S}_m^{\text{blur}\top} + \lambda_m \mathbf{R}_m^{\text{luma}} = \mathbf{A}_m(\mathbf{S}_m^{\text{blur}} \mathbf{S}_m^{\text{blur}\top} + \lambda_m \mathbf{I}_{C+1}), \quad (23)$$

where

$$\mathbf{R}_m^{\text{luma}} = \begin{bmatrix} \text{diag}(r_m^{\text{luma}}, r_m^{\text{luma}}, r_m^{\text{luma}}) & \mathbf{0} \end{bmatrix}, \quad (24)$$

that uses a single luma gain $r_m^{\text{luma}} \in \mathbb{R}$ applied uniformly to all RGB channels:

$$\begin{aligned} r_m^{\text{luma}} &= \frac{[0.25, 0.5, 0.25] \cdot \mathbf{T}_m^{\text{blur}}[1{:}C, C{+}1] + \lambda^{\text{glob, blur}}}{[0.25, 0.5, 0.25] \cdot \mathbf{S}_m^{\text{blur}}[1{:}C, C{+}1] + \lambda^{\text{glob, blur}}}, \\ \lambda^{\text{glob, blur}} &= \lambda\left(\sum_n c_n^{\text{blur}} + 1\right), \end{aligned}$$
$$\quad (25)$$

where $\mathbf{T}_m^{\text{blur}}$ and $\mathbf{S}_m^{\text{blur}}$ are the blurred version of $\mathbf{T}_m$ and $\mathbf{S}_m$ respectively, and $c_m^{\text{blur}}$ is the number of samples in the blurred grid $m$.

Intuitively, this achromatic assumption causes color crosstalk. Although empty cells are handled via spatial grid blurring, this process can oversmooth edges and fine details.

**Our Gated Regularization.** Instead of applying grid blurring that may create artifacts, we adopt the original BGU [29] objective (Eq. 22) and use per-channel adaptive gains:

$$\mathbf{R}_m^{\text{gate}} = \begin{bmatrix} \text{diag}(\mathbf{r}_m) & \mathbf{0} \end{bmatrix}, \quad (26)$$

where $\mathbf{r}_m = [r_{m,1}, r_{m,2}, r_{m,3}]$ is defined as:

$$r_{m,c} = \begin{cases} \dfrac{\mathbf{T}_m[c, C+1]}{\mathbf{S}_m[c, C+1] + \lambda_m}, & \text{if } c_m > 0, \\ r_c^{\text{glob}}, & \text{if } c_m = 0, \end{cases} \quad (27)$$

with global per-channel fallback gains:

$$\begin{aligned} r_c^{\text{glob}} &= \frac{\sum_n \mathbf{T}_n[c, C+1]}{\sum_n \mathbf{S}_n[c, C+1] + \lambda^{\text{glob}}}, \quad (28) \\ \lambda^{\text{glob}} &= \lambda\left(\sum_n c_n + 1\right). \quad (29) \end{aligned}$$

11

This gated regularization applies accurate local per-channel constraints when sufficient samples exist in a grid cell, while sparsely populated cells are robustly handled using global statistics. Unlike the original Halide scheme, our regularization eliminates the need for grid blurring, avoids color crosstalk, and preserves sharp details in the upsampled result (see Sec. K.1.2 for ablation studies).

## F. Loss Function Details

We provide the implementation details of the loss terms used to train our photofinishing module. Each loss is designed to complement a specific stage of the photofinishing module, ensuring perceptually faithful reproduction of tone, color, and structure. The overall objective is given in Eq. 9 of the main paper.

**SSIM Loss.** We adopt a differentiable variant of the structural similarity index [81], computed independently per color channel. Given our output and the ground-truth sRGB images $\mathbf{I}_{\text{gamma}}^{\downarrow}$ and $\mathbf{I}_{\text{GT}}^{\downarrow}$, we compute local statistics over an $11 \times 11$ Gaussian window (standard deviation $\sigma=1$) using depthwise convolution:

$$\mu_1 = \mathbf{I}_{\text{gamma}}^{\downarrow} * \mathbf{w}, \quad \mu_2 = \mathbf{I}_{\text{GT}}^{\downarrow} * \mathbf{w},$$
$$\sigma_1^2 = (\mathbf{I}_{\text{gamma}}^{\downarrow}{}^2 * \mathbf{w}) - \mu_1^2, \quad \sigma_2^2 = (\mathbf{I}_{\text{GT}}^{\downarrow}{}^2 * \mathbf{w}) - \mu_2^2,$$
$$\sigma_{12} = (\mathbf{I}_{\text{gamma}}^{\downarrow} \mathbf{I}_{\text{GT}}^{\downarrow} * \mathbf{w}) - \mu_1 \mu_2.$$
(30)

In Eq. 30, $\mathbf{w}$ denotes the Gaussian weighting kernel. The SSIM map is then computed as:

$$\text{SM}(\mathbf{I}_{\text{gamma}}^{\downarrow}, \mathbf{I}_{\text{GT}}^{\downarrow}) = \frac{(2\mu_1\mu_2 + C_1)(2\sigma_{12} + C_2)}{(\mu_1^2 + \mu_2^2 + C_1)(\sigma_1^2 + \sigma_2^2 + C_2)}, \quad (31)$$

where $C_1 = 0.01^2$ and $C_2 = 0.03^2$. The final loss is defined as $\ell_{\text{SSIM}} = 1 - \text{mean}(\text{SM})$.

**Perceptual Loss.** We compute the perceptual loss using a pretrained VGG-19 network [73] with frozen weights. We extract features from layers relu3_3 and relu4_2, following the convention of perceptual supervision in photo enhancement tasks. Both predicted and ground-truth RGB images are resized to $224 \times 224$ and normalized by the ImageNet mean and standard deviation. The loss is given by the following equation:

$$\ell_{\text{perc}} = \sum_{l \in \{\text{relu3\_3}, \text{relu4\_2}\}} \|\phi_l(\mathbf{I}_{\text{gamma}}^{\downarrow}) - \phi_l(\mathbf{I}_{\text{GT}}^{\downarrow})\|_1, \quad (32)$$

where $\phi_l(\cdot)$ denotes the activation map of layer $l$.

$\Delta E$ **Color Loss.** To enforce perceptual color accuracy, we compute a differentiable approximation of the CIE 1976 $\Delta E_{ab}^*$ (i.e., $\Delta E_{76}$) metric [71] in the CIE Lab color space. Both the predicted and ground-truth images are evaluated before gamma correction. For the ground-truth, we first remove the sRGB gamma encoding using the gamma scale, $\gamma$, predicted by our gamma-correction network ($\mathcal{D}_{\text{gamma}}$), so that both images are evaluated under a consistent, image-specific gamma setting. The resulting pre-gamma, quasi-linear sRGB images are then converted to the CIE Lab color space through a differentiable transformation chain:

$$\text{linear sRGB} \rightarrow \text{XYZ} \rightarrow \text{Lab}.$$

The XYZ conversion uses the standard transformation matrix:

$$\mathbf{M}_{\text{LsRGB} \rightarrow \text{XYZ}} = \begin{bmatrix} 0.4124564 & 0.3575761 & 0.1804375 \\ 0.2126729 & 0.7151522 & 0.0721750 \\ 0.0193339 & 0.1191920 & 0.9503041 \end{bmatrix}.$$
(33)

For the XYZ → Lab mapping, the standard CIE Lab definition uses a piecewise cubic-root function as follows:

$$f_{\text{cubic}}(t) = \begin{cases} t^{1/3}, & t > \delta^3, \\ \frac{t}{3\delta^2} + \frac{4}{29}, & \text{otherwise.} \end{cases}$$
(34)

where $\delta = 6/29$. Since this function is non-differentiable at $t = \delta^3$, we replace it with a smooth blend defined as:

$$f_{\text{soft}}(t) = \sigma_s\big(\alpha_s(t - \delta^3)\big) t^{1/3} + \big[1 - \sigma_s\big(\alpha_s(t - \delta^3)\big)\big] \left(\frac{t}{3\delta^2} + \frac{4}{29}\right),$$
(35)

where $\sigma_s(\cdot)$ denotes the sigmoid function used in our differentiable Lab conversion, and $\alpha_s = 150$ controls the sharpness of the transition. The resulting Lab tensors, denoted as $\mathbf{L}^{\text{out}}$ and $\mathbf{L}^{\text{GT}}$, represent the full three-channel $(L^*, a^*, b^*)$ values for our output and the ground-truth images, respectively. The per-pixel $\Delta E_{76}$ distance is then computed as:

$$\ell_{\Delta E} = \frac{1}{N} \sum_p \big\|\mathbf{L}_p^{\text{out}} - \mathbf{L}_p^{\text{GT}}\big\|_2, \quad (36)$$

where $N = H \times W$ is the total number of pixels. If the 3D LuT (pre-chroma mapping network) is included as a learnable component of the photofinishing module, we additionally incorporate an auxiliary color-difference term equal to half the $\Delta E_{76}$ between the de-gammaed ground-truth image and the 3D LuT output. This auxiliary term encourages the 3D LuT to approximate the target color rendering prior to subsequent chroma mapping.

**Chroma Loss.** This loss enforces chroma consistency by comparing the CbCr channels of the predicted output (after chroma mapping) with those of the de-gammaed ground truth, where the de-gammaing is performed using the image-specific gamma factor $\gamma$ predicted by our gamma-correction network. The chroma loss, $\ell_{\text{CbCr}}$, is defined as:

$$\ell_{\text{CbCr}} = \left\| \hat{\mathbf{C}}_{\text{LuT}} - \mathbf{C}_{\text{GT}} \right\|_1, \tag{37}$$

where $\mathbf{C}_{\text{GT}}$ denotes the CbCr channels of the de-gammaed ground-truth image $(\mathbf{I}_{\text{GT}}^{\downarrow}{}^{\gamma})$, and $\hat{\mathbf{C}}_{\text{LuT}}$ represents the mapped CbCr channels obtained using the predicted 2D LuT $\mathbf{L}_{\text{chroma}}$.

**Tone-Mapping Loss.** The tone-mapping loss enforces joint consistency between the global and local tone-mapping modules and the ground-truth tone-mapped image. Let $\mathbf{Y}_{\text{GT}}$ denote the ground-truth luminance (Y) channel, de-gammaed using our image-specific gamma factor $(\gamma)$. Similarly, let $\mathbf{Y}_{\text{GTM}}$ and $\mathbf{Y}_{\text{LTM}}$ denote the luminance channels obtained after the global and local tone-mapping stages, respectively. The tone-mapping loss is defined as:

$$\ell_{\text{TM}} = 0.6 \left\| \mathbf{Y}_{\text{GTM}}^{\downarrow 8} - \mathbf{Y}_{\text{GT}}^{\downarrow 8} \right\|_1 + \left\| \mathbf{Y}_{\text{LTM}} - \mathbf{Y}_{\text{GT}} \right\|_1, \tag{38}$$

where $\downarrow 8$ indicates $1/8$ downsampling via bilinear interpolation. The downsampled term penalizes deviations in the global tone and overall luminance of the tone-mapped output produced by the GTM network, while the second term encourages the LTM network to refine local tonal and luminance variations to better align with the ground truth.

**Luminance Consistency Loss.** We regularize the tone-mapping process to preserve the average image brightness:

$$\ell_{\text{luma}} = \left| \mathbb{E}[\mathbf{Y}_{\text{GTM}}] - \mathbb{E}[\mathbf{Y}_{\text{gain}}] \right|, \tag{39}$$

where $\mathbf{Y}_{\text{gain}}$ denotes the Y channel of the linear sRGB image after applying the predicted image-specific digital gain $d_g$ in its YCbCr representation. This loss prevents exposure drift and encourages the tone-mapping modules to refine contrast while preserving global brightness.

**Total Variation Regularization Loss.** The $\ell_{\text{LuT-s}}$ and $\ell_{\text{LTM-s}}$ terms promote spatial coherence in the chroma LuT and LTM maps. We implement these terms as isotropic total variation (TV) losses to encourage smoothness in the learned chroma LuT and local tone-mapping coefficient maps:

$$\ell_{\text{TV}}(\mathbf{X}) = \tfrac{1}{2} \left( \|\nabla_h \mathbf{X}\|_1 + \|\nabla_w \mathbf{X}\|_1 \right), \tag{40}$$

where $\nabla_h$ and $\nabla_w$ denote horizontal and vertical spatial gradients, respectively.

**Overall Loss Combination.** The total training objective is defined as a weighted sum of all loss components:

$$\mathcal{L}_{\text{total}} = \sum_j \lambda_j \, \ell_j, \tag{41}$$

where each weighting coefficient $\lambda_j$ balances the relative contribution of its corresponding loss term. Empirically, we found that combining low-level ($\ell_1$, $\ell_{\text{SSIM}}$, $\ell_{\text{CbCr}}$), perceptual ($\ell_{\Delta E}$, $\ell_{\text{perc}}$), and regularization ($\ell_{\text{LuT-s}}$, $\ell_{\text{LTM-s}}$, $\ell_{\text{luma}}$, $\ell_{\text{TM}}$) losses yields stable convergence and perceptually consistent results across different picture styles (see Sec. K.1.3 for ablation studies).

## G. Training Details

As mentioned in the main paper, we train each module separately. In this section, we provide further details of training our framework.

### G.1. Denoising Training

We trained the denoising model ($\mathcal{D}_{\text{raw}}$) to minimize the $\ell_1$ loss between the pseudo denoising ground truth on $256 \times 256$ patches and the model outputs. The network was optimized using AdamW [64] with $\beta_1 = 0.9$, $\beta_2 = 0.9$, and no weight decay. The learning rate was initialized at $10^{-3}$ and decayed to $10^{-5}$ using cosine annealing [63]. Training was performed for 100 epochs with a batch size of 32, and gradient clipping with a maximum norm of 0.01 was applied to stabilize optimization. For each dataset, including the generic denoisers (Sec. H), we trained three model variants—lite, base, and large. The architectural configurations of these variants are provided in Sec. C.1.

### G.2. Photofinishing Training

We trained the photofinishing module ($\mathcal{D}_{\text{gain}}$, $\mathcal{D}_{\text{GTM}}$, $\mathcal{D}_{\text{LTM}}$, $\mathcal{D}_{\text{chroma}}$, $\mathcal{D}_{\text{gamma}}$, $\mathbf{L}_{\text{chroma-base}}$, and optionally the 3D LuT, $\mathbf{L}_{\text{RGB}}$) for 600 epochs with a batch size of 8. The input images were the pseudo ground-truth denoised raw images mapped to linear sRGB space ($\mathbf{I}_{\text{LsRGB}}$), and the corresponding ground-truth images were the final sRGB targets ($\mathbf{I}_{\text{GT}}$). Both input and ground-truth images were resized to $512 \times 512$, with standard geometric augmentations applied during training.

Optimization was performed using Adam [54] with parameters $\beta_1 = 0.9$, $\beta_2 = 0.999$, and an $\ell_2$ regularization term of $10^{-7}$. The learning rate was initialized at $10^{-4}$ and decayed to $10^{-6}$ using a cosine annealing schedule [63]. The loss function combined multiple terms as described in the main paper and in Sec. F. The weights of these losses were empirically set as follows: $\lambda_1 = 2.5$ for $\ell_1$, $\lambda_{\text{SSIM}} = 0.5$ for $\ell_{\text{SSIM}}$, $\lambda_{\Delta E} = 0.02$ for $\ell_{\Delta E}$, $\lambda_{\text{perc}} = 0.01$ for the perceptual loss $\ell_{\text{perc}}$, $\lambda_{\text{CbCr}} = 1.0$ for the CbCr loss $\ell_{\text{CbCr}}$, $\lambda_{\text{LuT-s}} = 0.06$ for the LuT smoothness loss $\ell_{\text{LuT-s}}$,

$\lambda_{\text{TM}} = 0.5$ for the tone mapping loss $\ell_{\text{TM}}$, $\lambda_{\text{LTM-s}} = 0.6$ for the LTM smoothness loss $\ell_{\text{LTM-s}}$, and $\lambda_{\text{luma}} = 0.2$ for the luminance consistency loss $\ell_{\text{luma}}$. Note that in the loss function details section (Sec. F), the resized photofinishing output and ground truth are denoted as $\mathbf{I}^{\downarrow}_{\text{gamma}}$ (after applying the predicted gamma correction) and $\mathbf{I}^{\downarrow}_{\text{GT}}$, respectively, at a resolution of $512 \times 512$. In all other contexts, $\mathbf{I}^{\downarrow}_{\text{gamma}}$ and $\mathbf{I}^{\downarrow}_{\text{GT}}$ refer to the photofinishing output and ground-truth images at the native $1/4$ raw resolution.

### G.3. Detail-Enhancement Training

After training the denoising models and the photofinishing module, we trained our final detail-enhancement network ($\mathcal{D}_{\text{enh}}$), which refines the perceptual quality of the final output. To construct the training data, we first applied our trained denoisers (lite, base, and large variants, used jointly to enlarge the training set as an augmentation strategy). The denoised raw images were converted to linear sRGB using the metadata of each image (illuminant color vector and CCM), downsampled to one quarter of the resolution, processed by the trained photofinishing module, and then upsampled using the bilateral guided upsampling method [29] with our regularization (Sec. E). We extracted non-overlapping $512 \times 512$ patches from the guided-upsampled results and their corresponding ground-truth images. We used normalized image values in $[0, 1]$ (without quantization) to construct the paired training set.

For the Zurich dataset (Sec. B.2), we did not apply the downsampling and guided-upsampling procedure since the dataset provides only cropped $448 \times 448$ patches, and thus we trained directly on these patches.

Training was performed for 50 epochs with a batch size of 16. The network was optimized using Adam [54] with parameters $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\ell_2$ regularization $10^{-7}$. The learning rate was initialized at $10^{-4}$ and decayed to $10^{-6}$ using a cosine annealing schedule [63]. The pixel-wise $\ell_1$ loss between predicted and ground-truth sRGB patches was used as the loss function.

### H. Generalization Across Cameras

One of the main advantages of our modular design is the ability to decouple the camera-specific modules of our pipeline from the generic ones that can operate across cameras. Specifically, the raw image denoising module is one of the primary camera-specific components, as noise characteristics differ from camera to camera. Other camera-specific factors include color correction (i.e., auto white balancing and color correction matrices), which can be addressed using recent cross-camera illuminant estimation methods (e.g., [11]) or by reading metadata from DNG files—see Sec. I.3 for details on how we handle cross-camera auto white-balance correction.

To overcome the limitations of camera-dependent denoising, we trained a generic denoiser in addition to our camera-specific models trained for the S24 main camera, as reported in the main paper. To train the generic denoiser, we first modeled the noise characteristics of several cameras and then synthesized noise using these models on the pseudo ground-truth denoised raw images from the S24 dataset [16] to generate noisy raw inputs with synthetic noise from different cameras. These synthetic noisy raw images, along with the original noisy raw images from S24, were used to train our generic denoiser models under the same configurations described in Sec. G.1. Specifically, we constructed a hybrid S24 dataset comprising 20% clean pseudo ground-truth denoised images, 40% real noisy S24 images, and 40% synthetic noisy images generated using the aforementioned cross-camera noise models.

For each synthetic noisy sample, the ISO value was inherited from its corresponding real image to maintain realistic noise-exposure relationships. We adopted the widely used heteroscedastic Gaussian noise model [37, 38, 61]:

$$\mathbf{n}_{\text{gd}} \sim \mathcal{N}\left(\mathbf{0}, \sigma^2_{\text{gd}}\right), \quad \sigma^2_{\text{gd}} = \beta_{\text{gd},1}\,\mathbf{x}_{\text{gd}} + \beta_{\text{gd},2}, \quad (42)$$

where the noise variance $\sigma^2_{\text{gd}}$ increases linearly with the underlying clean pixel intensity $\mathbf{x}_{\text{gd}}$. The parameters $\beta_{\text{gd},1}$ and $\beta_{\text{gd},2}$ represent the shot- and read-noise components, respectively. To synthesize noise under this model, we collected $(\mathbf{x}_{\text{gd}}, \sigma^2_{\text{gd}})$ pairs and calibrated $\beta_{\text{gd},1}$ and $\beta_{\text{gd},2}$ for each ISO and color channel via linear regression. The shot-noise parameter $\beta_{\text{gd},1}$ corresponds to the slope, while the read-noise parameter $\beta_{\text{gd},2}$ corresponds to the intercept of the linear fit. During calibration, we performed 14-18 captures of an X-Rite ColorChecker chart (24 color patches). Each capture was taken at a different ISO setting, ranging from ISO 50 to ISO 3200. For each ISO, we extracted the 24 color patches and subdivided them into smaller sub-patches to increase the number of samples for fitting. This process yielded multiple small patches of size ($35 \times 35$) per ISO and color channel, with the total number varying depending on the chart's coverage within the image.

Next, we computed the mean and variance of each small patch. Using these means as $\mathbf{x}_{\text{gd}}$ and their corresponding variances as $\sigma^2_{\text{gd}}$, we fitted a linear regression model to estimate the slope ($\beta_{\text{gd},1}$) and intercept ($\beta_{\text{gd},2}$). This produced one pair of $\beta_{\text{gd},1}$ and $\beta_{\text{gd},2}$ values for each ISO in the calibration data. For ISOs not present, we interpolated the existing parameters (using linear interpolation for $\beta_{\text{gd},1}$ and cubic interpolation for $\beta_{\text{gd},2}$) to enable noise synthesis for all intermediate ISOs within the calibration range.

To synthesize diverse noise patterns during training, we used noise models calibrated for a total of 12 sensors: S20 FE (main), S20 Ultra (main), Pixel 6 (main), Pixel 9 Pro (main), Pixel 9 Pro (telephoto), S22 Plus (front), S22 Plus (telephoto), S22 Plus (ultra-wide), S22 Plus (main), S22

Table 5. Comparison of camera-specific and generic denoising models (lite, base, and large) on the S24 [16] and MIT Adobe 5K [26] test sets (noisy/denoised raw). The best results are highlighted in **yellow**.

| Denoising Model | S24 Test Set (Noisy/Denoised) | | Adobe 5K Test Set (Noisy/Denoised) | |
|---|---|---|---|---|
| | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ |
| Camera-specific (lite) | 53.39 | 0.998 | 55.57 | 0.998 |
| Camera-specific (base) | 55.97 | **0.999** | 57.36 | **0.999** |
| Camera-specific (large) | **57.33** | **0.999** | **59.85** | **0.999** |
| Generic (lite) | 51.19 | 0.997 | 51.84 | 0.996 |
| Generic (base) | 55.15 | **0.999** | 52.86 | 0.997 |
| Generic (large) | 56.46 | **0.999** | 53.07 | 0.997 |

Table 6. Comparison of denoising models (lite, base, large) trained on S24 [16], MIT Adobe 5K [26], and on a wide range of noise profiles for generic denoising, evaluated on the SSID dataset [1]. The best results are highlighted in **yellow**.

| Denoising Model | SSID Test Set (Raw Noisy/Raw Clean) | |
|---|---|---|
| | PSNR ↑ | SSIM ↑ |
| S24 (lite) | 46.53 | 0.940 |
| S24 (base) | 47.28 | 0.952 |
| S24 (large) | 47.47 | 0.956 |
| Adobe (lite) | 48.05 | 0.966 |
| Adobe (base) | 48.28 | 0.967 |
| Adobe (large) | 48.12 | 0.971 |
| Generic (lite) | 49.09 | 0.971 |
| Generic (base) | **50.00** | 0.978 |
| Generic (large) | 49.55 | **0.979** |

(ultra-telephoto), S22 (ultra-wide), and S22 Ultra (main). These sensors were selected to represent a diverse range of sensor types (main, telephoto, ultra-wide, and front), enabling the generic denoiser to learn from heterogeneous noise distributions.

Similar to the camera-specific denoiser, we trained three variants of the generic denoiser (lite, base, and large). Table 5 compares the accuracy of camera-specific models—each trained and evaluated on its corresponding dataset (S24 or MIT Adobe 5K [26])—against the generic denoisers evaluated on the same datasets. The denoised "ground-truth" images for Adobe 5K were obtained by running Adobe Lightroom's AI denoiser on the raw images, following the same procedure used for the S24 dataset. As expected, camera-specific models performed better on their corresponding cameras, while the generic models were designed to generalize to unseen devices. To further validate this, we evaluated both model types on the Smartphone Image Denoising Dataset (SIDD) [1], which contains raw images from smartphones not included in the S24, Adobe 5K, or generic model training sets. Results in Table 6 show that the generic models outperformed the camera-specific ones, demonstrating stronger cross-camera generalization.

Figure 18 shows qualitative results of our entire pipeline on unseen cameras (Samsung Galaxy S9 and iPhone X)

from the Raw-to-Raw dataset [6], using the photofinishing module trained on the S24 dataset. We report results with both the S24-trained denoiser and a generic denoiser, combined with either the camera AWB or the cross-camera AWB model (C5) [11], as well as with the post-AWB user-preference mapping [92] (see Sec. I.3 for more details). For comparison, we include outputs from Adobe Lightroom, both with and without its built-in auto enhancement and AI-based denoiser features. Our method achieves high-quality cross-camera rendering, competitive with commercial software. For the iPhone X, we observed that raw images are substantially darker than those from most other cameras; thus, we applied the auto-exposure adjustment before the digital gain step (see Sec. I.2 for details). This adjustment was facilitated by the modularity of our pipeline, as further discussed in Sec. I.

Figure 19 shows a qualitative comparison between our method (using the generic denoiser along with the photofinishing and detail-enhancement modules trained on the S24 dataset [16]) and the LiteISP [91] and ISPDiffuser [70] models, both trained on the same dataset. The comparison is conducted on images captured by various iPhone devices, none of which were included in the training of any model, including ours. The results clearly demonstrate the superior cross-camera generalization ability of our approach, producing consistent and visually pleasing results with greater controllability (e.g., the ability to adjust white balance), while requiring significantly fewer parameters. Specifically, LiteISP contains approximately 9 M parameters and ISPDiffuser around 21 M, whereas our lite version–which was used in this qualitative comparison–requires only about 0.5 M parameters.

Additional results on cameras not seen during training are provided in Figs. 1 and 5 of the main paper, as well as in Fig. 7 of this supplemental material. Further qualitative examples are shown in Figs. 20, 22, 31, 32, and 49 in the following sections.

## I. User Interface for Modular ISP Control

As described in the main paper, we implemented a graphical user interface (GUI) to enable intuitive interaction with the modular components of our pipeline. Our photo-editing tool allows users to control all stages of the pipeline by enabling or disabling modules, adjusting their strengths, selecting from different picture styles, and even blending between styles or intermediate stages. In addition, it provides optional operators (e.g., contrast, highlights, shadows, and others) that can be applied within the pipeline, offering full flexibility to adjust the look and feel of the rendered images. Users can also export their manual edits as new styles and apply them to multiple images in batches (see Fig. 20). Figure 21 illustrates where these additional operators are integrated within the main pipeline described in the paper.
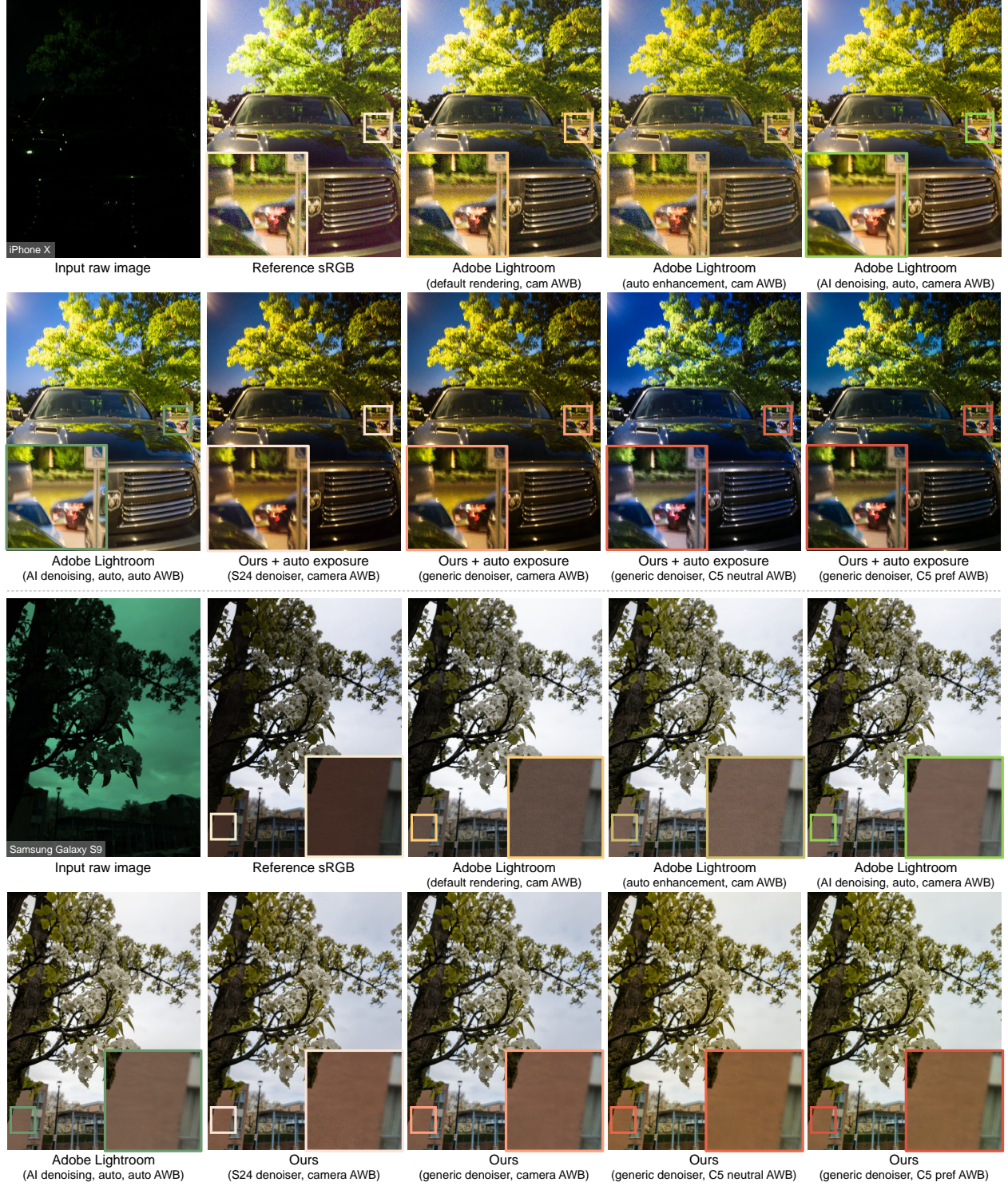
Figure 18. Qualitative results on unseen cameras from the Raw-to-Raw dataset [6]. The first row shows an iPhone X example and the second a Samsung Galaxy S9 example. For each case, we include the input raw image, the reference sRGB rendered by the capturing app, and Adobe Lightroom outputs (default, built-in auto enhancement with/without AI denoising, and results with camera AWB and Lightroom auto AWB). Our outputs are shown with two denoisers (trained on the S24 dataset [16] and a generic one), using camera AWB, C5 (neutral) AWB [11], and user-preference mapping [92]. All results are generated with the photofinishing module trained on the S24 dataset.
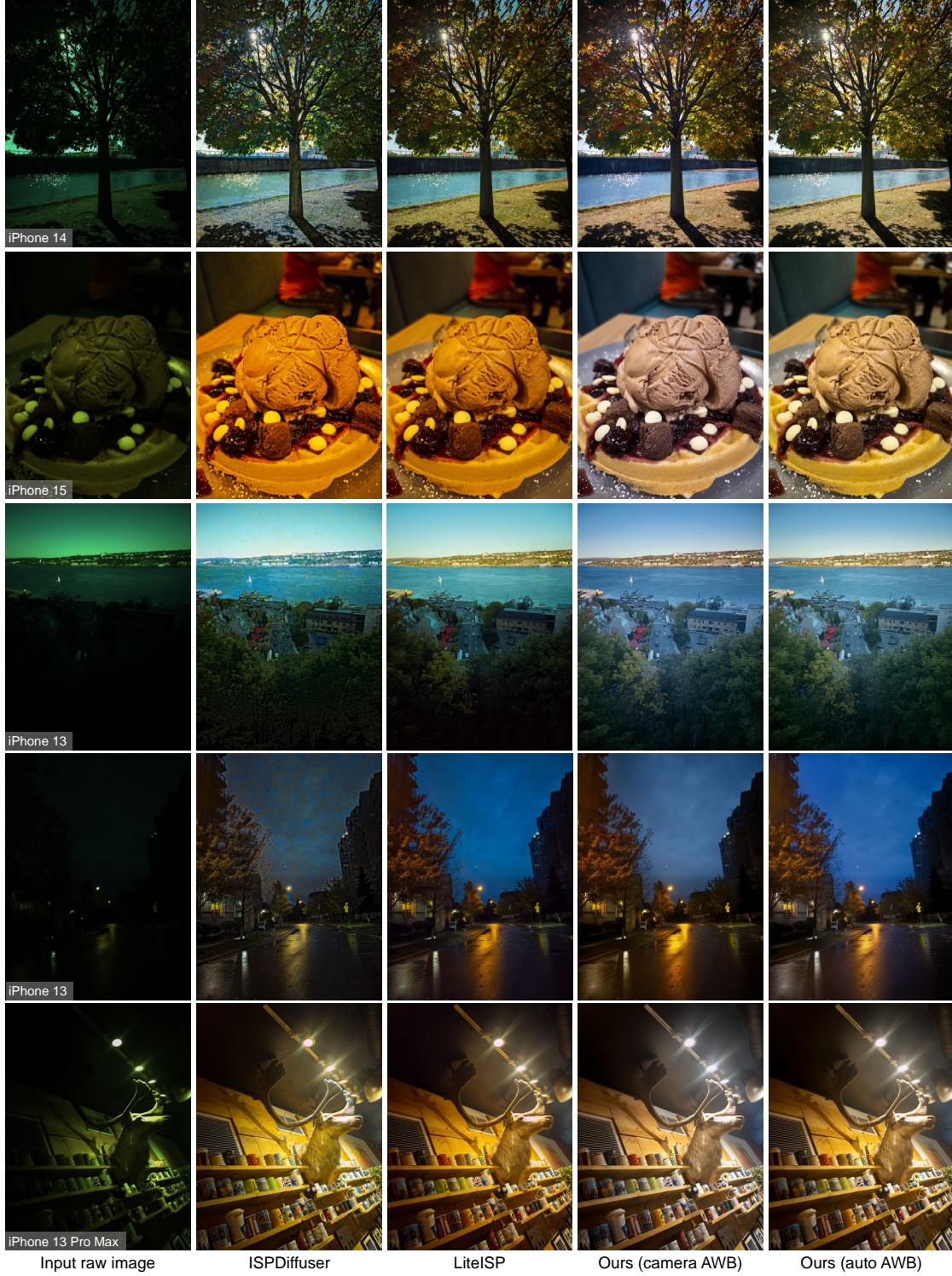
16

Figure 19. Qualitative comparison on unseen cameras (iPhone devices) between our method and the LiteISP [91] and ISPDiffuser [70] models, all trained on the S24 dataset [16]. Our method uses the generic denoiser, and none of the training data included iPhone captures. In this example, we show our results under the camera's automatic white balance (AWB) and after applying learned AWB models [11, 92] (see Sec. I.3 for details).
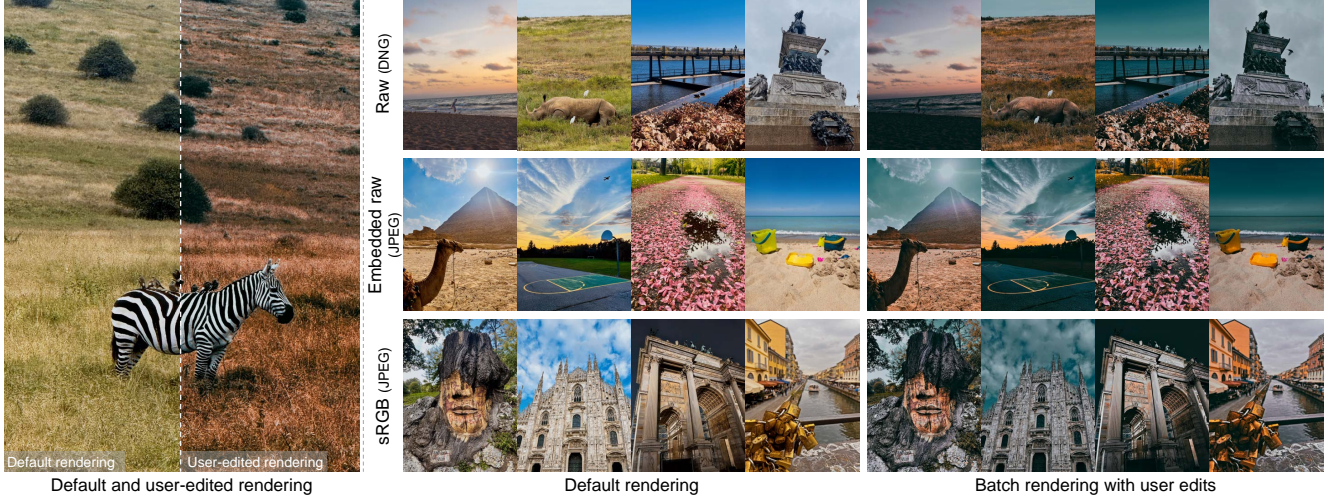
Figure 20. Our photo-editing tool offers full control over the modular neural ISP pipeline and includes several interactive editing operators. Our tool supports DNG raw inputs from any camera, as well as rendered sRGB images previously processed by our tool, where compressed raw data is embedded within the JPEG for accurate re-rendering. For sRGB inputs from other cameras, the tool synthesizes raw-like data on the fly for rendering. Users can save custom edits and batch-apply them to multiple images. The shown reference image was captured with a Canon EOS 90D, while other raw and rendered images (with embedded raw) were taken using iPhone 13 and Canon EOS 90D; all sRGB images were captured by the iPhone 13 main camera.
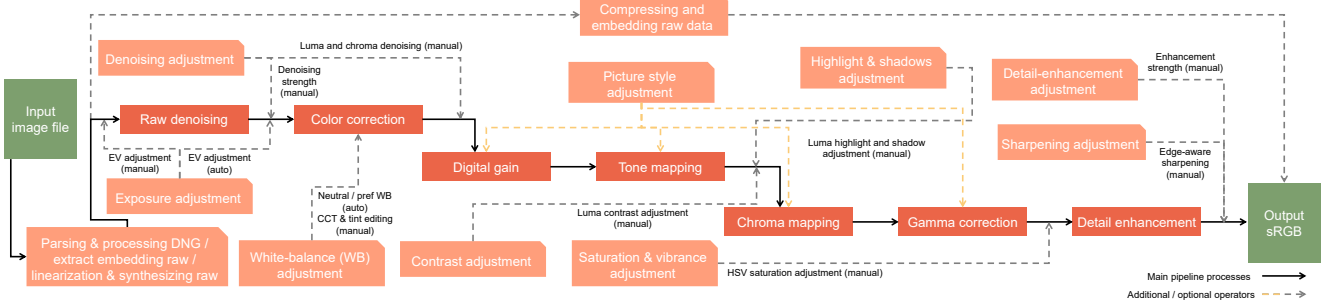


Figure 21. Overview of the operators available in our graphical user interface. In addition to the main pipeline processes, the tool provides optional user-controlled operators (e.g., exposure, contrast, highlights, shadows, saturation, and sharpening), allowing fine-grained control over the final rendered image.

In our implementation, the tool uses the photofinishing models trained on the S24 dataset [16] for the default picture style, as we found this model produces more visually plausible results (see Fig. 22). For the LTM, we adopt the proposed multi-scale processing and refinement (Sec. B.1) as the default configuration, while still allowing the user to disable this option if desired. In addition to the default picture style, we employ models trained on the S24 dataset that support five additional artistic picture styles (Style #1–5). In our GUI, we refer to them as `Warm`, `Moody`, `Cinematic`, `Greenish`, and `Retro`, respectively. For denoising, we rely on the S24-specific denoising base model for images captured by the S24 main camera, and a generic base denoising model to support images captured by other cameras. For white balancing (Sec. I.3), the tool includes two models: a camera-specific model for the S24 main camera and a cross-camera model to support unseen cameras.

Our tool accepts DNG raw images, where we first apply black-level normalization and then perform Menon demosaicing [65] using a parallel, tile-based strategy with partial overlaps between tiles. In practice, we divide the Bayer image into overlapping tiles of $512 \times 512$ pixels with a 16-pixel overlap to avoid boundary artifacts. Each tile is demosaiced independently using multiple threads and seamlessly blended into the final full-resolution RGB raw image. This tiled parallel processing improves runtime efficiency (up to a $5\times$ speed-up on 12-megapixel images compared to the standard full-frame implementation).

The tool further supports post-editable image re-rendering, where compressed raw data is embedded into the final rendered JPEG image (Sec. I.10). This design allows recovering the raw data for unlimited post-editable re-rendering under different settings, always starting from the raw image and thus avoiding accumulated degrada-

Input raw image     Ours (trained on MIT-Adobe 5K)     Ours (trained on S24)

Figure 22. Output of our method using photofinishing and detail-enhancement models trained on the S24 dataset [16] (default picture style; Style #0) and the MIT Adobe 5K dataset [26] (Expert C style; see Sec. K.2.2 for details). The image was captured with the iPhone 14 main camera. In both cases, our generic base denoising model was used.

tions. Before saving, the raw data is pre-processed by a learned compression model to improve efficiency and keep the JPEG file size compact.

To broaden applicability, the tool can also process sRGB images not originally rendered with our system (and therefore lacking embedded raw data). For this case, we include a learning-based linearization module that maps camera-rendered sRGB images into a camera-agnostic space and synthesizes a raw-like representation, allowing the tool to handle sRGB input images produced outside our pipeline (Sec. I.11). This capability to process sRGB images from third parties, together with the batch-processing option, enables the tool to handle extracted sRGB video frames in a frame-by-frame manner (see Fig. 23). To improve temporal consistency across video frames, we apply a simple temporal smoothing over the current rendering parameters using information from the previous frames within a 9-frame window. More advanced temporal smoothing approaches are beyond the scope of this paper and are left for future work.

In total, our tool requires loading 3,902,604 parameters (approximately 14.9 MB / 0.015 GB) into GPU memory. This includes the denoising models (camera-specific and generic), photofinishing models for six picture styles (including the default), white balance models (cross-camera and camera-specific), a single detail-enhancement model (we use the one trained for Style #0 for simplicity, instead of loading style-specific models), the raw pre-processing model (for raw embedding), and the linearization model. As a result, the system remains practical to deploy without heavy GPU memory requirements—especially compared to other methods such as ISPDiffuser [70], which alone requires 20,938,890 parameters to render raw to a single picture style (6× for six styles), or LiteISP [91], which requires 9,094,000 parameters (6× for six styles). The remainder of this section describes the functionalities of our tool.

## I.1. Denoising Adjustment

The tool provides flexible control over the raw denoising strength. This control is implemented as a linear interpolation between the image before and after the denoising operation, modulated by a user-specified scalar parameter in the range $[0, 1]$. A value of $0$ preserves the unprocessed raw image, a value of $1$ applies the full denoising operation, and intermediate values yield proportionally mixed results.

In addition to this raw-domain adjustment, we introduce two denoising operators that act in the color-corrected domain—specifically on the linear sRGB image—after the color correction stage. These two operators are applied in the YCbCr representation of the linear sRGB image (where, for simplicity and similar to the photofinishing module, we use the standard RGB-to-YCbCr conversion matrix, which is originally defined assuming the input is in the non-linear sRGB space).

**Luma Denoising.** To reduce luminance noise while preserving structural edges, we apply an edge-preserving guided filter [43] to the luma channel, $\mathbf{Y}_{\text{LsRGB}}$, controlled by a radius $r_{\text{y-d}}$ and a regularization term $\epsilon_{\text{y-d}}$. The output $\mathbf{Y}'_{\text{LsRGB}}$ is combined with the original $\mathbf{Y}_{\text{LsRGB}}$ through:

$$\mathbf{Y}_{\text{y-d}} = (1 - \lambda_{\text{y-d}})\,\mathbf{Y}_{\text{LsRGB}} + \lambda_{\text{y-d}}\,\mathbf{Y}'_{\text{LsRGB}}, \qquad (43)$$

where $\lambda_{\text{y-d}} \in [0, 1]$ determines the luma denoising strength. Both $r_{\text{y-d}}$ and $\epsilon_{\text{y-d}}$ scale proportionally with $\lambda_{\text{y-d}}$, allowing a smooth transition from minimal smoothing at low strengths to stronger spatial denoising as $\lambda_{\text{y-d}}$ increases. This preserves local contrast and fine structure while progressively suppressing luminance grain. In our implementation, we used $r_{\text{y-d}} = 2 + 10\,\lambda_{\text{y-d}}$ and $\epsilon_{\text{y-d}} = (0.001 + 0.03\,\lambda_{\text{y-d}})^2$, which provide a balanced trade-off between structure preservation and noise reduction.

**Chroma Denoising.** For the chroma components, we apply spatial Gaussian smoothing to suppress color blotches while retaining perceptual color consistency. The filtered channels $\mathbf{Cb}'_{\text{LsRGB}}$ and $\mathbf{Cr}'_{\text{LsRGB}}$ are defined as

$$\begin{aligned}\mathbf{Cb}'_{\text{LsRGB}} &= \mathcal{G}(\mathbf{Cb}_{\text{LsRGB}};\,\sigma_{\text{c-d}}),\\ \mathbf{Cr}'_{\text{LsRGB}} &= \mathcal{G}(\mathbf{Cr}_{\text{LsRGB}};\,\sigma_{\text{c-d}}),\end{aligned} \qquad (44)$$

where $\mathcal{G}(\cdot;\sigma_{\text{c-d}})$ denotes a separable Gaussian convolution with standard deviation $\sigma_{\text{c-d}}$. The blur scale $\sigma_{\text{c-d}}$ is adaptively determined based on the user-defined chroma denoising strength $\lambda_{\text{c-d}}$ and the image resolution:

$$\sigma_{\text{c-d}} = \left(3 + 12\,\lambda_{\text{c-d}}\right)\left(\frac{H}{3000}\right)^{0.9}, \qquad (45)$$

where $H$ refers to the image height. This adaptive scaling maintains consistent perceptual smoothing across varying image resolutions. For high chroma denoising strength

19

Figure 23. Our tool can process sRGB input images (or video frames) and perform in-batch rendering with specific picture styles and edits. Shown are example frames from sRGB video sequences processed by our method, where each row corresponds to a distinct picture style and edits. See the supplemental video (click to view).
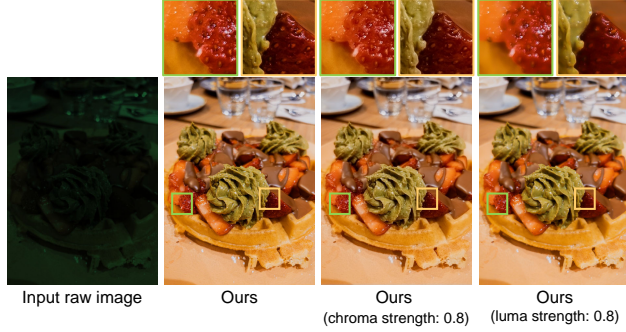


Figure 24. Example of luma and chroma denoising. The shown image was captured using a Samsung S24 main camera and rendered with our tool. We show our output without luma/chroma denoising and with both enabled at 80% strength.

$(\lambda_{\text{c-d}} > 0.4)$, an additional Gaussian pass is applied to further enhance color uniformity in low-texture regions. The final chroma outputs are blended with the original channels using a non-linear mixing rule:

$$
\begin{aligned}
\mathbf{Cb}_{\text{c-d}} &= (1 - \lambda_{\text{c-d}}^{1.5})\, \mathbf{Cb}_{\text{LsRGB}} + \lambda_{\text{c-d}}^{1.5}\, \mathbf{Cb}'_{\text{LsRGB}}, \\
\mathbf{Cr}_{\text{c-d}} &= (1 - \lambda_{\text{c-d}}^{1.5})\, \mathbf{Cr}_{\text{LsRGB}} + \lambda_{\text{c-d}}^{1.5}\, \mathbf{Cr}'_{\text{LsRGB}}.
\end{aligned}
\tag{46}
$$

The denoised luma and chroma channels are merged to form the denoised YCbCr version of the linear sRGB image, which is then transformed back to the linear sRGB space to be processed by the next procedure in our pipeline. Figure 24 shows an example.

## I.2. Exposure Adjustment

In addition to manual exposure adjustment (defined by an exposure value, EV, applied directly to the raw image and computed as $2^{\text{EV}}$), we provide an optional auto-exposure (AE) adjustment, implemented as a digital brightness compensation applied after conversion to linear sRGB and before the digital gain in the first stage of our photofinishing module. This step is particularly useful in cross-camera scenarios, where some devices produce raw images that are considerably darker than those encountered during training. Without such correction, these inputs may be misinterpreted as near-black, leading to insufficient boosting by the photofinishing network.

Our AE approach follows a histogram-based formulation inspired by prior work on exposure correction and tone reproduction [21, 75]. To reduce computation and suppress noise, we first downsample the input to $128 \times 128$ pixels using area pooling. From the downsampled linear sRGB image, we compute a luminance representation $\mathbf{Y}_E \in [0, 1]^{128 \times 128}$ and construct a 1D histogram with 96 bins. We then define a Gaussian target histogram over bin centers $b \in [0, 1]$ as follows:

$$
H_t^E(b) \propto \exp\left( -\frac{1}{2} \left( \frac{b - g_E}{\sigma_E} \right)^2 \right),
\tag{47}
$$

centered at a mid-gray value $g_E = 0.08$ with a fixed spread $\sigma_E = 0.05$. This target distribution encodes the desired tonal balance, with mid-tones dominating while highlights and shadows are preserved. To determine the exposure correction, we search over a discrete set of EV candidates

20

$\Delta \in [-E, +E]$ with $E=1.8$, corresponding to multiplicative scalings of $2^{\Delta}$. For each candidate, the luminance matrix $\mathbf{Y}_E$ is scaled, a histogram is recomputed, and the $\ell_2$ distance to the target distribution is evaluated. The EV that minimizes this distance is selected:

$$\Delta_E^* = \arg \min_{\Delta \in [-E, +E]} \left\| H^E(2^{\Delta} \mathbf{Y}_E) - H_t^E \right\|_2^2. \quad (48)$$

The final output is obtained by scaling the input image by $2^{\Delta_E^*}$.

We apply the EV correction in linear sRGB space, since this stage precedes the non-linear operations in the photofinishing module. Applying exposure adjustment after non-linear mapping can limit the ability to adjust exposure effectively [13]. Figure 25 illustrates a comparison between applying our AE in linear sRGB versus applying it at the end of the pipeline.

## I.3. White Balance Estimation and Adjustment

In all of our main experiments, we relied on the camera's auto white balance (AWB)—the scene illuminant color produced by the on-device AWB estimator embedded in the conventional ISP. Nevertheless, our tool also offers the ability to compute AWB using state-of-the-art illuminant estimation methods. For images captured by the Samsung S24 main camera, we adopted the camera-specific illuminant estimator of [16], which we re-trained after removing time and location features since location information is not always available in DNG files. For images captured by other cameras, we employed the cross-camera illuminant estimator C5 [11]. Following [52], we used a modified C5 model that required only the testing image, without the need for additional images from the same camera as in the original C5 method. Specifically, we employed a single-encoder variant that received a $48\times48$ histogram as input. This model was trained on the NUS dataset [31], the Cube++ dataset [36], and the S24 dataset. We followed the training instructions in the original works for both the camera-specific and cross-camera estimators.

Both AWB models (camera-specific and cross-camera) target a neutral white balance, where the goal is for achromatic surfaces to appear gray. However, such neutral estimates do not always match human perceptual preference [8], where a bias is often desirable. To account for this, we optionally integrate a post-processing step that applies a non-linear mapping from the neutral AWB estimate to a perceptual preference-aware AWB, following [92].

Our tool also allows manual adjustment of the correlated color temperature (CCT) and tint. For this, we rely on Robertson's method [84] to map between CIE XYZ, CCT, and tint coordinates, combined with an interpolation of the Planckian locus. In this context, the CCT corresponds to the point on the Planckian locus that is closest to the illuminant chromaticity, while the tint represents the signed perpendicular offset of the illuminant from this locus in the CIE 1960 $(u, v)$ chromaticity diagram. By parameterizing AWB in terms of both CCT and tint, we enable intuitive user control: CCT adjusts the overall warmth or coolness of the image, while tint provides fine-grained correction of residual color casts that cannot be addressed by CCT alone.

Finally, for the case where the AWB gains are taken from the camera (i.e., the ISP-provided AWB stored in the DNG files), we use the interpolated color correction matrix (CCM) already provided in the DNG files. For all other AWB settings (whether estimated using the camera-specific model or the cross-camera C5 model, in either neutral or preference-biased form, or manually adjusted via CCT and/or tint), we re-compute the CCM that maps the white-balanced raw RGB to linear sRGB. This is achieved by interpolating between the pre-calibrated CCMs available in the DNG metadata, following the DNG specification [5]. Figure 26 shows an example output from our tool, illustrating different AWB renderings.

## I.4. Contrast Adjustment

After applying local tone mapping, we optionally include a contrast adjustment step that can be controlled by the user before the chroma mapping stage in our photofinishing module. This operation modifies the luminance channel $\mathbf{Y}_{\text{LTM}}$ around mid-gray (0.5) to either increase or decrease the perceived image contrast. Specifically, given a user-specified adjustment factor $\alpha_{\text{contrast}} \in [-1, 1]$, we compute:

$$\hat{\mathbf{Y}}_{\text{LTM}} = \left(\mathbf{Y}_{\text{LTM}} - 0.5\right)\left(1 + 0.5 \ \alpha_{\text{contrast}}\right) + 0.5. \quad (49)$$

Positive values of $\alpha_{\text{contrast}}$ increase image contrast by expanding the difference between dark and bright regions, while negative values reduce contrast by compressing luminance values toward mid-gray. This adjustment allows users to fine-tune the global contrast of the image in addition to the automatic tone mapping provided by the pipeline. We find that applying contrast adjustment in the luminance channel at this stage leads to more visually pleasing results compared to applying it directly to the three channels of the processed linear sRGB either after local tone mapping or after the full photofinishing module (see Fig. 27).

## I.5. Highlights and Shadows Adjustment

We allow the user to optionally adjust highlights and shadows after the local tone-mapping stage. Both operations are applied to the luminance channel $(\mathbf{Y}_{\text{LTM}})$ of the YCbCr local tone-mapped image.

We first construct smooth masks that isolate highlight and shadow regions:

$$\mathbf{M}_{\text{high}} = \varsigma(\mathbf{Y}_{\text{LTM}}; \tau_{\text{high}}, \tau_{\text{high}} + \epsilon_{\text{high}}), \quad (50)$$

$$\mathbf{M}_{\text{shad}} = 1 - \varsigma(\mathbf{Y}_{\text{LTM}}; \tau_{\text{shad}}, \tau_{\text{shad}} + \epsilon_{\text{shad}}), \quad (51)$$

| Input raw image | Ours | Ours + auto exposure in sRGB | Ours + auto exposure in LsRGB | Ground truth |

Figure 25. We allow applying auto exposure (AE) after conversion to linear sRGB and directly before the photofinishing module. This figure shows a qualitative example from the S24 validation set [16], including the input raw image, our output without AE, and with AE applied either at the end of the pipeline or in linear sRGB. The ground-truth sRGB image from the S24 dataset is shown for reference.
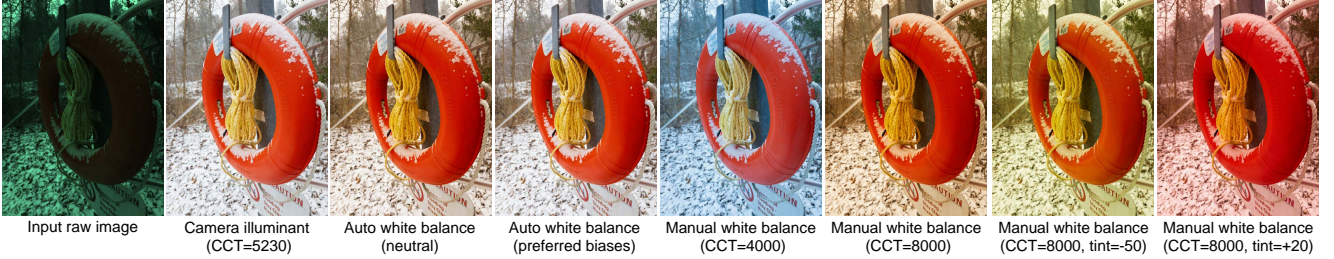


| Input raw image | Camera illuminant (CCT=5230) | Auto white balance (neutral) | Auto white balance (preferred biases) | Manual white balance (CCT=4000) | Manual white balance (CCT=8000) | Manual white balance (CCT=8000, tint=-50) | Manual white balance (CCT=8000, tint=+20) |

Figure 26. Our tool provides different white-balance options. The image can be rendered using the camera's AWB (referred to as "as shot" in our tool), or we can re-estimate the illuminant color using a learning-based auto white-balance model. For the latter, we provide two modes: a neutral white-balance correction that enforces achromatic objects to appear gray, and a preference-aware AWB that applies a bias to better match human perceptual preference. In addition, the user may manually adjust the CCT and tint. Example shown from the test set of the S24 dataset [16].

where $\tau_{\mathrm{high}}$ and $\tau_{\mathrm{shad}}$ are thresholds for highlight and shadow regions, $\epsilon_{\mathrm{high}}$ and $\epsilon_{\mathrm{shad}}$ control the transition smoothness, and $\varsigma$ denotes the smoothstep interpolation defined as:

$$\varsigma(x; e_0, e_1) = \left(\frac{x - e_0}{e_1 - e_0}\right)^2, \qquad (52)$$

with $e_0$ and $e_1$ denoting the lower and upper transition edges.

To prevent over-amplification near saturation, adjustments are limited to the valid luminance range [0.1, 0.9]. In our implementation, we set $\tau_{\mathrm{high}} = 0.7$, $\tau_{\mathrm{shad}} = 0.3$, and $\epsilon_{\mathrm{high}} = \epsilon_{\mathrm{shad}} = 0.1$. The adjusted luminance channels are then computed as follows:

$$\mathbf{Y}'_{\mathrm{high}} = \mathbf{Y}_{\mathrm{LTM}} + \frac{\alpha_{\mathrm{high}}}{20} \ \mathbf{Y}_{\mathrm{LTM}} \ \mathbf{M}_{\mathrm{high}}, \qquad (53)$$

$$\mathbf{Y}'_{\mathrm{shad}} = \mathbf{Y}_{\mathrm{LTM}} + \frac{\alpha_{\mathrm{shad}}}{20} \ (1 - \mathbf{Y}_{\mathrm{LTM}}) \ \mathbf{M}_{\mathrm{shad}}, \qquad (54)$$

where $\alpha_{\mathrm{high}}, \alpha_{\mathrm{shad}} \in [-1, 1]$ are user-controlled parameters. Positive $\alpha_{\mathrm{high}}$ values boost highlights, while negative values compress highlights to recover detail. Positive $\alpha_{\mathrm{shad}}$ values lift dark regions, whereas negative values deepen them.

Finally, the adjusted luminance is combined with the original chroma channels to reconstruct the full YCbCr image, which is then propagated through the remainder of the pipeline. This provides intuitive, user-controlled tonal adjustments while preserving color fidelity. See Fig. 28 for visual examples.

### I.6. Saturation and Vibrance Adjustment

We optionally allow user-controlled color saturation and vibrance adjustments in the HSV color space after gamma correction in our photofinishing module. Working in the HSV color space is numerically stable because its saturation channel is bounded in $[0, 1]$, and it provides an intuitive axis for perceptual editing.

For each pixel, given its HSV representation $(p_h, p_s, p_v)$, the saturation adjustment rescales the saturation channel uniformly across all colors as:

$$p'_s = \min\left(p_s \ (1 + \alpha_{\mathrm{sat}}), 1\right), \qquad (55)$$

where $\alpha_{\mathrm{sat}}$ is the saturation control parameter. Positive values of $\alpha_{\mathrm{sat}}$ increase global saturation, while negative values move the image toward grayscale.

| Input raw image | Ours (contrast: 0.0) | Ours (contrast: -1.0 in LsRGB) | Ours (contrast: - 1.0 in sRGB) | Ours (contrast: -1.0) |

Figure 27. We implement a user-controlled contrast adjustment on the luminance channel after local tone mapping. This figure shows a qualitative example from the S24 validation set [16], including the input raw image, our output without adjustment ($\alpha_{\text{contrast}} = 0.0$), and with reduced contrast ($\alpha_{\text{contrast}} = -1.0$). For comparison, results are shown when the adjustment is applied after local tone mapping in linear sRGB, after the full photofinishing module, and to the Y channel after local tone mapping ($\mathbf{Y}_{\text{LTM}}$) in YCbCr (our choice).



| Input raw image | highlights: 0.0, shadows: 0.0 | highlights: 0.75, shadows: 0.0 | highlights: -0.75, shadows: 0.0 | highlights: 0.0, shadows: 0.75 |
| highlights: 0.0, shadows: -0.75 | highlights: 0.75, shadows: 0.75 | highlights: -0.75, shadows: 0.75 | highlights: 0.75, shadows: -0.75 | highlights: -0.75, shadows: -0.75 |

Figure 28. We allow the user to optionally adjust highlights and shadows after applying local tone mapping. Shown here is an example from the S24 validation set [16], including the input raw image and our outputs with different values of the highlight and shadow control parameters, $\alpha_{\text{high}}$ and $\alpha_{\text{shad}}$, respectively.

In contrast, vibrance selectively boosts muted colors while preserving already saturated regions. The vibrance transformation is defined as:

$$p'_s = \min\left(p_s \ (1 + \alpha_{\text{vib}} \ (1 - p_s)), 1\right), \qquad (56)$$

where $\alpha_{\text{vib}}$ is the vibrance strength. Positive values of $\alpha_{\text{vib}}$ increase vibrance by amplifying low-saturation colors more strongly than highly saturated ones, whereas negative values reduce vibrance.

After applying either transformation, the adjusted HSV values $(p_h, p'_s, p_v)$ are converted back to the sRGB space to continue processing in the pipeline. Figure 29 illustrates an example from the S24 test set [16], rendered with our photofinishing pipeline using different saturation and vibrance settings. As shown, the saturation adjustment uniformly scales the chroma across all colors, whereas vibrance selectively boosts muted colors while preserving already saturated regions.

## I.7. Detail-Enhancement Control

Similar to the denoising strength control, our tool provides continuous control over the detail-enhancement strength. This is implemented as a linear interpolation between the image before and after the detail enhancement operation, modulated by a scalar parameter in the range $[0, 1]$. A value of 0 preserves the unprocessed image, a value of 1 applies the full operation, and intermediate values yield proportional effects.

## I.8. Sharpening Adjustment

We apply an optional, user-controlled, edge-aware sharpening step at the final stage of our pipeline, operating on the full-resolution sRGB image. Let $\mathbf{I}_{\text{in}}$ denote the input to the sharpening operator (equivalent to $\mathbf{I}_{\text{out}}$, the final output of our pipeline). We first compute a smoothed base layer $\mathbf{B}_{\text{sharp}}$ using a Gaussian blur with kernel size $\rho_{\text{sharp}} = 3$
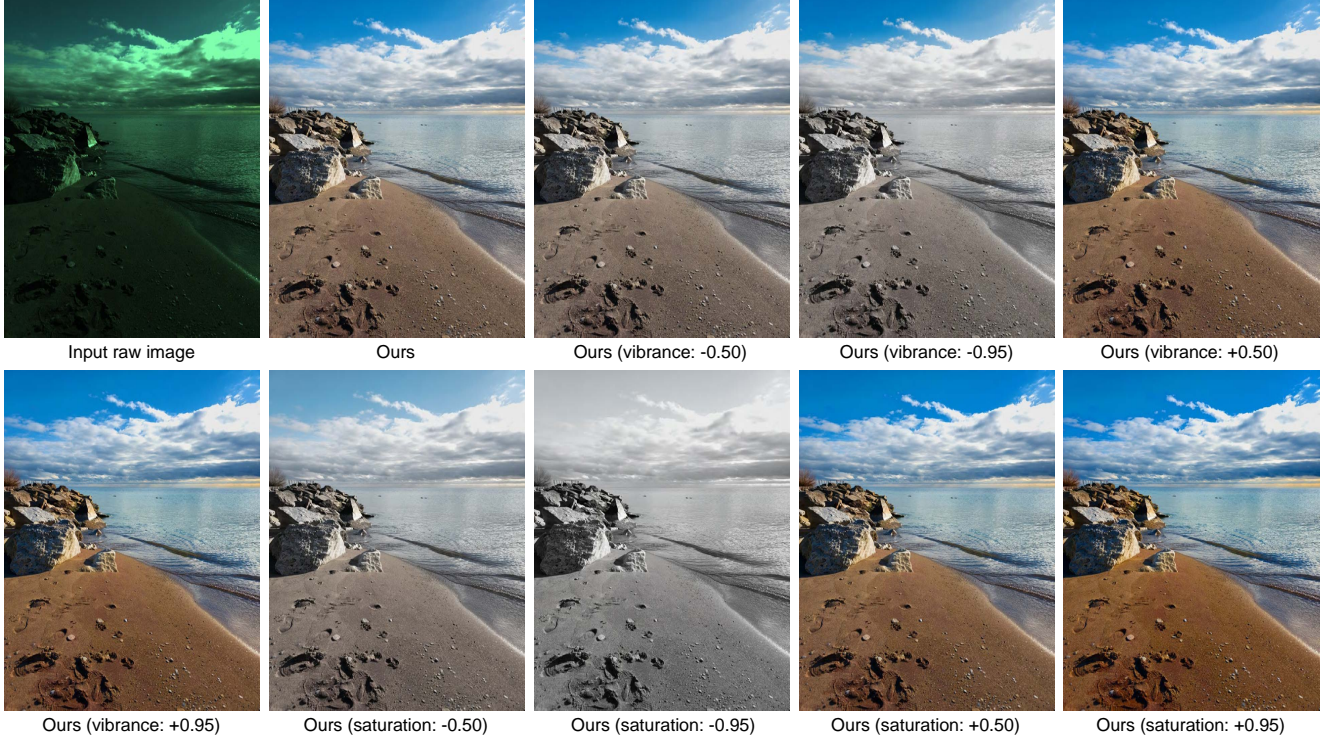
23

Figure 29. We allow the user to optionally adjust the saturation and vibrance in our pipeline. Shown here is an example from the S24 test set [16], including the input raw image, our output without adjustment ($\alpha_{\text{sat}} = \alpha_{\text{vib}} = 0.0$), and results with different saturation and vibrance settings ($\alpha \in \{-0.5, +0.5, -0.95, +0.95\}$). Saturation uniformly scales chroma across all colors, whereas vibrance selectively enhances muted colors while preserving already saturated regions.

and standard deviation $\sigma_{\text{sharp}} = 1.0$. The high-frequency detail layer is then defined as:

$$\mathbf{D}_{\text{sharp}} = \mathbf{I}_{\text{in}} - \mathbf{B}_{\text{sharp}}. \tag{57}$$

To prevent uniform sharpening of noisy flat regions, we construct an edge-aware mask. Horizontal and vertical gradients are estimated by convolving the input image with simple derivative kernels, producing response maps $\mathbf{G}_x$ and $\mathbf{G}_y$. Gradients are computed independently per channel using $[-1, 0, 1]$ derivative kernels, and the resulting magnitudes are averaged to form a normalized edge mask. The edge magnitude map is then computed as:

$$E_{\text{sharp}(x,y)} = \sqrt{G^2_{x(x,y)} + G^2_{y(x,y)}}, \tag{58}$$

which is normalized to form the edge mask $\mathbf{M}_{\text{sharp}} \in [0, 1]$. This mask ensures that sharpening is applied predominantly along edges. The final sharpened image is obtained as:

$$\mathbf{I}'_{\text{sharp}} = \mathbf{I}_{\text{in}} + \alpha_{\text{sharp}}\mathbf{D}_{\text{sharp}}\mathbf{M}_{\text{sharp}}, \tag{59}$$

where $\alpha_{\text{sharp}}$ controls the sharpening strength. Higher values of $\alpha_{\text{sharp}}$ increase edge contrast, while setting $\alpha_{\text{sharp}} = 0$ disables sharpening. Figure 30 shows qualitative results of the sharpening operator applied to an example image.

## I.9. Editing Picture Styles

As our photofinishing module governs the overall appearance and perceived "look and feel" of the final image, and given its modular design, we provide users with fine-grained control over each operator and enable flexible style editing. Specifically, our tool allows: 1) disabling individual operators (e.g., tone mapping, gamma correction) within a selected picture style, 2) interpolating operators' parameters between multiple picture styles, or 3) replacing specific operators of one style with those of another target style. We implement style mixing through simple linear interpolation of the operator parameters predicted by the corresponding photofinishing networks of each selected style. The interpolated parameters are then applied once through our photofinishing module, ensuring consistent and efficient rendering.

Figure 31 shows an example image captured with an iPhone 13 (an unseen camera in our training) and processed using our generic denoiser and the photofinishing module (trained on images taken by the S24's main camera) with different picture styles. The figure demonstrates three types of style mixing: 1) operator replacement (e.g., rendering the image with Style #0 while substituting chroma mapping from another style), 2) operator interpolation (e.g., perform-
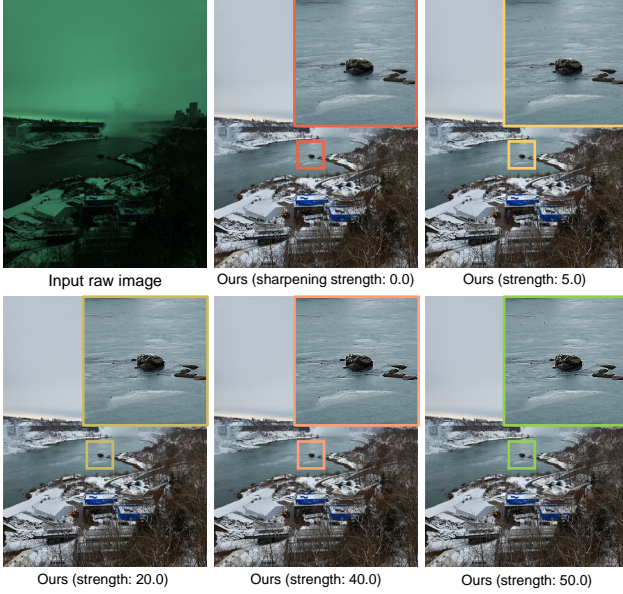
Figure 30. We allow the user to apply sharpening as the final step of our pipeline. The example shown is from the validation set of S24 [16], including the input raw image and our rendered outputs with different sharpening strengths $\alpha_{\text{sharp}}$. Increasing $\alpha_{\text{sharp}}$ enhances edge contrast, with $\alpha_{\text{sharp}} = 0$ corresponding to no sharpening.

ing chroma mapping by interpolating 50% from Style #3 and 50% from Style #5), and 3) operator disabling (e.g., rendering with all operators of Style #5 except global tone mapping, which is disabled).

### I.10. Re-Rendering with Embedded Raw

After rendering the raw image to sRGB, our tool embeds the original raw data within the JPEG container of the rendered output. This was achieved by first compressing the raw image as a JPEG using the Raw-JPEG Adapter [15], and then appending this raw-JPEG file (along with the operator parameters, DNG metadata, and the current editing configuration of the image in the tool) after the end-of-image (EOI) marker of the sRGB JPEG. In this way, the visible JPEG remains fully standards-compliant and decodable by any conventional viewer, while our tool can still access the embedded raw data for re-rendering and further adjustments.

This design offers several advantages. First, it ensures backward compatibility, as the sRGB JPEG remains viewable on all platforms without modification. At the same time, it enables the distribution of a single JPEG file that contains both the standard viewable image (compatible with any JPEG decoder) and the corresponding raw data (accessible by our tool), with a small increase in storage size compared to saving the raw data separately without compression (e.g., in DNG format). Lastly, by embedding the operator parameters, the user can later reset all manual adjustments and reapply a new set exactly as applied during

the initial rendering, allowing unlimited post-saving edits without needing to access the original raw DNG file.

For the Raw-JPEG Adapter [15], we provide two quality settings: raw-JPEG quality = 95, which typically adds about 2-3 MB to the sRGB JPEG file size, and raw-JPEG quality = 75, which adds around 1-2 MB. These values represent a modest storage overhead compared to alternatives such as HEIC on iPhone, which also supports post-capture re-rendering through Apple's Photos editing tool but often produces larger files (exceeding 10 MB in some scenes or picture styles) and DNG files, which typically range from 12 MB to 35 MB per 12-megapixel image. In contrast, for 12-megapixel images, our approach of appending the JPEG-compressed raw data to the final image requires, on average, only about 5-6 MB in total at raw-JPEG quality = 95 (including approximately 3 MB for the sRGB JPEG) and about 3-4 MB at raw-JPEG quality = 75.

To generate a raw-JPEG compressed file, we pre-process the raw images before JPEG compression using trained models corresponding to each target JPEG quality, as described in [15]. During decoding, and prior to JPEG decompression, we inverted the Raw-JPEG Adapter operators using the stored metadata of the operator parameters embedded in the raw JPEG file.

In our implementation, we employed the Raw-JPEG Adapter model variant without the DCT component, which was recommended for better generalization to unseen cameras during training [15]. The Raw-JPEG Adapter model is lightweight, with only 32,076 parameters, and introduces an average overhead of about 0.3 seconds for raw image processing before encoding and 0.1 seconds for restoration.

With this design, our tool can re-render or edit images with the same level of functionality as when operating directly on the original DNG raw files during the post-render stage, even after multiple re-rendering operations without cumulative degradation in accuracy. This is achieved by preserving the complete raw data within the final JPEG file, with only a slight loss in accuracy due to compression. Figure 32 demonstrates this capability, showing an example where an image saved using our tool is later re-rendered both with the same settings and with different styles and parameter configurations. A quantitative evaluation is provided in Sec. K.2.3, where we compare our design against prior re-rendering alternatives, showing that embedding the raw data within the JPEG enables higher-quality re-rendering results.

### I.11. Processing Input sRGB Images

To broaden the applicability of our framework beyond raw images, we also support standard 8-bit sRGB inputs (e.g., JPEG/PNG) that are not accompanied by embedded raw data. Since our pipeline is designed to operate on raw images, we first linearize the input sRGB image using the

Figure 31. Our modular photofinishing design provides users with full control over the picture style of the final image. In addition to selecting from pre-defined picture styles (each rendered by a pre-trained photofinishing module), users can also mix operators from different styles or disable specific operators entirely. The shown example was captured with the iPhone 13 main camera, which is an unseen camera for our pre-trained photofinishing modules.

method of [10], where a lightweight network maps the input sRGB image to the CIE XYZ linear space.

We then follow the CIE XYZ-to-raw mapping strategy described in Sec. 4.2.4 of [10]. In this step, we fix the CIE XYZ-to-raw $3{\times}3$ matrix to the one corresponding to the S24 main camera under D50 illumination (CCT=5000K). The mapped raw image is then multiplied by the D65 light color in the raw space, yielding a synthetic raw-like representation of the input sRGB image within the S24 main

camera's raw space. Unlike the original design in [10], which employed two networks to map between sRGB and CIE XYZ in both directions, we only use a single network to map from sRGB to CIE XYZ. To further reduce computational cost and memory usage, we use a lightweight variant of the original architecture. Specifically, we reduced the local sub-network depth in [10] from 16 to 8 convolutional blocks, and the number of channels per convolutional layer from 32 to 24. For the global sub-network in [10], we
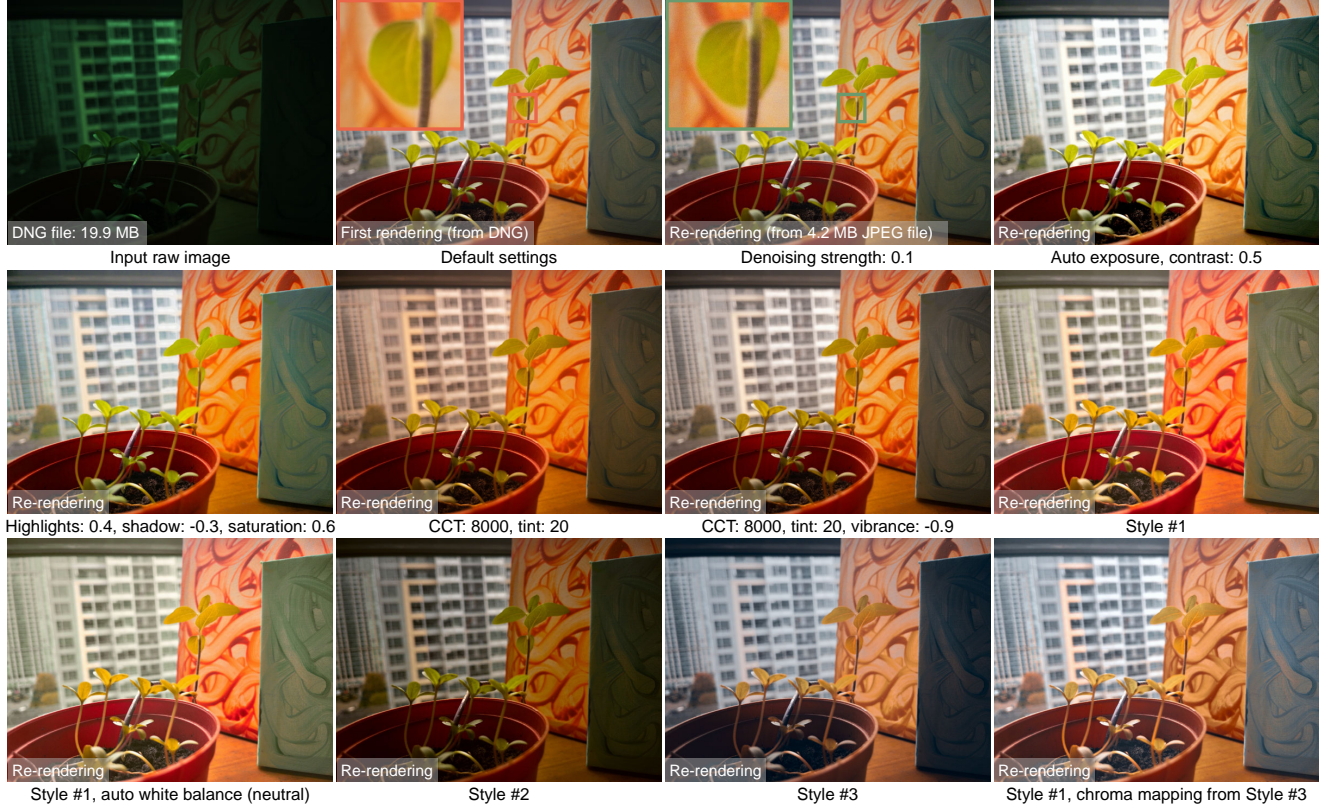
Figure 32. Our framework enables post-editable re-rendering of saved JPEG images with the same functionality available to the original raw data. This is achieved by embedding the raw image in a compact form alongside the final sRGB rendering. Shown here is an example image captured with an iPhone 13 and processed by our tool. The first result shows the default rendering, while the remaining results are obtained by re-rendering from the restored raw image using different rendering settings.

use a depth of 4 (instead of 5) with 24 channels (instead of 64), and reduce the final three fully connected layers from [1024, 1024, 1024] to [512, 256, 256], while disabling the dropout layer. After these modifications, the total number of parameters in the linearization model is 693,493 ($\sim$2.7 MB), resulting in only a light additional load on the GPU.

We trained this lightweight linearization model for 300 epochs using the Adam optimizer [54] on 512×512 non-overlapping patches from the sRGB-to-CIE XYZ dataset [10], with a mini-batch size of 8 and an L2 regularization factor of $10^{-6}$ (in contrast to $10^{-3}$ used in the original work).

With this synthetic raw conversion process, our tool remains applicable to sRGB inputs rendered by any camera or software ISP, as well as to synthetically generated sRGB images—albeit with reduced accuracy due to residual non-linearities that cannot be fully removed by linearization. Figure 33 illustrates two examples: the first is an sRGB image captured and rendered by the iPhone's native camera ISP, and the second is generated by the Gemini 2.5 Flash Image model. Both images are linearized, mapped to the S24 camera raw space, and subsequently rendered with dif-

ferent settings using our tool.

## J. User Study Details

As discussed in the main paper, we conducted a user study to compare our method with the Samsung S24 native camera ISP and Adobe Lightroom (using its built-in auto-enhancement feature). We captured 45 scenes with the S24 main camera in Pro mode (which outputs DNG raw files) and re-captured the same scenes using the native camera application.

The images covered a variety of conditions, including indoor, outdoor daylight, sunset, and low-light scenes. The DNG files were processed with Adobe Lightroom (auto enhancement) and with our method, resulting in three versions for each scene: ours, native ISP, and Lightroom.

Participants were presented with these three versions side-by-side in a custom user-study GUI. The order of the images was randomized for each trial to prevent bias, and the internal order of the methods was randomized as well. For each scene, participants were asked to select the best image under four criteria: 'color quality' — `How good and natural do the colors look?`, 'bright-

Input sRGB image          Synthesized raw image          Default settings

AWB (preferred biases)          Highlights: 0.7, contrast: 0.5, saturation: 0.3          Style #1

Input sRGB image (generated by Gemini)          Synthesized raw image          Default settings

Default settings w/o tone mapping          Style #4 w/ AWB (neutral)          Style #5

Figure 33. To broaden the applicability of our tool to sRGB images saved without embedded raw data (i.e., outside our framework), we linearize the input sRGB images (whether produced by unknown ISPs or generated by AI models) into synthetic raw-like representations, thereby enabling the full functionality of our pipeline. The top example shows an image captured and rendered by the iPhone camera ISP, from which we generate a raw-like image and process it with different settings using our tool. The second example is an image generated by the Gemini 2.5 Flash Image model with the prompt: `Generate an ancient Egyptian scene with pyramids, the Nile, and palm trees glowing under a golden sunset`. We convert it into a raw-like image and apply our tool with different rendering options.
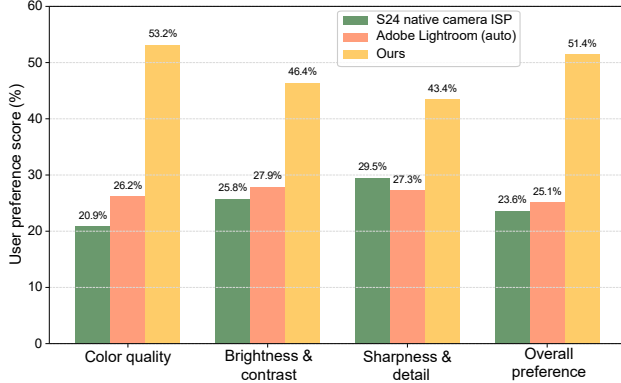
Figure 34. User study results. Preference rates (%) for each criterion comparing our method, the Samsung S24 native ISP, and Adobe Lightroom (auto enhancement). Our method is consistently preferred by participants across all four criteria.

ness & contrast' — `Is the image clear, with a good balance between light and dark areas?`, 'sharpness & detail' — `How crisp and detailed does the image look?`, and 'overall preference' — `Which image do you simply like more overall?`. A total of 20 participants took part in the study.

As shown in Fig. 34, our method was consistently preferred across all four criteria: 53.2% in 'color quality', 46.4% in 'brightness & contrast', 43.4% in 'sharpness & detail', and 51.4% in 'overall preference', with margins of +13.9–27.0% over the closest competitor. These results demonstrate that users consistently favor our rendering pipeline over both the commercial Samsung ISP and Lightroom auto enhancement. See Fig. 35 for examples from the captured images and the rendering results from the native camera ISP, Adobe Lightroom (auto enhancement), and our method.

## K. Ablation and Additional Results

### K.1. Ablation Studies

We conducted a series of ablation studies to validate the effectiveness of each component in our method. In this subsection, we describe these experiments in detail and present the corresponding results.

### K.1.1. Denoising Ablations

As discussed in the main paper, we evaluated three NafNet variants [30] for raw image denoising: lite (245 K parameters), base (933 K), and large (3.63 M). To compare NafNet with an alternative architecture for raw image denoising, Table 7 shows results for the large NafNet model and Restormer [88] (26.1 M), both trained on the S24 raw/denoised pairs [16]. The input to each model is a noisy raw image, and the output is a denoised raw image, eval-

Table 7. Raw image denoising results on the S24 noisy/denoised test set [16] using NafNet [30] and Restormer [88]. The best results are highlighted in yellow.

| Method | S24 Test Set (Noisy/Denoised) | |
| --- | --- | --- |
| | PSNR↑ | SSIM↑ |
| NafNet (large, ours) [30] | 57.33 | 0.999 |
| Restormer [88] | 55.42 | 0.999 |

Table 8. Ablation study of raw denoising and detail enhancement on the S24 [16] test set. Inputs are noisy raw images, and outputs are compared against ground-truth sRGB images. Denoising was performed using the base model (∼933K parameters). The best results are highlighted in yellow.

| Framework Variation | S24 Test Set | |
| --- | --- | --- |
| | PSNR↑ | SSIM↑ |
| w/o raw denoising | 24.48 | 0.731 |
| w/ raw denoising | 26.48 | 0.883 |
| w/ raw denoising + enhancement | 27.52 | 0.922 |

uated against the pseudo ground truth provided in the S24 dataset. The large NafNet variant achieves superior results while requiring substantially fewer parameters.

Table 8 reports results on the S24 test set, where the inputs are noisy raw images and both the outputs and ground truth are sRGB-rendered images. We compare three configurations of our framework (with the photofinishing module enabled in all cases): without denoising, with denoising only (excluding the detail-enhancement network), and with both denoising and detail enhancement, using the base denoising model. The results show that image denoising yields significant improvements.

### K.1.2. Upsampling Ablations

In both the main paper and this supplementary material (Sec. E), we have introduced our regularized variant of BGU [29]. Here, we present ablation results comparing different guided upsampling approaches, including BGU without our regularization. For completeness, we also compare against the regularization used in the original BGU Halide implementation, which we denote as Halide regularization.

A key challenge in this evaluation arises from the ground-truth sRGB images in the S24 test set. These ground-truth images include sharpening and detail enhancements that are not present in the pseudo ground-truth denoised raw images. Since the guide image (denoised raw or any of its variants, e.g., linear sRGB) lacks such details, no guided upsampling method can fully recover them. To address this, we conducted experiments under two ground-truth settings:

1. Original S24 ground truth: sRGB images with sharpening, detail enhancement, and compression (as provided in the original dataset).
2. Alternative ground truth: sRGB images rendered di-

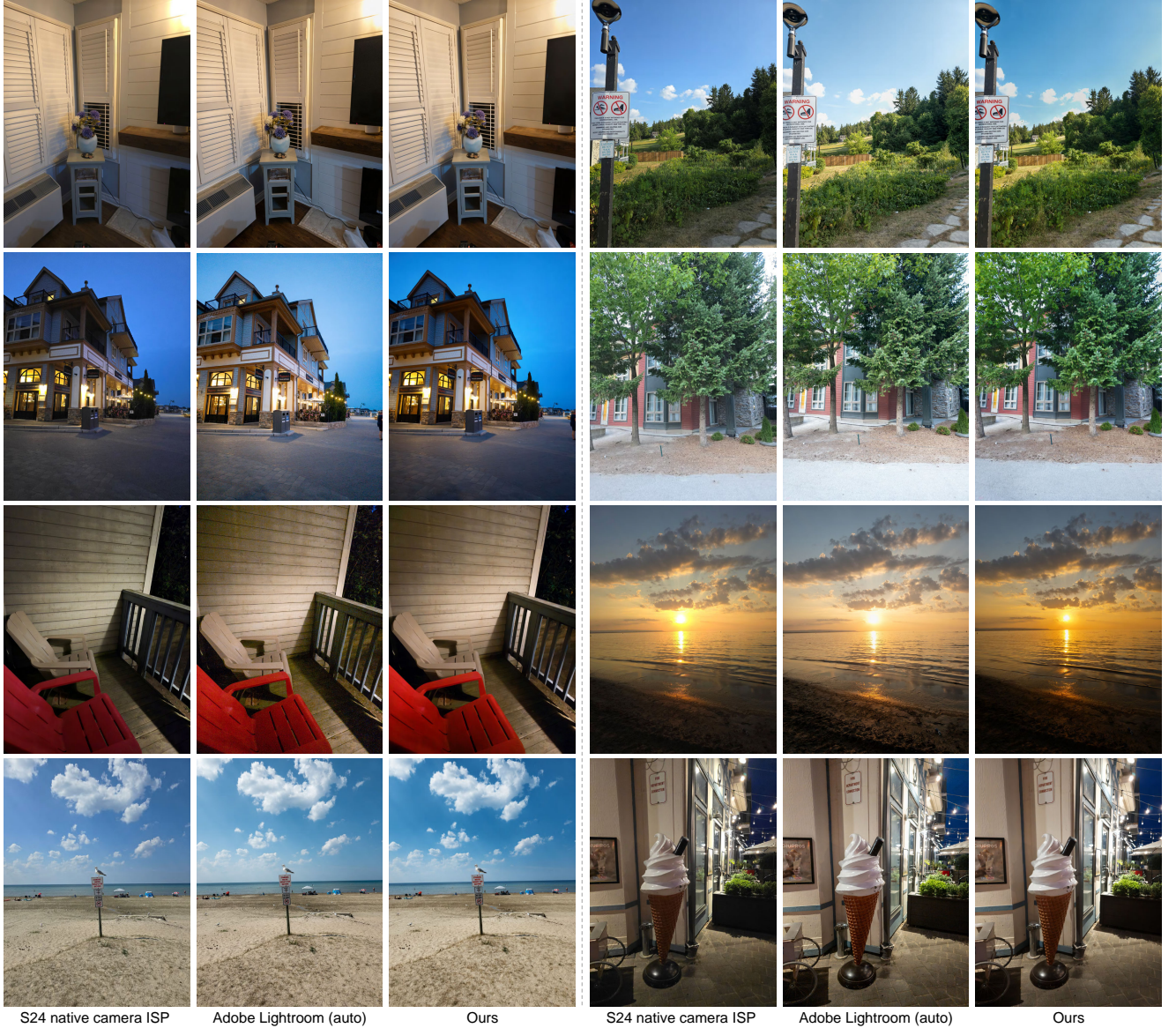| S24 native camera ISP | Adobe Lightroom (auto) | Ours | S24 native camera ISP | Adobe Lightroom (auto) | Ours |

Figure 35. Example comparisons shown to participants during the user study. For each scene, participants evaluated three versions: our method, the Samsung S24 native ISP, and Adobe Lightroom (auto enhancement). The order was randomized in the GUI for every trial.

rectly from the denoised raw DNG files using Adobe Photoshop, with sharpening and detail enhancements disabled and no compression applied.

For both ground-truth settings, we evaluated two guide image versions: 1) the denoised raw image and 2) its linear sRGB conversion. The corresponding source image was generated by downsampling the guide image to one-quarter of its original resolution, while the target image was obtained by downsampling the ground-truth sRGB image to the same resolution.

We utilized ground-truth sRGB images in these experiments, rather than outputs from our photofinishing module, to isolate the evaluation of the guided upsampling meth-

ods themselves under the best-case scenario, without indirect effects from the photofinishing stage. Table 9 reports results comparing BGU [29] with our regularization, BGU without regularization, and BGU with the Halide regularization, along with two alternative guided upsampling methods: guided linear upsampling (GLU) [74] and classical guided image filtering (GF) [43]. Figure 36 provides a visual comparison of BGU results without and with our proposed regularization. Overall, our regularization produces reconstructions that are both visually and quantitatively closer to the ground truth.

30

Table 9. Comparison of different guided upsampling methods on the S24 test set [16] under two ground-truth settings: original S24 ground truth (with detail enhancements) and alternative ground truth (without enhancements). Results are reported for two guide image versions: denoised raw and linear sRGB (LsRGB). The best results are highlighted in **yellow**.

| Method | Original S24 GT (Denoised Raw/sRGB) | | Alt. S24 GT (Denoised Raw/sRGB) | | Original S24 GT (LsRGB/sRGB) | | Alt. S24 GT (LsRGB/sRGB) | |
|---|---|---|---|---|---|---|---|---|
| | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ |
| GLU [74] | 33.46 | 0.919 | 37.95 | 0.966 | 33.47 | 0.919 | 38.42 | 0.971 |
| GF [43] | 31.58 | 0.907 | 33.74 | 0.943 | 31.97 | 0.912 | 34.84 | 0.959 |
| BGU (w/o regularization) [29] | 33.73 | **0.928** | 40.82 | 0.986 | 33.77 | 0.928 | 41.00 | 0.988 |
| BGU (w/ Halide regularization) [29] | 32.49 | 0.916 | 37.28 | 0.977 | 32.78 | 0.918 | 38.26 | 0.980 |
| BGU (w/ our regularization) [29] | **33.85** | 0.927 | **41.69** | **0.987** | **33.92** | **0.928** | **42.01** | **0.989** |



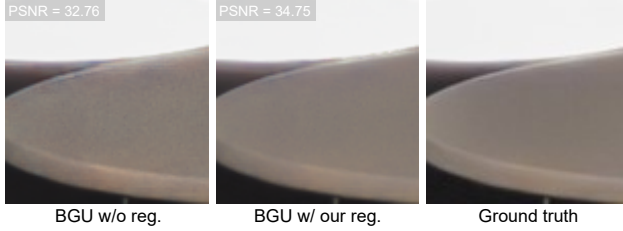BGU w/o reg.    BGU w/ our reg.    Ground truth

Figure 36. Visual comparison of guided upsampling on an example from the S24 dataset [16]. We compare BGU [29] without regularization against our proposed regularization. Our regularization better preserves fine details and produces results closer to the ground truth.

### K.1.3. Photofinishing Loss Ablations

To analyze the contribution of each loss term in our photofinishing module, we conducted a detailed ablation study on the default picture style (Style #0) of the S24 dataset [16] (Tables 10 and 11). For these experiments, we chose to process images starting from the pseudo ground-truth denoised raw inputs, which were subsequently mapped to the linear sRGB color space. This setup allows us to isolate and evaluate the photofinishing module independently, without the influence of potential degradations introduced by the denoising stage in our full pipeline.

When training with all loss terms enabled, our model achieves the best accuracy in both PSNR and SSIM, as highlighted in Table 10. Removing any single loss term slightly decreases the accuracy. Conversely, progressively adding the losses one by one (Table 11) results in steady improvements, indicating that each term contributes positively to the overall quality.

Figure 37 shows that the tone-mapping loss ($\ell_{\text{TM}}$) and luminance-consistency loss ($\ell_{\text{luma}}$) lead to more balanced tone reproduction across the image, encouraging the global and local tone-mapping stages of the photofinishing module to function as intended. The global stage refines the overall tone distribution of the gain-adjusted linear sRGB image, while the local stage further adjusts contrast in spatially varying regions. Both stages maintain a comparable overall luminance level, which encourages the local tone-mapping stage to focus on localized adjustments rather than global

Table 10. Ablation study showing the impact of removing each loss term individually on the results of our photofinishing module, evaluated on the S24 test set [16]. The input images are linear sRGB frames generated from pseudo ground-truth denoised images in the S24 test set, downsampled to one-quarter resolution. The outputs are compared against the ground truth at the same resolution. The best results are highlighted in **yellow**.

| | | | | S24 Test Set (1/4 LsRGB/sRGB) | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| $\ell_1$ | $\ell_{\text{SSIM}}$ | $\ell_{\text{perc}}$ | $\ell_{\Delta E}$ | $\ell_{\text{CbCr}}$ | $\ell_{\text{LuT-s}}$ | $\ell_{\text{luma}}$ | $\ell_{\text{TM}}$ | $\ell_{\text{LTM-s}}$ | PSNR | SSIM |
| | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | 27.01 | 0.934 |
| ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | 27.01 | 0.931 |
| ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | 26.89 | 0.934 |
| ✓ | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | 26.96 | 0.934 |
| ✓ | ✓ | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | 26.99 | 0.934 |
| ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ | ✓ | ✓ | 27.18 | 0.934 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ | ✓ | 27.20 | 0.934 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ | 27.13 | 0.934 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | 27.05 | 0.934 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | **27.49** | **0.939** |

Table 11. Ablation study showing the cumulative effect of adding each loss term progressively to our photofinishing module, evaluated on the S24 test set [16]. Starting from a baseline trained with $\ell_1$ only, additional loss terms are introduced one at a time. Input images are linear sRGB frames generated from pseudo ground-truth denoised images in the S24 test set, downsampled to one-quarter resolution. The outputs are compared against the ground truth at the same resolution. The best results are highlighted in **yellow**.

| | | | | S24 Test Set (1/4 LsRGB/sRGB) | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| $\ell_1$ | $\ell_{\text{SSIM}}$ | $\ell_{\text{perc}}$ | $\ell_{\Delta E}$ | $\ell_{\text{CbCr}}$ | $\ell_{\text{LuT-s}}$ | $\ell_{\text{luma}}$ | $\ell_{\text{TM}}$ | $\ell_{\text{LTM-s}}$ | PSNR | SSIM |
| ✓ | | | | | | | | | 26.80 | 0.930 |
| ✓ | ✓ | | | | | | | | 27.03 | 0.935 |
| ✓ | ✓ | ✓ | | | | | | | 27.07 | 0.935 |
| ✓ | ✓ | ✓ | ✓ | | | | | | 27.08 | 0.935 |
| ✓ | ✓ | ✓ | ✓ | ✓ | | | | | 27.09 | 0.934 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | | | 27.09 | 0.934 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | | 27.13 | 0.935 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | 27.05 | 0.934 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | **27.49** | **0.939** |

shifts. Similarly, Fig. 38 shows that our chosen loss weight for the local tone-mapping smoothness term ($\lambda_{\text{LTM-s}}=0.6$)

Table 12. Effect of different values of the local tone-mapping smoothness loss weight ($\lambda_{\text{LTM-s}}$) on the S24 test set [16]. Input images are pseudo ground-truth denoised raw images mapped to the linear sRGB space at one-quarter of the original raw resolution. The best results are highlighted in yellow.

| Loss weight ($\lambda_{\text{LTM-s}}$) | S24 Test Set (1/4 LsRGB/sRGB) | |
| --- | --- | --- |
| | PSNR↑ | SSIM↑ |
| 0.06 | 27.13 | 0.935 |
| 0.6 | **27.49** | **0.939** |
| 6.0 | 26.19 | 0.933 |

Table 13. Ablation on the design of our photofinishing module. For this ablation, we remove one network of our photofinishing module at both training and inference to evaluate its contribution. Results are shown on the S24 test set [16] (default style; Style #0), where the input images are pseudo ground-truth denoised raw images mapped to the linear sRGB space and downsampled to one-quarter of the original resolution, and the ground truth is the corresponding sRGB image at the same resolution. The best results are highlighted in yellow.

| S24 Test Set (1/4 LsRGB/sRGB) | | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| Gain | GTM | LTM | Chroma | Gamma | PSNR↑ | SSIM↑ |
| | ✓ | ✓ | ✓ | | 26.68 | 0.931 |
| ✓ | | ✓ | ✓ | | 26.84 | 0.934 |
| ✓ | ✓ | | ✓ | | 24.66 | 0.898 |
| ✓ | ✓ | ✓ | | | 26.45 | 0.928 |
| ✓ | ✓ | ✓ | ✓ | | 26.91 | 0.935 |
| ✓ | ✓ | ✓ | ✓ | ✓ | **27.49** | **0.939** |

achieves a good balance between suppressing local artifacts and preserving fine details (see Table 12 for the quantitative comparison).

### K.1.4. Photofinishing Design Ablations

We conducted a set of ablation studies to evaluate our design against alternative configurations. As described earlier, the proposed photofinishing module consists of five networks: digital gain, GTM, LTM, chroma mapping, and gamma correction. Each network predicts the parameters of its corresponding operator. The rest of this subsection details the contribution of each component.

**Overall Module Composition.** Table 13 shows the results of different variants of our photofinishing module, where one network is removed at both training and testing to evaluate the contribution of each stage. As shown, the proposed full design not only provides greater modularity and user control but also achieves the best quantitative results. Figure 39 shows a qualitative comparison.

**Local Tone-Mapping and Chroma Design.** We conducted additional experiments to ablate the design choices of our LTM and chroma mapping components. The LTM network predicts five spatially varying coefficient maps:

Table 14. Ablation on the local tone-mapping (LTM) design. We progressively add each predicted coefficient map ($\mathbf{A}_{\text{LTM}}$-$\mathbf{G}_{\text{LTM}}$) to evaluate their contribution. The "All (w/o multi-scale guide)" variant includes all five parameters but omits the multi-scale guidance network, using a single guidance network in its place. Results are on the S24 test set [16] using pseudo ground-truth denoised raw inputs mapped to linear sRGB at one-quarter resolution. The best results are highlighted in yellow.

| LTM Design | S24 Test Set (1/4 LsRGB/sRGB) | |
| --- | --- | --- |
| | PSNR↑ | SSIM↑ |
| $\mathbf{A}_{\text{LTM}}$ | 26.49 | 0.927 |
| $\mathbf{A}_{\text{LTM}}, \mathbf{B}_{\text{LTM}}$ | 26.54 | 0.928 |
| $\mathbf{A}_{\text{LTM}}, \mathbf{B}_{\text{LTM}}, \mathbf{C}_{\text{LTM}}$ | 26.58 | 0.927 |
| $\mathbf{A}_{\text{LTM}}, \mathbf{B}_{\text{LTM}}, \mathbf{C}_{\text{LTM}}, \mathbf{W}_{\text{LTM}}$ | 26.86 | 0.931 |
| All (w/o multi-scale guide) | 26.87 | 0.933 |
| Ours (all) | **27.49** | **0.939** |

$\mathbf{A}_{\text{LTM}}$, $\mathbf{B}_{\text{LTM}}$, $\mathbf{C}_{\text{LTM}}$, $\mathbf{W}_{\text{LTM}}$, and $\mathbf{G}_{\text{LTM}}$. We progressively enabled these coefficients to evaluate their contribution, and we further examined the effect of the LTM grid size and the design of the chroma LuT used for chroma mapping.

Table 14 reports results of training our photofinishing module with different subsets of the LTM parameters. We began by predicting a single-channel pixel-wise map representing $\mathbf{A}_{\text{LTM}}$ to modulate the tone-mapping exponent. Next, we introduced $\mathbf{B}_{\text{LTM}}$ to allow adaptive control over the compression strength in the denominator of the tone-mapping function, followed by adding $\mathbf{C}_{\text{LTM}}$ to model asymmetric behavior between bright and dark regions. We then incorporated the blending map $\mathbf{W}_{\text{LTM}}$ to locally mix the LTM output with the globally tone-mapped result $\mathbf{I}_{\text{GTM}}^{\downarrow}$, and subsequently added the gain map $\mathbf{G}_{\text{LTM}}$, which adjusts local exposure before tone mapping.

This full set of parameters was tested both without and with the multi-scale guidance network to evaluate its contribution. See Fig. 40 for a qualitative example.

Table 15 complements this analysis by showing the results of modifying the LTM grid size and the chroma mapping design. We first reduced the LTM grid size from $64\times64\times18\times5$ to $32\times32\times9\times5$ to examine the impact of spatial resolution of the LTM parameter grid. Next, we evaluated predicting the chroma LuT directly instead of learning it as a residual to the image-independent learnable LuT.

Lastly, we evaluated a smaller chroma LuT with $N_h=12$ bins (i.e., a LuT size of $12 \times 12 \times 2$) instead of $N_h=24$ (i.e., $24 \times 24 \times 2$) to study the effect of chroma quantization granularity. The results demonstrate that both the residual formulation and higher grid/LuT resolution contribute to improved reconstruction accuracy, with our full configuration achieving the best overall results.

**Multi-Branch and Coordinate Attention.** We further evaluated the contribution of the multi-branch convolutional (MBConv) and coordinate attention (CA) [45] blocks
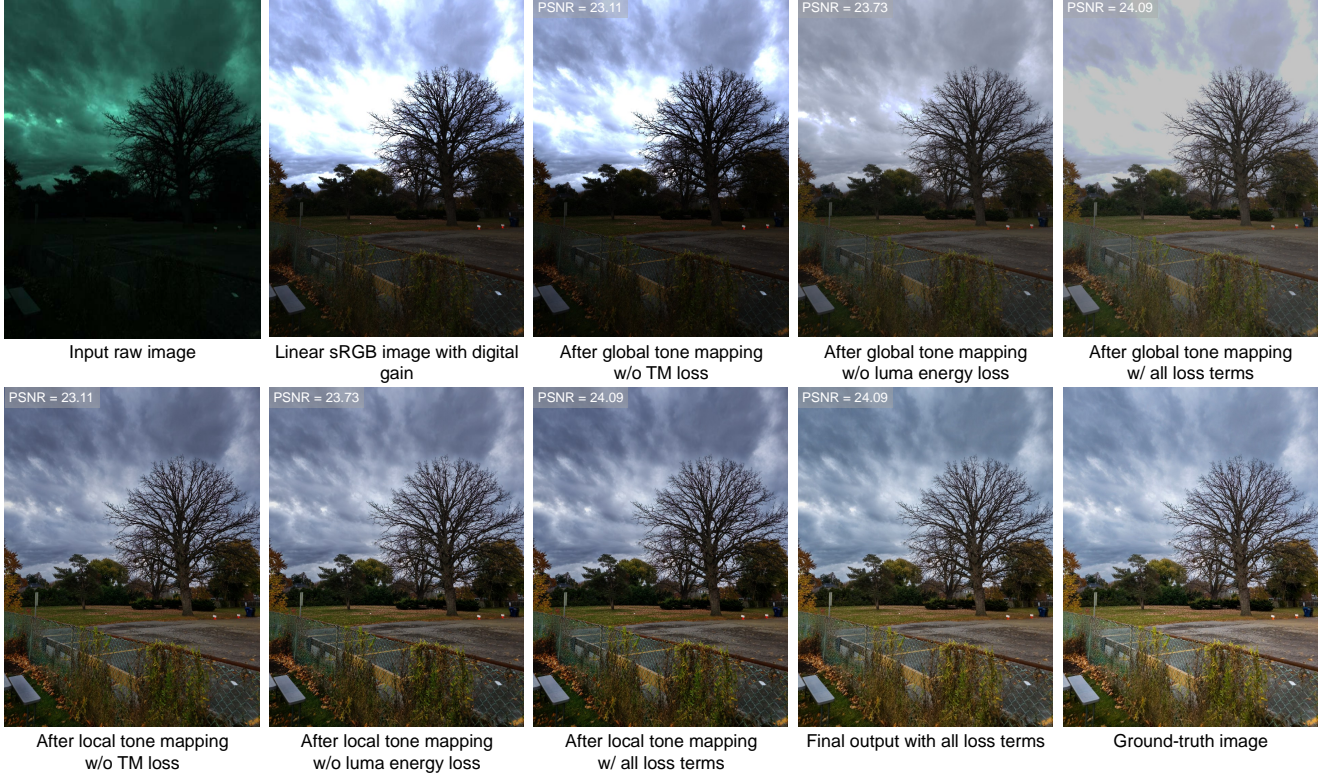
Figure 37. Visual comparison showing the impact of training with and without the tone mapping loss ($\ell_{\text{TM}}$) and the luma energy loss ($\ell_{\text{luma}}$). Shown are the input raw image, the color-corrected linear sRGB image after digital gain, the results after global and local tone mapping, and the final photofinishing output compared against the ground truth at the same resolution (i.e., one-quarter of the original raw resolution). The input sample is from the validation split of the S24 dataset [16], where the photofinishing result shown here is obtained by processing the pseudo ground-truth denoised raw image mapped to the linear sRGB space.

Table 15. Ablation on the design of the chroma mapping and LTM networks. We evaluate 1) predicting the chroma LuT directly instead of as a residual, 2) using a smaller number of chroma bins ($N_h{=}12$ instead of $N_h{=}24$), and 3) reducing the LTM grid size to $N_g{=}32$ and depth to $N_c{=}9$, instead of $N_g{=}64$ and $N_c{=}18$. Results are shown on the S24 test set [16], where the input images are pseudo ground-truth denoised raw images mapped to linear sRGB and downsampled to one-quarter of the original resolution. The best results are highlighted in yellow.

| Variant | S24 Test Set (1/4 LsRGB/sRGB) | |
| --- | --- | --- |
| | PSNR ↑ | SSIM ↑ |
| Chroma LuT (no residual) | 27.07 | 0.935 |
| Chroma LuT bins ($N_h = 12$) | 26.88 | 0.932 |
| LTM grid size ($N_c = 9$, $N_g = 32$) | 26.94 | 0.934 |
| Ours | **27.49** | **0.939** |

Table 16. Ablation on the effect of multi-branch (MBConv) and coordinate attention (CA) [45] in our photofinishing module. Results are reported on the S24 test set [16] (1/4 linear sRGB/sRGB). The best results are highlighted in yellow.

| Photofinishing Networks | S24 Test Set (1/4 LsRGB/sRGB) | |
| --- | --- | --- |
| | PSNR ↑ | SSIM ↑ |
| w/o MBConv | 26.95 | 0.934 |
| w/o CA | 26.93 | 0.933 |
| Ours (w/ MBConv & CA) | **27.49** | **0.939** |

(Sec. C.2) used across the networks of our photofinishing module. For this ablation, we trained and tested variants with each component removed independently. Table 16 shows that both the multi-branch structure and the coordinate attention blocks improve the result of our photofinishing module.

**Tone-Mapping Domain.** In our tone-mapping design, we apply both the global and local tone-mapping operators directly to the linear sRGB representation instead of to the luma (Y) channel of the YCbCr representation at each stage. While operating in the luma domain may seem more intuitive (since it isolates luminance from chroma), we found that this separation leads to performance degradation. Specifically, applying tone mapping (global and local) to the Y channel and then allowing the chroma-mapping network to process the CbCr components resulted in noticeably poorer qualitative and quantitative results. Table 17 compares our proposed design against the luma-based vari-
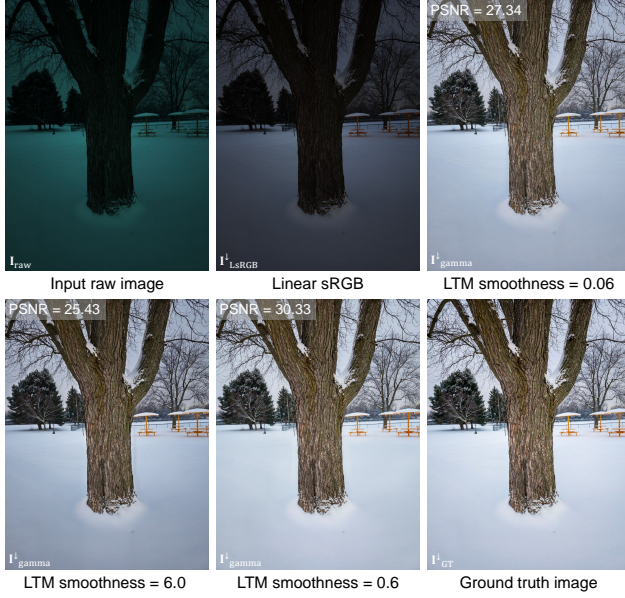
Figure 38. The impact of different values of $\lambda_{\texttt{LTM-s}}$ on the final output of the photofinishing module. Our chosen value ($\lambda_{\texttt{LTM-s}}$=0.6) achieves a good balance between smoothing the local tone-mapping maps and maintaining both qualitative and quantitative quality. The shown image is from the test set of the S24 dataset [16]. The displayed result corresponds to the final output of our photofinishing module, obtained by processing the pseudo ground-truth denoised raw image mapped to the linear sRGB space. Both the prediction and ground truth are shown at one-quarter of the original raw resolution.

Table 17. Comparison between our tone-mapping design and a luma-based tone-mapping design. In the luma-based design, global and local tone-mapping are applied only to the Y (luminance) channel of the YCbCr color space. In our design, tone-mapping is applied jointly to all RGB channels in the linear sRGB domain. Results are reported on the S24 test set [16] (1/4 linear sRGB/sRGB). The best results are highlighted in <mark>yellow</mark>.

| | S24 Test Set (1/4 LsRGB/sRGB) | |
|---|---|---|
| **Design** | **PSNR ↑** | **SSIM ↑** |
| Luma-based TM design | 25.80 | 0.922 |
| Ours | **27.49** | **0.939** |

ant on the S24 test set [16]. As shown, applying tone mapping in the linear RGB domain yields superior PSNR and SSIM values. See Fig. 41 for a qualitative comparison.

**3D Lookup Table.** In the main paper, we described an optional learnable global 3D LuT designed to improve the rendering of *artistic* picture styles. Specifically, we learn a global, image-independent $11{\times}11{\times}11$ RGB LuT, $\mathbf{L}_{\text{RGB}}$, which is applied before the chroma-mapping operator. This 3D LuT helps capture more aggressive color transforma-
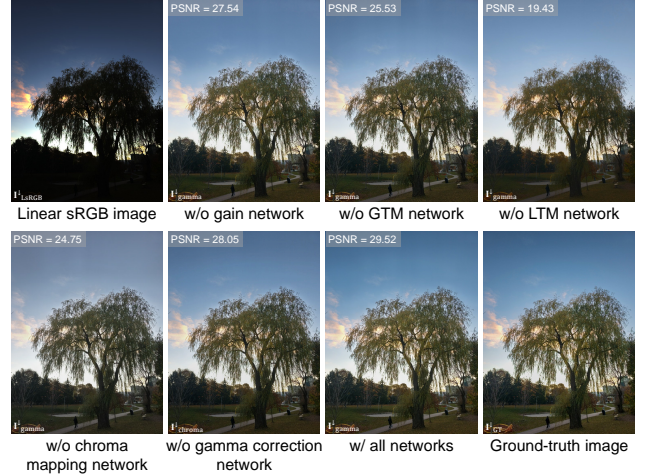


Figure 39. Output of the photofinishing module trained without one of the processing stages (digital gain, global tone mapping, local tone mapping, chroma mapping, and gamma correction) compared to the full model trained with all stages enabled. The shown image is from the validation set of the S24 dataset [16]. The results are obtained by processing the pseudo ground-truth denoised raw image provided in the dataset, mapped to the linear sRGB space at one-quarter of the original raw resolution, and compared against the ground-truth sRGB image at the same resolution.
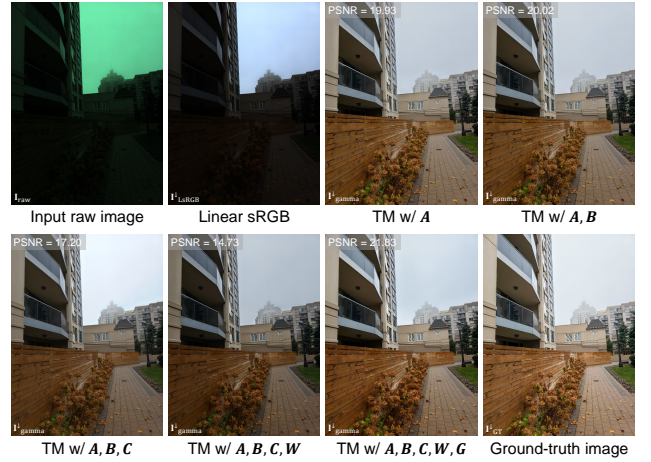


Figure 40. Qualitative results from our photofinishing module trained with progressively added LTM parameters ($\mathbf{A}_{\text{LTM}}$, $\mathbf{B}_{\text{LTM}}$, $\mathbf{C}_{\text{LTM}}$, $\mathbf{W}_{\text{LTM}}$, and $\mathbf{G}_{\text{LTM}}$) during both training and testing. The input is a pseudo ground-truth denoised raw image mapped to linear sRGB (1/4 resolution), and the output is compared against the corresponding sRGB ground truth. Examples are taken from the S24 validation set [16].

tions that the 2D chroma LuT, $\mathbf{L}_{\texttt{chroma}}$, constructed by our chroma-mapping network is limited in modeling (see Fig. 42 for a qualitative example).

However, we observed that when learning simpler picture styles–such as the default color style used in the S24
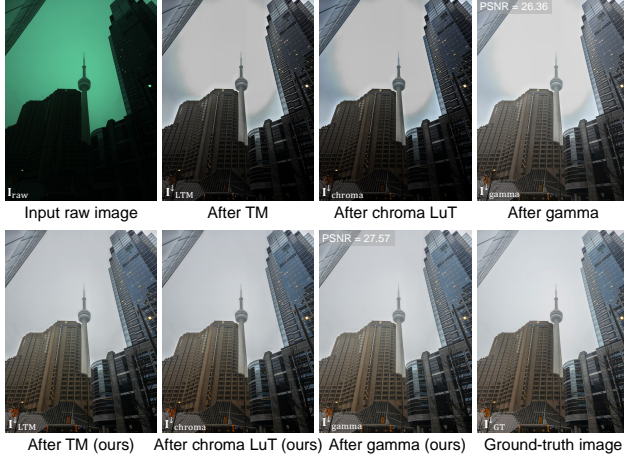
Figure 41. Qualitative comparison between the luma-based tone-mapping design and our proposed design. Inputs to the photofinishing module are pseudo ground-truth denoised raw images mapped to linear sRGB and downsampled to one-quarter of the original resolution. Results are shown from the validation set of the S24 dataset [16].
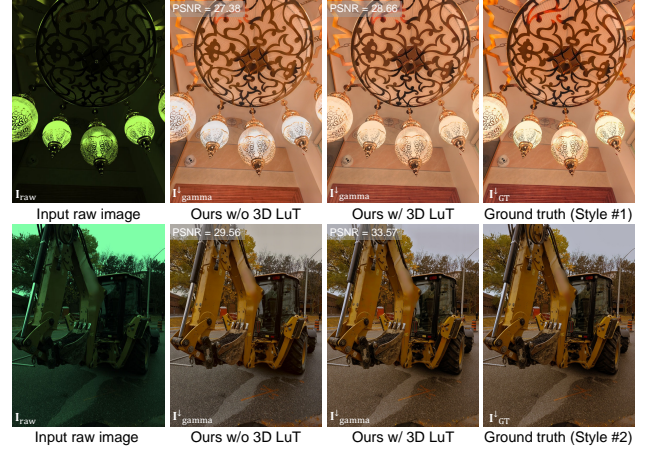


Figure 42. Impact of incorporating a global 3D LuT when learning artistic picture styles. Results are shown on the S24 test set [16], where the input to the photofinishing module is the linear sRGB version of the pseudo ground-truth denoised raw image downsampled to one-quarter of the original resolution.

dataset [16] (Style #0)–the 3D LuT provides negligible benefit. We evaluated our photofinishing module with and without the 3D LuT, where the input is the linear sRGB image derived from the pseudo ground-truth denoised raw data at one-quarter of the original raw resolution. Similar to our other ablations, we perform this analysis in isolation from other components such as guided upsampling or denoising to focus on the photofinishing performance.

Table 18 reports the results on the artistic picture styles of the S24 dataset (Styles #1-5). As shown, incorporating the 3D LuT provides a consistent improvement across styles, with only a minor increase of approximately 4,000 parameters per style. In contrast, Table 19 shows that for the default picture style (Style #0), the 3D LuT yields no significant improvement and even slightly degrades SSIM.

**Comparison with HDRNet [39].** While both our LTM network and HDRNet [39] aim to predict locally adaptive tone-mapping operators, their formulations and architectures are fundamentally different. Both methods employ a low-resolution bilateral grid of coefficients that are sliced using a high-resolution guidance image to produce pixel-wise operator parameters. However, the tone-mapping operators modeled by our LTM differ from those in HDRNet, and our network design is entirely distinct (Sec. C.5). To better highlight these differences, we conducted two ablation experiments. In the first, we replaced our LTM network and its operators with HDRNet, while keeping all other photofinishing networks unchanged. In the second, we replaced the entire photofinishing module with a single HDRNet model. All models were trained jointly following



Figure 43. Qualitative comparison between our method and HDR-Net [39]. We show results when HDRNet replaces 1) only the LTM network/operator and 2) the entire photofinishing module. The example is from the S24 test set [16]. In our result, multiscale processing and refinement (Sec. B.1) are applied within the LTM stage.

the same configuration used for our photofinishing module on the S24 dataset [16]. Table 20 reports the quantitative results on the S24 test set. As shown, our design achieves higher PSNR and SSIM while requiring fewer parameters. See Fig. 43 for a qualitative comparison.

35

Table 18. Ablation on the effect of the learned 3D LuT for rendering artistic picture styles using our photofinishing module. Results are shown on the S24 test set [16], where the input consists of pseudo ground-truth denoised raw images mapped to the linear sRGB space and downsampled to one-quarter of the original resolution, and the ground truth is the corresponding sRGB image in the target picture style at the same resolution. The evaluation covers five distinct artistic picture styles (Styles #1–5). The best results are highlighted in yellow.

| Photofinishing Module | S24 Test Set (1/4 LsRGB/sRGB) | | | | | | | | | | # params (per style) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Style #1 | | Style #2 | | Style #3 | | Style #4 | | Style #5 | | |
| | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | |
| Without 3D LuT | 26.08 | 0.926 | 28.70 | 0.923 | 26.36 | 0.924 | 26.24 | 0.916 | 28.03 | 0.949 | 207,224 |
| With 3D LuT | 27.61 | 0.935 | 30.02 | 0.935 | 27.78 | 0.937 | 27.52 | 0.932 | 29.45 | 0.968 | 211,217 |

Table 19. Ablation on the effect of using the 3D LuT in our photofinishing module for the default style (Style #0) on the S24 dataset [16]. The input consists of pseudo ground-truth denoised raw images mapped to the linear sRGB space and downsampled to one-quarter of the original resolution, and the ground truth is the corresponding sRGB image at the same resolution. The 3D LuT mainly benefits target picture styles with *artistic* appearance, whereas for simpler styles that mainly enhance color and overall tone, the gain is marginal. The best results are highlighted in yellow.

| Photofinishing Module | S24 Test Set (1/4 LsRGB/sRGB) | |
|---|---|---|
| | PSNR↑ | SSIM↑ |
| Without 3D LuT | 27.49 | 0.939 |
| With 3D LuT | 27.53 | 0.935 |

Table 20. Alternative designs using HDRNet [39] as a replacement for our local tone-mapping network and operator, or by substituting the entire photofinishing module with a single HDRNet. Results are reported on the S24 test set [16], where the input consists of pseudo ground-truth denoised raw images mapped to the linear sRGB space and downsampled to one-quarter of the original resolution, and the ground truth is the corresponding sRGB image at the same resolution. The best results are highlighted in yellow.

| Method | S24 Test Set (1/4 LsRGB/sRGB) | | |
|---|---|---|---|
| | PSNR↑ | SSIM↑ | # params |
| LTM → HDRNet | 25.35 | 0.920 | 570,462 |
| Photofinishing → HDRNet | 24.87 | 0.911 | 483,453 |
| Ours | 27.49 | 0.939 | 207,224 |

## K.2. Additional Results

In this subsection, we provide additional results, including visual outputs of intermediate stages, quantitative results on an additional dataset, detailed evaluations across the S24 dataset styles, and further qualitative examples.

### K.2.1. Visualization of Photofinishing Stages

Figure 44 visualizes intermediate stages of our pipeline (after denoising and color correction, digital gain, tone mapping, chroma mapping, gamma correction, and detail enhancement) using the default picture style (Style #0). Another example of the default style is shown in Fig. 45. Figure 46 shows two examples rendered with different artistic styles (Styles #3 and #4), illustrating how the module's internal representations adapt across styles. Figure 47 compares the intermediate outputs of the same scene across dif-

Table 21. Results on the MIT-Adobe 5K dataset [26]. We report PSNR, SSIM [81], LPIPS [89], and ΔE 2000 [72], along with the total number of parameters for each method. The best results are highlighted in yellow, and the second best in green.

| Method | Adobe 5K Test Set | | | | |
|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | ΔE 2000↓ | # params |
| ISPDiffuser [70] | 19.44 | 0.711 | 0.333 | 11.030 | 20,938,890 |
| PyNet [48] | 19.18 | 0.695 | 0.219 | 12.300 | 47,548,170 |
| CIE XYZ Net [10] | 20.25 | 0.715 | 0.189 | 10.352 | 1,348,789 |
| Invertible-ISP [85] | 20.06 | 0.693 | 0.205 | 10.800 | 1,413,760 |
| LAN [69] | 16.46 | 0.618 | 0.275 | 14.817 | 46,847 |
| MicroISP [49] | 19.84 | 0.708 | 0.217 | 12.188 | 13,560 |
| ParamISP [53] | 21.17 | 0.735 | 0.184 | 9.729 | 1,420,000 |
| LiteISP [91] | 21.45 | 0.740 | 0.169 | 9.687 | 9,094,000 |
| FourierISP [44] | 20.29 | 0.737 | 0.208 | 10.582 | 7,589,736 |
| Ours (lite, w/o enhancement) | 21.10 | 0.733 | 0.223 | 11.021 | 452,447 |
| Ours (base, w/o enhancement) | 21.10 | 0.733 | 0.226 | 11.020 | 1,139,907 |
| Ours (large, w/o enhancement) | 21.11 | 0.736 | 0.225 | 11.005 | 3,841,547 |
| Ours (lite, w/ enhancement) | 21.27 | 0.739 | 0.219 | 10.765 | 503,082 |
| Ours (base, w/ enhancement) | 21.28 | 0.740 | 0.221 | 10.762 | 1,190,542 |
| Ours (large, w/ enhancement) | 21.29 | 0.742 | 0.220 | 10.743 | 3,892,182 |

ferent picture styles, highlighting how the tone and color transformations learned by each network component vary with style.

### K.2.2. Results on MIT-Adobe 5K Dataset

In the main paper, we reported our results on the S24 dataset, which includes multiple picture styles, including artistic styles. To further evaluate our method against other alternatives, we trained and tested it, along with competing methods, on the MIT-Adobe FiveK dataset [26], using Expert C as the ground truth. For our method, we generated pseudo ground-truth denoised images following the same procedure used for the S24 dataset, employing the AI-based denoiser available in Adobe Lightroom. Further details on how we trained and evaluated other methods on this dataset are provided in Sec. L. The results are reported in Table 21. Our method achieves promising results, ranking second in PSNR and first in SSIM, while maintaining a fully modular design with controllable rendering stages. Moreover, the lite version of our method surpasses several recent methods despite having significantly fewer parameters (21.27 dB with ∼0.5 M parameters vs. ISPDiffuser [70] at 19.44 dB with ∼21 M parameters, and ParamISP [53] at 21.17 dB with ∼1.4 M parameters).

### K.2.3. Re-Rendering Results with Embedded Raw

In Sec. I.10, we described our approach for embedding JPEG-compressed raw data alongside the rendered sRGB

Figure 44. Outputs of the intermediate stages in our pipeline, including raw denoising, color correction (CC), and the photofinishing module, which comprises digital gain, global tone mapping (GTM), local tone mapping (LTM), chroma mapping, and detail enhancement (enh.). We show enlarged patches of the final result when raw denoising is disabled, when detail enhancement is disabled, and when both are enabled. The example is taken from the S24 test set [16].
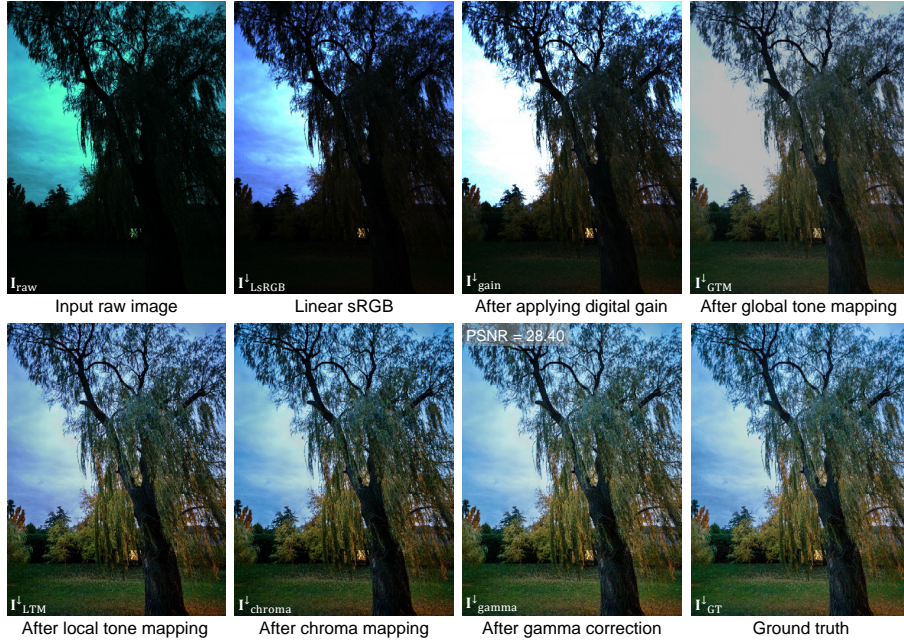


Figure 45. Intermediate outputs of our photofinishing module for the default picture style (Style #0). The example is from the S24 test set [16]. The input to the photofinishing module is the pseudo ground-truth denoised raw image mapped to the linear sRGB space at one-quarter of the original resolution, and the ground-truth reference is the corresponding Style #0 sRGB image at the same resolution.

image within the saved file. This design enables an unlimited number of post-capture re-rendering operations without cumulative degradation in accuracy or the need to reconstruct raw data from the rendered sRGB image, while introducing only a moderate increase in the final JPEG file size compared to alternatives such as iPhone's HEIC format for picture styles or saving the original DNG file.

Here, we present a quantitative comparison against alternative methods that enable image re-rendering through raw reconstruction.

Specifically, we compare our embedded-raw approach (Sec. I.10) with InvISP [85] and ParamISP [53], both of which reconstruct raw data from the rendered image to allow re-rendering with different configurations (i.e., using different model weights corresponding to distinct picture styles). Our evaluation scheme is as follows. Each test image in the S24 dataset [16] is first rendered into a target style (referred to here as the *source picture style*) using InvISP or ParamISP, followed by raw reconstruction and subsequent re-rendering into the remaining five picture styles using the models trained for those *target picture styles*.

For our method, we directly extract and decode the embedded raw data, then re-render it using our pipeline with the photofinishing and detail-enhancement models of the target picture style. In this experiment, we used the Raw-JPEG Adapter [15] with a raw-JPEG quality setting of 75, which introduces approximately 1-2 MB of overhead for the embedded encoded raw data. As shown in Tables 22-27, our method achieves consistently superior accuracy and avoids the variability observed in other methods, whose re-
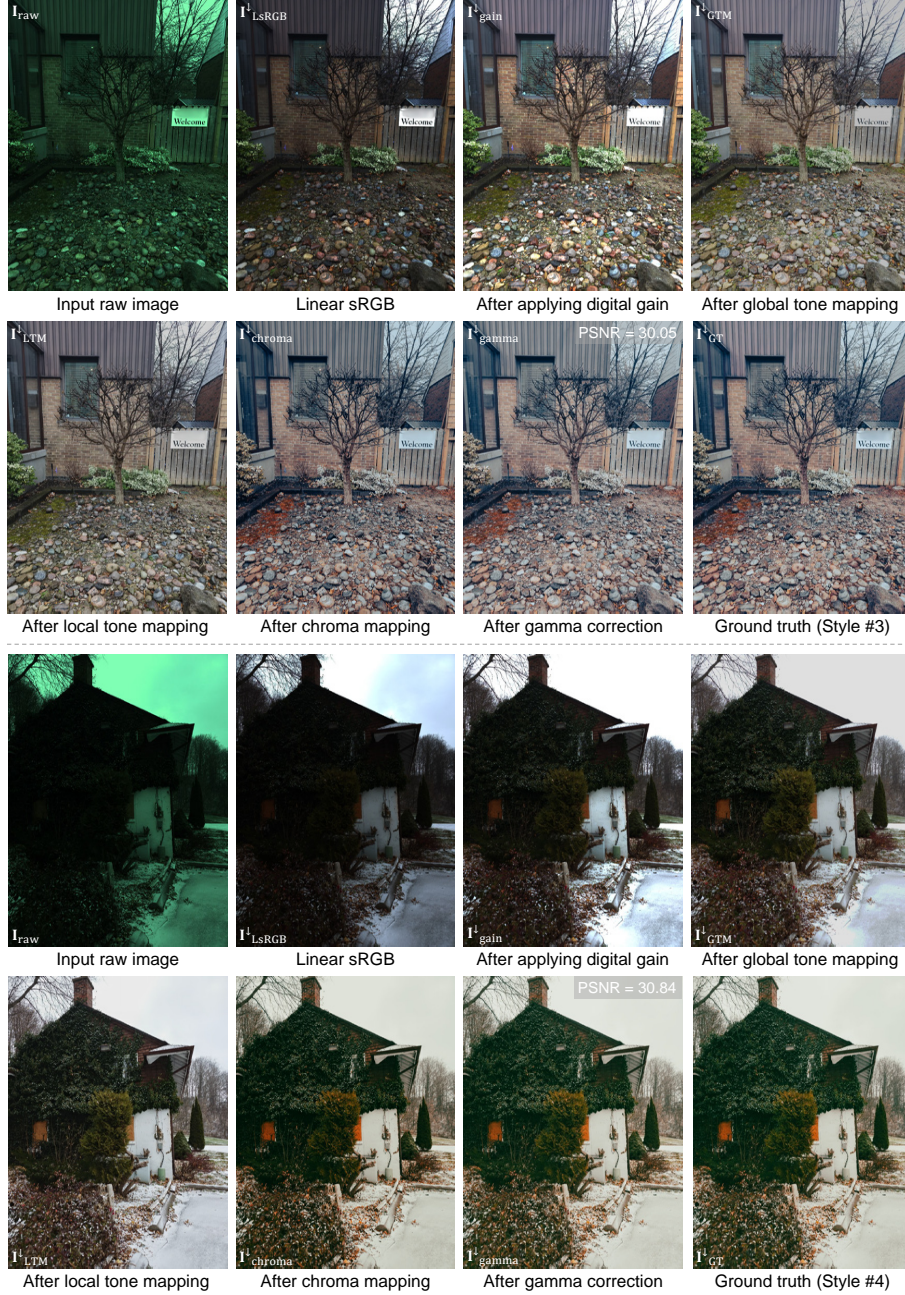
Figure 46. Intermediate outputs of our photofinishing module for two artistic picture styles (Style #3 and Style #4). Each example is taken from the S24 test set [16]. Inputs are pseudo ground-truth denoised raw images mapped to linear sRGB space at one-quarter of the original resolution, and the corresponding sRGB ground-truth references are shown at the same resolution.

rendering quality depends on the source picture style of the initially rendered sRGB image.

### K.2.4. Detailed Results on S24 Styles #1–5

In the main paper, we reported the PSNR scores of our method compared to alternative approaches across different picture styles available in the S24 dataset. Tables 28– 32 provide a more comprehensive evaluation, including

SSIM [81], LPIPS [89], and $\Delta$E 2000 [72]. The detailed results for Styles #1–5.

### K.2.5. Additional Qualitative Results

We present additional qualitative results obtained using our method. Figure 48 shows a qualitative comparison between our method and recent state-of-the-art neural ISP methods on the S24 test set [16], along with the PSNR results of

38

Figure 47. Comparison of intermediate-stage outputs for the same scene rendered with different picture styles. Images are taken from the S24 test set [16]. The examples illustrate how the photofinishing stages adapt their tone and chroma adjustments depending on the target picture style.

each method for every style. As shown, our method consistently produces high-quality results while using fewer parameters overall, since only the photofinishing and detail-enhancement modules employ style-specific weights, owing to the modularity of our design. Lastly, Figure 49 presents additional results on images captured by different iPhone devices. None of our training data includes iPhone images, further demonstrating the cross-device generalization ability of our method.
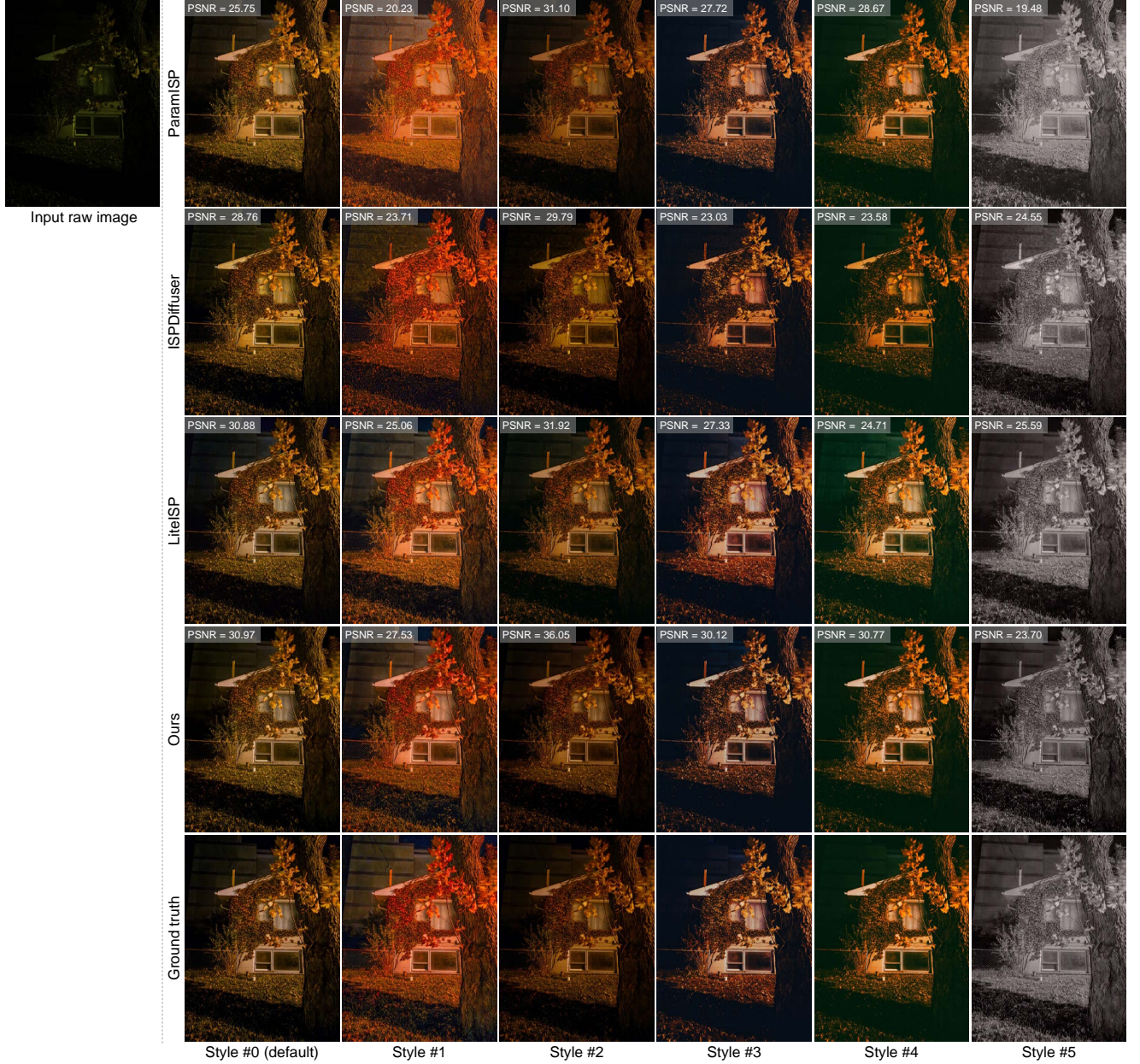
Figure 48. Additional qualitative comparison between our method and recent neural ISP methods (ISPDiffuser [70], LiteISP [91], and ParamISP [53]). Results are shown for the default style of the S24 dataset [16] (Style #0) and the remaining artistic picture styles (Styles #1–5). The shown images are from the S24 test set.

## L. Evaluation Details

Except for the reported results on the Zurich Raw-to-sRGB dataset [48] (Sec. B.2), we trained and evaluated all other methods ourselves, following their original papers and released codebases. We generally followed the training procedures and hyperparameters described in the original works, with minor modifications where necessary (as elaborated in the following). For the S24 dataset [16], only demosaiced linear RGB raw images are available. For methods requiring single-channel mosaiced inputs, we simulated the mosaicing process by applying a fixed Bayer pattern to the RGB images. For three-channel input methods (such as ours), the raw input images from the S24 [16], MIT-Adobe FiveK [26], and Zurich Raw-to-sRGB [48] datasets were generated by applying black-level normalization (using black and saturation levels specified in the DNG files or dataset metadata), followed by demosaicing with the Menon algorithm [65] when RGB demosaiced raw images were not provided.

Figure 49. Additional results produced by our method. The shown images were rendered from raw captures taken with different iPhone devices (13, 13 Pro Max, 14, and 15) and across various camera modules (main, ultra-wide, and telephoto). For these results, we used the photofinishing and detail-enhancement models trained on the S24 dataset [16], along with the generic denoising and cross-camera AWB models. None of the models were trained on data from any iPhone device.

Table 22. Re-rendering results for target picture style #0 on the S24 dataset [16]. We report PSNR, SSIM [81], LPIPS [89], and ΔE 2000 [72]. For Invertible-ISP [85] and ParamISP [53], images were first rendered to sRGB using the source picture style (as indicated in the table) and then reconstructed to raw before re-rendering to the target style. For our method, the embedded raw data were used directly for re-rendering, with the base denoising model applied. The best results are highlighted in yellow.

| Method | S24 Test Set (Target: Style #0) | | | |
| --- | --- | --- | --- | --- |
| | PSNR↑ | SSIM↑ | LPIPS↓ | ΔE 2000↓ |
| InvISP (source: Style #1) | 22.69 | 0.817 | 0.152 | 7.319 |
| InvISP (source: Style #2) | 22.51 | 0.822 | 0.154 | 7.419 |
| InvISP (source: Style #3) | 22.74 | 0.812 | 0.160 | 7.438 |
| InvISP (source: Style #4) | 22.50 | 0.809 | 0.160 | 7.550 |
| InvISP (source: Style #5) | 22.82 | 0.801 | 0.185 | 7.882 |
| ParamISP (source: Style #1) | 21.71 | 0.789 | 0.181 | 8.072 |
| ParamISP (source: Style #2) | 22.67 | 0.805 | 0.154 | 7.478 |
| ParamISP (source: Style #3) | 21.50 | 0.786 | 0.217 | 9.844 |
| ParamISP (source: Style #4) | 22.40 | 0.798 | 0.170 | 7.937 |
| ParamISP (source: Style #5) | 16.88 | 0.730 | 0.402 | 16.532 |
| Ours (base denoiser) | 26.90 | 0.901 | 0.066 | 4.256 |

Table 23. Re-rendering results for target picture style #1 on the S24 dataset [16]. For Invertible-ISP [85] and ParamISP [53], images were first rendered to sRGB using the source picture style and then reconstructed to raw before re-rendering. For our method, the embedded raw data were used directly for re-rendering, with the base denoising model applied. The best results are highlighted in yellow.

| Method | S24 Test Set (Target: Style #1) | | | |
| --- | --- | --- | --- | --- |
| | PSNR↑ | SSIM↑ | LPIPS↓ | ΔE 2000↓ |
| InvISP (source: Style #0) | 21.38 | 0.810 | 0.211 | 8.676 |
| InvISP (source: Style #2) | 23.27 | 0.833 | 0.164 | 6.509 |
| InvISP (source: Style #3) | 23.46 | 0.830 | 0.166 | 6.530 |
| InvISP (source: Style #4) | 23.29 | 0.824 | 0.165 | 6.609 |
| InvISP (source: Style #5) | 22.85 | 0.828 | 0.181 | 7.881 |
| ParamISP (source: Style #0) | 23.11 | 0.814 | 0.169 | 7.065 |
| ParamISP (source: Style #2) | 22.80 | 0.814 | 0.166 | 7.135 |
| ParamISP (source: Style #3) | 21.84 | 0.796 | 0.215 | 8.687 |
| ParamISP (source: Style #4) | 22.74 | 0.810 | 0.180 | 7.517 |
| ParamISP (source: Style #5) | 17.36 | 0.754 | 0.374 | 14.479 |
| Ours (base denoiser) | 26.68 | 0.897 | 0.088 | 4.417 |
| Ours (base denoiser + 3D LuT) | 26.04 | 0.890 | 0.100 | 4.989 |

Table 24. Re-rendering results for target picture style #2 on the S24 dataset [16]. The best results are highlighted in yellow.

| Method | S24 Test Set (Target: Style #2) | | | |
| --- | --- | --- | --- | --- |
| | PSNR↑ | SSIM↑ | LPIPS↓ | ΔE 2000↓ |
| InvISP (source: Style #0) | 26.17 | 0.843 | 0.124 | 5.301 |
| InvISP (source: Style #1) | 26.14 | 0.853 | 0.121 | 5.260 |
| InvISP (source: Style #3) | 26.23 | 0.851 | 0.129 | 5.280 |
| InvISP (source: Style #4) | 26.19 | 0.847 | 0.125 | 5.283 |
| InvISP (source: Style #5) | 25.91 | 0.841 | 0.152 | 5.611 |
| ParamISP (source: Style #0) | 25.67 | 0.843 | 0.131 | 5.647 |
| ParamISP (source: Style #1) | 24.73 | 0.832 | 0.157 | 6.179 |
| ParamISP (source: Style #3) | 24.09 | 0.827 | 0.179 | 7.455 |
| ParamISP (source: Style #4) | 24.94 | 0.834 | 0.148 | 6.423 |
| ParamISP (source: Style #5) | 19.86 | 0.743 | 0.330 | 11.855 |
| Ours (base denoiser) | 29.48 | 0.923 | 0.067 | 3.605 |
| Ours (base denoiser + 3D LuT) | 28.74 | 0.915 | 0.082 | 4.216 |

For the S24 and Zurich Raw-to-sRGB datasets, training was performed using each dataset's respective training

Table 25. Re-rendering results for target picture style #3 on the S24 dataset [16]. The best results are highlighted in yellow.

| Method | S24 Test Set (Target: Style #3) | | | |
| --- | --- | --- | --- | --- |
| | PSNR↑ | SSIM↑ | LPIPS↓ | ΔE 2000↓ |
| InvISP (source: Style #0) | 23.82 | 0.851 | 0.146 | 7.316 |
| InvISP (source: Style #1) | 23.67 | 0.853 | 0.149 | 7.327 |
| InvISP (source: Style #2) | 23.56 | 0.858 | 0.147 | 7.394 |
| InvISP (source: Style #4) | 23.64 | 0.852 | 0.151 | 7.413 |
| InvISP (source: Style #5) | 23.39 | 0.840 | 0.181 | 7.651 |
| ParamISP (source: Style #0) | 23.31 | 0.834 | 0.151 | 7.126 |
| ParamISP (source: Style #1) | 22.49 | 0.821 | 0.179 | 7.781 |
| ParamISP (source: Style #2) | 23.15 | 0.836 | 0.150 | 7.267 |
| ParamISP (source: Style #4) | 22.89 | 0.828 | 0.160 | 7.469 |
| ParamISP (source: Style #5) | 18.64 | 0.768 | 0.303 | 12.829 |
| Ours (base denoiser) | 26.89 | 0.905 | 0.085 | 4.660 |
| Ours (base denoiser + 3D LuT) | 26.35 | 0.897 | 0.101 | 5.318 |

Table 26. Re-rendering results for target picture style #4 on the S24 dataset [16]. The best results are highlighted in yellow.

| Method | S24 Test Set (Target: Style #4) | | | |
| --- | --- | --- | --- | --- |
| | PSNR↑ | SSIM↑ | LPIPS↓ | ΔE 2000↓ |
| InvISP (source: Style #0) | 23.30 | 0.842 | 0.149 | 7.647 |
| InvISP (source: Style #1) | 23.24 | 0.843 | 0.148 | 7.686 |
| InvISP (source: Style #2) | 23.27 | 0.846 | 0.146 | 7.605 |
| InvISP (source: Style #3) | 23.31 | 0.841 | 0.151 | 7.639 |
| InvISP (source: Style #5) | 22.85 | 0.828 | 0.181 | 7.881 |
| ParamISP (source: Style #0) | 22.90 | 0.824 | 0.141 | 6.585 |
| ParamISP (source: Style #1) | 22.01 | 0.806 | 0.171 | 7.398 |
| ParamISP (source: Style #2) | 22.88 | 0.825 | 0.140 | 6.622 |
| ParamISP (source: Style #3) | 22.09 | 0.811 | 0.167 | 7.441 |
| ParamISP (source: Style #5) | 18.64 | 0.744 | 0.323 | 12.482 |
| Ours (base denoiser) | 26.44 | 0.897 | 0.080 | 4.428 |
| Ours (base denoiser + 3D LuT) | 26.21 | 0.893 | 0.087 | 4.692 |

Table 27. Re-rendering results for target picture style #5 on the S24 dataset [16]. The best results are highlighted in yellow.

| Method | S24 Test Set (Target: Style #5) | | | |
| --- | --- | --- | --- | --- |
| | PSNR↑ | SSIM↑ | LPIPS↓ | ΔE 2000↓ |
| InvISP (source: Style #0) | 24.89 | 0.870 | 0.168 | 5.844 |
| InvISP (source: Style #1) | 24.91 | 0.870 | 0.165 | 5.840 |
| InvISP (source: Style #2) | 24.82 | 0.872 | 0.162 | 5.902 |
| InvISP (source: Style #3) | 24.83 | 0.863 | 0.175 | 5.941 |
| InvISP (source: Style #4) | 24.78 | 0.863 | 0.174 | 5.960 |
| ParamISP (source: Style #0) | 24.47 | 0.841 | 0.160 | 4.963 |
| ParamISP (source: Style #1) | 24.20 | 0.834 | 0.180 | 5.122 |
| ParamISP (source: Style #2) | 24.50 | 0.847 | 0.153 | 5.000 |
| ParamISP (source: Style #3) | 23.95 | 0.833 | 0.170 | 5.088 |
| ParamISP (source: Style #4) | 24.18 | 0.840 | 0.159 | 5.098 |
| Ours (base denoiser) | 27.75 | 0.916 | 0.102 | 3.537 |
| Ours (base denoiser + 3D LuT) | 28.27 | 0.921 | 0.094 | 3.332 |

split. For the MIT-Adobe FiveK dataset, we followed the train/validation/test split provided in [47]; specifically, the split consists of 493 validation images, 989 test images, and the remaining images are used for training.

Since the MIT-Adobe FiveK dataset [26] (Sec. K.2.2) contains images captured with multiple DSLR cameras, training and evaluating black-box neural ISP methods can often be unstable. To ensure fairness and stability, we trained and tested other methods on linear sRGB inputs, obtained by converting the raw images using the illuminant and CCM provided in the DNG metadata. For methods that require single-channel mosaiced inputs, we used mosaiced

Table 28. Detailed results for Style #1 on the S24 test set [16]. Our method is evaluated with different denoising model capacities (lite, base, and large), with and without the enhancement network and optional 3D LuT. The best results are highlighted in yellow and the second best in green.

| Method | S24 Test Set (Style #1) | | | |
|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | ΔE 2000↓ |
| ISPDiffuser [70] | 25.60 | 0.910 | 0.117 | 5.100 |
| PyNet [48] | 24.36 | 0.875 | 0.097 | 5.759 |
| CIE XYZ Net [10] | 22.40 | 0.859 | 0.176 | 8.291 |
| PP (cmKAN) [67] | 20.93 | 0.875 | 0.139 | 8.436 |
| FlexISP [60] | 24.86 | 0.891 | 0.108 | 6.150 |
| Invertible-ISP [85] | 23.48 | 0.832 | 0.157 | 6.642 |
| LAN [69] | 22.98 | 0.812 | 0.126 | 6.669 |
| MicroISP [49] | 20.30 | 0.747 | 0.183 | 11.298 |
| ParamISP [53] | 24.97 | 0.856 | 0.123 | 5.724 |
| LiteISP [91] | 26.66 | 0.915 | 0.067 | 4.702 |
| FourierISP [44] | 25.19 | 0.925 | 0.099 | 5.432 |
| Ours (lite, w/o enhancement) | 25.16 | 0.866 | 0.124 | 5.609 |
| Ours (base, w/o enhancement) | 25.28 | 0.873 | 0.116 | 5.483 |
| Ours (large, w/o enhancement) | 25.31 | 0.874 | 0.114 | 5.463 |
| Ours (lite, w/o enhancement, w/ 3D LuT) | 26.39 | 0.874 | 0.086 | 4.486 |
| Ours (base, w/o enhancement, w/ 3D LuT) | 26.52 | 0.880 | 0.078 | 4.348 |
| Ours (large, w/o enhancement, w/ 3D LuT) | 26.56 | 0.882 | 0.076 | 4.329 |
| Ours (lite, w/ enhancement, w/ 3D LuT) | 26.56 | 0.906 | 0.092 | 4.807 |
| Ours (base, w/ enhancement, w/ 3D LuT) | 26.71 | 0.914 | 0.085 | 4.692 |
| Ours (large, w/ enhancement, w/ 3D LuT) | 26.75 | 0.916 | 0.083 | 4.668 |

Table 29. Detailed results for Style #2 on the S24 test set [16]. The best results are highlighted in yellow and the second best in green.

| Method | S24 Test Set (Style #2) | | | |
|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | ΔE 2000↓ |
| ISPDiffuser [70] | 27.30 | 0.912 | 0.116 | 4.600 |
| PyNet [48] | 25.94 | 0.890 | 0.095 | 5.432 |
| CIE XYZ Net [10] | 24.05 | 0.872 | 0.120 | 6.634 |
| PP (cmKAN) [67] | 23.04 | 0.875 | 0.131 | 7.613 |
| FlexISP [60] | 27.47 | 0.911 | 0.089 | 4.554 |
| Invertible-ISP [85] | 26.35 | 0.861 | 0.115 | 5.229 |
| LAN [69] | 23.74 | 0.782 | 0.112 | 6.436 |
| MicroISP [49] | 23.66 | 0.801 | 0.155 | 8.885 |
| ParamISP [53] | 27.11 | 0.869 | 0.101 | 4.947 |
| LiteISP [91] | 28.33 | 0.922 | 0.064 | 4.284 |
| FourierISP [44] | 28.03 | 0.928 | 0.082 | 4.488 |
| Ours (lite, w/o enhancement) | 28.09 | 0.890 | 0.094 | 4.589 |
| Ours (base, w/o enhancement) | 28.19 | 0.893 | 0.090 | 4.516 |
| Ours (large, w/o enhancement) | 28.22 | 0.893 | 0.089 | 4.503 |
| Ours (lite, w/o enhancement, w/ 3D LuT) | 29.21 | 0.898 | 0.068 | 3.641 |
| Ours (base, w/o enhancement, w/ 3D LuT) | 29.29 | 0.901 | 0.065 | 3.594 |
| Ours (large, w/o enhancement, w/ 3D LuT) | 29.31 | 0.901 | 0.064 | 3.588 |
| Ours (lite, w/ enhancement, w/ 3D LuT) | 28.92 | 0.920 | 0.078 | 4.153 |
| Ours (base, w/ enhancement, w/ 3D LuT) | 28.99 | 0.923 | 0.075 | 4.135 |
| Ours (large, w/ enhancement, w/ 3D LuT) | 29.01 | 0.924 | 0.075 | 4.130 |

Table 30. Detailed results for Style #3 on the S24 test set [16]. The best results are highlighted in yellow and the second best in green.

| Method | S24 Test Set (Style #3) | | | |
|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | ΔE 2000↓ |
| ISPDiffuser [70] | 25.02 | 0.899 | 0.125 | 5.600 |
| PyNet [48] | 24.70 | 0.882 | 0.096 | 5.999 |
| CIE XYZ Net [10] | 22.00 | 0.854 | 0.162 | 8.906 |
| PP (cmKAN) [67] | 21.85 | 0.876 | 0.130 | 8.077 |
| FlexISP [60] | 25.23 | 0.901 | 0.106 | 5.807 |
| Invertible-ISP [85] | 23.84 | 0.852 | 0.144 | 7.347 |
| LAN [69] | 23.47 | 0.830 | 0.118 | 6.902 |
| MicroISP [49] | 21.45 | 0.792 | 0.169 | 10.751 |
| ParamISP [53] | 24.77 | 0.872 | 0.121 | 6.417 |
| LiteISP [91] | 26.31 | 0.913 | 0.073 | 5.220 |
| FourierISP [44] | 25.38 | 0.919 | 0.100 | 5.703 |
| Ours (lite, w/o enhancement) | 25.66 | 0.877 | 0.113 | 5.930 |
| Ours (base, w/o enhancement) | 25.73 | 0.881 | 0.110 | 5.894 |
| Ours (large, w/o enhancement) | 25.75 | 0.882 | 0.110 | 5.883 |
| Ours (lite, w/o enhancement, w/ 3D LuT) | 26.69 | 0.884 | 0.085 | 4.727 |
| Ours (base, w/o enhancement, w/ 3D LuT) | 26.79 | 0.889 | 0.081 | 4.675 |
| Ours (large, w/o enhancement, w/ 3D LuT) | 26.83 | 0.890 | 0.080 | 4.665 |
| Ours (lite, w/ enhancement, w/ 3D LuT) | 26.78 | 0.913 | 0.090 | 5.209 |
| Ours (base, w/ enhancement, w/ 3D LuT) | 26.78 | 0.913 | 0.090 | 5.213 |
| Ours (large, w/ enhancement, w/ 3D LuT) | 26.83 | 0.914 | 0.089 | 5.196 |

Table 31. Detailed results for Style #4 on the S24 test set [16]. The best results are highlighted in yellow and the second best in green.

| Method | S24 Test Set (Style #4) | | | |
|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | ΔE 2000↓ |
| ISPDiffuser [70] | 25.93 | 0.904 | 0.083 | 5.100 |
| PyNet [48] | 24.34 | 0.876 | 0.094 | 5.745 |
| CIE XYZ Net [10] | 22.26 | 0.846 | 0.147 | 7.624 |
| PP (cmKAN) [67] | 20.91 | 0.852 | 0.140 | 8.394 |
| FlexISP [60] | 23.96 | 0.878 | 0.108 | 5.965 |
| Invertible-ISP [85] | 23.33 | 0.842 | 0.145 | 7.694 |
| LAN [69] | 22.80 | 0.790 | 0.116 | 6.736 |
| MicroISP [49] | 20.34 | 0.750 | 0.175 | 11.499 |
| ParamISP [53] | 24.18 | 0.853 | 0.118 | 6.138 |
| LiteISP [91] | 25.04 | 0.894 | 0.082 | 5.517 |
| FourierISP [44] | 24.74 | 0.906 | 0.100 | 5.591 |
| Ours (lite, w/o enhancement) | 25.47 | 0.868 | 0.107 | 5.486 |
| Ours (base, w/o enhancement) | 25.55 | 0.872 | 0.104 | 5.437 |
| Ours (large, w/o enhancement) | 25.58 | 0.873 | 0.104 | 5.424 |
| Ours (lite, w/o enhancement, w/ 3D LuT) | 26.35 | 0.880 | 0.079 | 4.414 |
| Ours (base, w/o enhancement, w/ 3D LuT) | 26.47 | 0.884 | 0.076 | 4.365 |
| Ours (large, w/o enhancement, w/ 3D LuT) | 26.51 | 0.885 | 0.075 | 4.355 |
| Ours (lite, w/ enhancement, w/ 3D LuT) | 26.66 | 0.905 | 0.081 | 4.601 |
| Ours (base, w/ enhancement, w/ 3D LuT) | 26.79 | 0.910 | 0.077 | 4.564 |
| Ours (large, w/ enhancement, w/ 3D LuT) | 26.84 | 0.911 | 0.075 | 4.553 |

data from the MIT-Adobe FiveK dataset only after applying white balancing and color correction.

For FlexISP [60], which employs a three-stage framework (denoising, white balance, and Bayer-to-sRGB mapping), we first trained the denoising stage on the S24 dataset using paired noisy raw and denoised ground-truth images, and then fine-tuned it for each picture style in the S24 dataset, following the training scheme described in the original FlexISP paper [60].

Among all methods, ParamISP [53] and ISPDiffuser [70] required special handling during training compared to the default configurations described in their original papers. For ParamISP [53], originally designed for two-stage training (multi-camera pretraining followed by camera-specific fine-tuning), we found that omitting the pretraining stage yielded better results on the S24 dataset. For ISPDiffuser [70], we observed that training on image patches both reduced training time and improved accuracy. In both cases, we report the best-performing configurations we identified.

For the Zurich dataset (Sec. B.2), results are taken from prior publications, except for CIE XYZ Net [10], which we trained and tested ourselves. LAN [69] results are reported from Rawformer [68], while the remaining numbers are taken from [44, 48]. We report quantitative results for all methods using PSNR, SSIM [81], LPIPS [89], and $\Delta E_{2000}$ [72]. For LPIPS, both outputs and ground-truth images are resized to 1024×1024 prior to computation.

## Acknowledgments

Table 32. Detailed results for Style #5 on the S24 test set [16]. The best results are highlighted in <mark>yellow</mark> and the second best in <mark>green</mark>.

| Method | S24 Test Set (Style #5) | | | |
| --- | --- | --- | --- | --- |
| | PSNR↑ | SSIM↑ | LPIPS↓ | ΔE 2000↓ |
| ISPDiffuser [70] | 26.83 | 0.929 | 0.116 | 3.800 |
| PyNet [48] | 26.32 | 0.920 | 0.089 | 3.982 |
| CIE XYZ Net [10] | 24.67 | 0.897 | 0.133 | 5.270 |
| PP (cmKAN) [67] | 21.30 | 0.880 | 0.147 | 7.338 |
| FlexISP [60] | 24.65 | 0.912 | 0.124 | 8.448 |
| Invertible-ISP [85] | 24.90 | 0.875 | 0.163 | 5.927 |
| LAN [69] | 25.38 | 0.878 | 0.108 | 4.662 |
| MicroISP [49] | 22.68 | 0.823 | 0.221 | 10.500 |
| ParamISP [53] | 25.43 | 0.867 | 0.134 | 4.499 |
| LiteISP [91] | 28.07 | 0.935 | **0.071** | 3.434 |
| FourierISP [44] | 27.41 | **0.947** | 0.089 | 3.561 |
| Ours (lite, w/o enhancement) | 27.08 | 0.898 | 0.112 | 3.868 |
| Ours (base, w/o enhancement) | 27.19 | 0.906 | 0.108 | 3.831 |
| Ours (large, w/o enhancement) | 27.23 | 0.908 | 0.107 | 3.820 |
| Ours (lite, w/o enhancement, w/ 3D LuT) | 27.93 | 0.911 | 0.096 | 3.398 |
| Ours (base, w/o enhancement, w/ 3D LuT) | 28.13 | 0.921 | 0.088 | 3.336 |
| Ours (large, w/o enhancement, w/ 3D LuT) | 28.19 | 0.924 | 0.086 | 3.323 |
| Ours (lite, w/ enhancement, w/ 3D LuT) | 28.73 | 0.930 | 0.090 | 3.234 |
| Ours (base, w/ enhancement, w/ 3D LuT) | **28.95** | 0.938 | 0.084 | **3.177** |
| Ours (large, w/ enhancement, w/ 3D LuT) | **29.03** | **0.941** | **0.082** | **3.160** |

iments.

# References

[1] Abdelrahman Abdelhamed, Stephen Lin, and Michael S Brown. A high-quality denoising dataset for smartphone cameras. In *CVPR*, 2018. 2, 3, 15

[2] Abdullah Abuolaim and Michael S Brown. Defocus deblurring using dual-pixel data. In *ECCV*, 2020. 1

[3] Adobe Community. Samsung expert raw DNG output shows green / over-exposed photo in Lightroom. https://community.adobe.com/t5/lightroom-ecosystem-cloud-based-discussions/p-samsung-expert-raw-dng-output-shows-green-over-exposed-photo-in-lightroom/m-p/14714898, 2024. Accessed: 2025-11-07. 2

[4] Adobe Community. Denoise of HDR merge file results in color artifacts. https://community.adobe.com/t5/lightroom-classic-discussions/p-denoise-of-hdr-merge-file-results-in-color-artifacts/td-p/14915059, 2024. Accessed: 2025-11-07. 2

[5] Adobe Systems Inc. Digital negative (DNG) specification. https://helpx.adobe.com/camera-raw/digital-negative.html, 2004–2025. Accessed: 2025-11-05. 21

[6] Mahmoud Afifi and Abdullah Abuolaim. Semi-supervised raw-to-raw mapping. In *BMVC*, 2021. 1, 15, 16

[7] Mahmoud Afifi and Michael S Brown. Sensor-independent illumination estimation for DNN models. In *BMVC*, 2019. 10, 11

[8] Mahmoud Afifi and Michael S Brown. Deep white-balance editing. In *CVPR*, 2020. 21

[9] Mahmoud Afifi, Brian Price, Scott Cohen, and Michael S Brown. When color constancy goes wrong: Correcting improperly white-balanced images. In *CVPR*, 2019. 5

[10] Mahmoud Afifi, Abdelrahman Abdelhamed, Abdullah Abuolaim, Abhijith Punnappurath, and Michael S Brown. CIE XYZ Net: Unprocessing images for low-level computer vision tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):4688–4700, 2021. 2, 6, 7, 26, 27, 36, 43, 44

[11] Mahmoud Afifi, Jonathan T Barron, Chloe LeGendre, Yun-Ta Tsai, and Francois Bleibel. Cross-camera convolutional color constancy. In *ICCV*, 2021. 5, 9, 14, 15, 16, 17, 21

[12] Mahmoud Afifi, Marcus A Brubaker, and Michael S Brown. HistoGAN: Controlling colors of GAN-generated and real images via color histograms. In *CVPR*, 2021. 10, 11

[13] Mahmoud Afifi, Konstantinos G Derpanis, Bjorn Ommer, and Michael S Brown. Learning multi-scale photo exposure correction. In *CVPR*, 2021. 21

[14] Mahmoud Afifi, Zhenhua Hu, and Liang Liang. Optimizing illuminant estimation in dual-exposure HDR imaging. In *ECCV*, 2024. 1

[15] Mahmoud Afifi, Ran Zhang, and Michael S Brown. Raw-JPEG adapter: Efficient raw image compression with JPEG. *arXiv preprint arXiv:2509.19624*, 2025. 6, 25, 37

[16] Mahmoud Afifi, Luxi Zhao, Abhijith Punnappurath, Mohammed A Abdelsalam, Ran Zhang, and Michael S Brown. Time-aware auto white balance in mobile photography. In *ICCV*, 2025. 3, 5, 6, 7, 8, 1, 4, 14, 15, 16, 17, 18, 19, 21, 22, 23, 24, 25, 29, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44

[17] Pesala Bandara. Influencers complain about 'unflattering' iPhone 15 selfie camera. https://petapixel.com/2023/10/04/influencers-complain-about-unflattering-iphone-15-selfie-camera/, 2023. Accessed: 2025-11-07. 2

[18] Jonathan T Barron. Convolutional color constancy. In *ICCV*, 2015. 1

[19] Jonathan T Barron and Ben Poole. The fast bilateral solver. In *ECCV*, 2016. 2, 3, 4, 5

[20] Jonathan T Barron and Yun-Ta Tsai. Fast Fourier color constancy. In *CVPR*, 2017. 3

[21] Jarosław Bernacki. Automatic exposure algorithms for digital photography. *Multimedia Tools and Applications*, 79(19):12751–12776, 2020. 20

[22] Simone Bianco, Arcangelo R Bruna, Filippo Naccari, and Raimondo Schettini. Color correction pipeline optimization for digital cameras. *Journal of Electronic Imaging*, 22(2):023014–023014, 2013. 1

[23] Tim Brooks, Ben Mildenhall, Tianfan Xue, Jiawen Chen, Dillon Sharlet, and Jonathan T Barron. Unprocessing images for learned raw denoising. In *CVPR*, 2019. 1

[24] Michael S Brown. Color processing for digital cameras. *Fundamentals and applications of colour engineering*, pages 81–98, 2023. 1

[25] Gershon Buchsbaum. A spatial processor model for object colour perception. *Journal of the Franklin Institute*, 310(1):1–26, 1980. 5

[26] Vladimir Bychkovsky, Sylvain Paris, Eric Chan, and Frédo Durand. Learning photographic global tonal adjustment with a database of input/output image pairs. In *CVPR*, 2011. 6, 8, 4, 15, 19, 36, 40, 42

[27] DL Cade. Smartphones are ruining wildfire sky photos with auto white balance. https://petapixel.com/2020/09/10/smartphones-are-ruining-wildfire-sky-photos-with-auto-white-balance/, 2020. Accessed: 2025-11-07. 2

[28] Chen Chen, Qifeng Chen, Jia Xu, and Vladlen Koltun. Learning to see in the dark. In *CVPR*, 2018. 2

[29] Jiawen Chen, Andrew Adams, Neal Wadhwa, and Samuel W Hasinoff. Bilateral guided upsampling. *ACM Transactions on Graphics (TOG)*, 35(6):1–8, 2016. 5, 11, 14, 29, 30, 31

[30] Liangyu Chen, Xiaojie Chu, Xiangyu Zhang, and Jian Sun. Simple baselines for image restoration. In *ECCV*, 2022. 7, 10, 29

[31] Dongliang Cheng, Dilip K Prasad, and Michael S Brown. Illuminant estimation for color constancy: Why spatial-domain methods work and the role of the color distribution. *Journal of the Optical Society of America A*, 31(5):1049–1058, 2014. 21

[32] Marcos V Conde, Florin Vasluianu, and Radu Timofte. BSRAW: Improving blind raw image super-resolution. In *WACV*, 2024. 1

[33] Marcos V Conde, Zihao Lu, and Radu Timofte. PixTalk: Controlling photorealistic image processing and editing with language. In *ICCV*, 2025. 2

[34] Mauricio Delbracio, Damien Kelly, Michael S Brown, and Peyman Milanfar. Mobile computational photography: A tour. *Annual review of vision science*, 7(1):571–604, 2021. 1

[35] Akshay Dudhane, Syed Waqas Zamir, Salman Khan, Fahad Shahbaz Khan, and Ming-Husan Yang. Burst image restoration and enhancement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 2

[36] Egor Ershov, Alexey Savchik, Illya Semenkov, Nikola Banić, Alexander Belokopytov, Daria Senshina, Karlo Koščević, Marko Subašić, and Sven Lončarić. The Cube++ illumination estimation dataset. *IEEE access*, 8:227511–227527, 2020. 21

[37] Alessandro Foi. Clipped noisy images: Heteroskedastic modeling and practical denoising. *Signal Processing*, 89 (12):2609–2629, 2009. 14

[38] Alessandro Foi, Mejdi Trimeche, Vladimir Katkovnik, and Karen Egiazarian. Practical Poissonian-Gaussian noise modeling and fitting for single-image raw-data. *IEEE Transactions on Image Processing*, 17(10):1737–1754, 2008. 14

[39] 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. 1, 35, 36

[40] Clément Godard, Kevin Matzen, and Matt Uyttendaele. Deep burst denoising. In *ECCV*, 2018. 2

[41] Jeremy Gray. How Apple's Next-Gen photographic styles transform iPhone photography. https://petapixel.com/2024/10/24/how-apples-next-gen-photographic-styles-transform-iphone-photography/, 2024. Accessed: 2025-11-07. 2

[42] Samuel W Hasinoff, Dillon Sharlet, Ryan Geiss, Andrew Adams, Jonathan T Barron, Florian Kainz, Jiawen Chen, and Marc Levoy. Burst photography for high dynamic range and low-light imaging on mobile cameras. *ACM Transactions on Graphics (ToG)*, 35(6):1–12, 2016. 2

[43] Kaiming He, Jian Sun, and Xiaoou Tang. Guided image filtering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(6):1397–1409, 2012. 19, 30, 31

[44] Xuanhua He, Tao Hu, Guoli Wang, Zejin Wang, Run Wang, Qian Zhang, Keyu Yan, Ziyi Chen, Rui Li, Chengjun Xie, et al. Enhancing raw-to-sRGB with decoupled style structure in Fourier domain. In *AAAI*, 2024. 1, 6, 7, 36, 43, 44

[45] Qibin Hou, Daquan Zhou, and Jiashi Feng. Coordinate attention for efficient mobile network design. In *CVPR*, 2021. 7, 32, 33

[46] Yuanming Hu, Baoyuan Wang, and Stephen Lin. FC4: Fully convolutional color constancy with confidence-weighted pooling. In *CVPR*, 2017. 1

[47] 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. 42

[48] Andrey Ignatov, Luc Van Gool, and Radu Timofte. Replacing mobile camera ISP with a single deep learning model. In *CVPRW*, 2020. 1, 6, 7, 8, 3, 4, 36, 40, 43, 44

[49] Andrey Ignatov, Anastasia Sycheva, Radu Timofte, Yu Tseng, Yu-Syuan Xu, Po-Hsiang Yu, Cheng-Ming Chiang, Hsien-Kai Kuo, Min-Hung Chen, Chia-Ming Cheng, et al. MicroISP: Processing 32MP photos on mobile devices with deep learning. In *ECCV*, 2022. 1, 6, 7, 36, 43, 44

[50] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 7

[51] Eric Kee, Adam Pikielny, Kevin Blackburn-Matzen, and Marc Levoy. Removing reflections from raw photos. In *CVPR*, 2025. 1

[52] Dongyoung Kim, Mahmoud Afifi, Dongyun Kim, Michael S Brown, and Seon Joo Kim. CCMNet: Leveraging calibrated color correction matrices for cross-camera color constancy. In *ICCV*, 2025. 3, 21

[53] Woohyeok Kim, Geonu Kim, Junyong Lee, Seungyong Lee, Seung-Hwan Baek, and Sunghyun Cho. ParamISP: Learned forward and inverse ISPs using camera parameters. In *CVPR*, 2024. 2, 6, 7, 36, 37, 40, 42, 43, 44

[54] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 13, 14, 27

[55] Marc Levoy and Florian Kainz. Project Indigo: A computational photography camera app. https://research.adobe.com/articles/indigo/indigo.html, 2025. Accessed: 2025-10-15. 7, 8

[56] Leyi Li, Huijie Qiao, Qi Ye, and Qinmin Yang. Metadata-based raw reconstruction via implicit neural functions. In *CVPR*, 2023. 6

[57] Chih-Hung Liang, Yu-An Chen, Yueh-Cheng Liu, and Winston H Hsu. Raw image deblurring. *IEEE Transactions on Multimedia*, 24:61–72, 2020. 1

[58] Zhetong Liang, Jianrui Cai, Zisheng Cao, and Lei Zhang. CameraNet: A two-stage framework for effective camera

ISP learning. *IEEE Transactions on Image Processing*, 30: 2248–2262, 2021. 2

[59] Orly Liba, Kiran Murthy, Yun-Ta Tsai, Tim Brooks, Tianfan Xue, Nikhil Karnad, Qiurui He, Jonathan T Barron, Dillon Sharlet, Ryan Geiss, et al. Handheld mobile photography in very low light. *ACM Transactions on Graphics (TOG)*, 38 (6):164–1, 2019. 1

[60] Shuai Liu, Chaoyu Feng, Xiaotao Wang, Hao Wang, Ran Zhu, Yongqiang Li, and Lei Lei. Deep-FlexISP: A three-stage framework for night photography rendering. In *CVPRW*, 2022. 2, 6, 7, 43, 44

[61] Xinhao Liu, Masayuki Tanaka, and Masatoshi Okutomi. Practical signal-dependent noise parameter estimation from a single noisy image. *IEEE Transactions on Image Processing*, 23(10):4361–4371, 2014. 14

[62] Yi-Chen Lo, Chia-Che Chang, Hsuan-Chao Chiu, Yu-Hao Huang, Chia-Ping Chen, Yu-Lin Chang, and Kevin Jou. CLCC: Contrastive learning for color constancy. In *CVPR*, 2021. 3

[63] Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016. 13, 14

[64] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 13

[65] Daniele Menon, Stefano Andriani, and Giancarlo Calvagno. Demosaicing with directional filtering and a posteriori decision. *IEEE Transactions on Image Processing*, 16(1):132–141, 2007. 18, 40

[66] Ben Mildenhall, Jonathan T Barron, Jiawen Chen, Dillon Sharlet, Ren Ng, and Robert Carroll. Burst denoising with kernel prediction networks. In *CVPR*, 2018. 2

[67] Artem Nikonorov, Georgy Perevozchikov, Andrei Korepanov, Nancy Mehta, Mahmoud Afifi, Egor Ershov, and Radu Timofte. Color matching using hypernetwork-based Kolmogorov-Arnold networks. 2025. 7, 43, 44

[68] Georgy Perevozchikov, Nancy Mehta, Mahmoud Afifi, and Radu Timofte. Rawformer: Unpaired raw-to-raw translation for learnable camera ISPs. In *ECCV*, 2024. 1, 43

[69] Daniel Wirzberger Raimundo, Andrey Ignatov, and Radu Timofte. LAN: Lightweight attention-based network for raw-to-RGB smartphone image processing. In *CVPRW*, 2022. 1, 6, 7, 36, 43, 44

[70] Yang Ren, Hai Jiang, Menglong Yang, Wei Li, and Shuaicheng Liu. ISPDiffuser: Learning raw-to-sRGB mappings with texture-aware diffusion models and histogram-guided color consistency. In *AAAI*, 2025. 1, 6, 7, 15, 17, 19, 36, 40, 43, 44

[71] Alan R Robertson. The CIE 1976 color-difference formulae. *Color Research & Application*, 2(1):7–11, 1977. 12

[72] Gaurav Sharma, Wencheng Wu, and Edul N. Dalal. The CIEDE2000 color-difference formula: Implementation notes, supplementary test data, and mathematical observations. *Color Research & Application*, 30(1):21–30, 2005. 7, 36, 38, 42, 43

[73] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 12

[74] Shuangbing Song, Fan Zhong, Tianju Wang, Xueying Qin, and Changhe Tu. Guided linear upsampling. *ACM Transactions on Graphics (TOG)*, 42(4):1–12, 2023. 30, 31

[75] SaiKiran Tedla, Beixuan Yang, and Michael S Brown. Examining autoexposure for challenging scenes. In *ICCV*, 2023. 20

[76] Ethan Tseng, Yuxuan Zhang, Lars Jebe, Xuaner Zhang, Zhihao Xia, Yifei Fan, Felix Heide, and Jiawen Chen. Neural photo-finishing. *ACM Transactions on Graphics (TOG)*, 41 (6):238–1, 2022. 2

[77] Yael Vinker, Inbar Huberman-Spiegelglas, and Raanan Fattal. Unpaired learning for high dynamic range image tone mapping. In *ICCV*, 2021. 1

[78] Tengfei Wang, Jiaxin Xie, Wenxiu Sun, Qiong Yan, and Qifeng Chen. Dual-camera super-resolution with aligned attention modules. In *ICCV*, 2021. 1

[79] Yuzhi Wang, Haibin Huang, Qin Xu, Jiaming Liu, Yiqun Liu, and Jue Wang. Practical deep raw image denoising on mobile devices. In *ECCV*, 2020. 1

[80] Yufei Wang, Yi Yu, Wenhan Yang, Lanqing Guo, Lap-Pui Chau, Alex C Kot, and Bihan Wen. Raw image reconstruction with learned compact metadata. In *CVPR*, 2023. 6

[81] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: From error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004. 7, 12, 36, 38, 42, 43

[82] Yuxin Wu and Kaiming He. Group normalization. In *ECCV*, 2018. 7

[83] Yicheng Wu, Qiurui He, Tianfan Xue, Rahul Garg, Jiawen Chen, Ashok Veeraraghavan, and Jonathan T Barron. How to train neural networks for flare removal. In *ICCV*, 2021. 1

[84] Günther Wyszecki and Walter Stanley Stiles. *Color science: concepts and methods, quantitative data and formulae*. John wiley & sons, 2000. 21

[85] Yazhou Xing, Zian Qian, and Qifeng Chen. Invertible image signal processing. In *CVPR*, 2021. 1, 6, 7, 36, 37, 42, 43, 44

[86] David Young. Iterative methods for solving partial difference equations of elliptic type. *Transactions of the American Mathematical Society*, 76(1):92–111, 1954. 2

[87] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. CycleISP: Real image restoration via improved data synthesis. In *CVPR*, 2020. 1

[88] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Restormer: Efficient transformer for high-resolution image restoration. In *CVPR*, 2022. 29

[89] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 7, 36, 38, 42, 43

[90] Xuaner Zhang, Qifeng Chen, Ren Ng, and Vladlen Koltun. Zoom to learn, learn to zoom. In *CVPR*, 2019. 1

[91] Zhilu Zhang, Haolin Wang, Ming Liu, Ruohao Wang, Jiawei Zhang, and Wangmeng Zuo. Learning raw-to-sRGB mappings with inaccurately aligned supervision. In *ICCV*, 2021. 1, 6, 7, 15, 17, 19, 36, 40, 43, 44

[92] Luxi Zhao, Mahmoud Afifi, and Michael S Brown. Learning camera-agnostic white-balance preferences. In *ICCVW*, 2025. 5, 15, 16, 17, 21