Uni-Renderer: Unifying Rendering and Inverse Rendering Via Dual Stream Diffusion

Supplementary Material

1. Appendix / supplemental material

In this supplementary, we will first discuss the detailed network architecture and the detailed algorithm for calculating
different timesteps for reducing the tasks spaces. Then we
will provide a description of the configurations used for baseline comparison. We also include more qualitative cases
to demonstrate the capacity of our framework to perform
smooth rendering and inverse rendering.

009 1.1. Parallel Stream Diffusion

010 The design of our framework involves two parallel stable dif-011 fusions. The upper branch takes in channel-concatenated attributes. Its UNet input "Conv_in" and output "Conv_out" 012 013 layers are modified and extended to 24 channels for the corresponding input and output latents. The lower branch 014 remains unchanged. The communication between the up-015 016 per and lower branch are implemented through a crossconnected manner. We first take the intermediate feature 017 "mid_block_res_samples" from the upper encoder and add 018 019 it to the lower decoder through a zero convolution layer. We do the same for the lower encoder. Such design effectively 020 enables the communication between two stable diffusions 021 022 in a cross-conditioned manner. The introduction of the zero convolution layer maintains the pertaining weight not get 023 024 disrupted during training.

025 1.2. Model Training

During training, we adopted the x0 prediction into our loss calculation. It effectively helps to solve the channel bandwidth overhead problem. The training of diffusion models is performed on eight A800 GPUs, with a batch size of 4, a learning rate of 1e-5, and a total training iteration number of 150,000. We utilize the Adam optimizer with $adam_beta1$ and $adam_beta2$ equal to 0.8 and 0.999, respectively.

1.3. Modeling conditional distributions with TwoTimesteps

035To achieve rendering and inverse rendering within a single036model, we introduced a reduced timestep strategy to elim-037inate redundant tasks and thereby speed up convergence038time and generation quality. In algorithm 1, we showed the039pseudo-code for generating different timesteps.

1.4. Configuration for rendering baseline compari son

We show comparison against GAN-based material editing[3], Null-text inversion with prompt-to-prompt [2], Instruct-

Algorithm 1 Compute time steps matrix

- **Require:** *len_t* Length of timesteps matrix, defaults to 2 in our case.
- **Require:** *num_timesteps* Number of timesteps, ranging from 0 to *T*.

Require: *bs* Batch size

1: $timesteps \leftarrow$ initialize a zero matrix of size $len_t \times bs$

- 2: $idx \leftarrow random integer from 0 to len_t 1$
- 3: $all_t[idx] \leftarrow$ random integers from 0 to $num_timesteps-1$ for each column
- 4: for $i \leftarrow 0$ to $len_t 1$ do
- 5: **if** $i \neq idx$ **then**
- 6: **for** $j \leftarrow 0$ to bs 1 **do**
- 7: $all_t[i][j]$ \leftarrow random choice of $\{0, num_timesteps - 1\}$
- 8: end for
- 9: **end if**
- 10: end for
- 11: return timesteps

Pix2Pix [1], and InstructPix2Pix prompt-only version trained 044 on our data. For rendering baseline comparison, we first per-045 formed inverse rendering to acquire the necessary intrinsic at-046 tributes and used those to re-render with swapping attributes. 047 By doing this, we ensured our setting was the same as other 048 material editing pipelines. Next, we will go over each of the 049 baselines, and briefly discuss the testing configuration for 050 each of the methods. 051

InstructPix2Pix with our data [1]. This method takes 052 an input image and a text prompt for material editing. We 053 compared our method with two versions of InstructPix2Pix: 054 the finetuned version and the original version. For training, 055 we finetuned the model using a training set of 300 objects 056 and trained for 1000 epochs To create pairs for the image 057 editing framework, we varied two attributes-roughness and 058 metallicity-from 0 to 1 while keeping the other attribute at 059 0. Prompts are built as "make it more/less rough/metallic." 060 For example, for the prompt "make it rougher," the roughness 061 of the input and ground truth would be 0 and 1, respectively. 062 For testing, we evaluated both the finetuned model and the 063 original model on our validation set using the default settings: 064 inference steps set to 100, text CFG set to 7.5, and image 065 CFG set to 1.0. 066

Null text inversion [2]. This method, also a prompt-only067version, first optimizes the null text embedding to recover068

CVPR 2025 Submission #4621. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

078

079

080

081

082



Figure 1. More Qualitative on Inverse Rendering We included more qualitative cases to demonstrate the ability of Uni-renderer to perform inverse. Best viewed in color.



Figure 2. More Qualitative on Inverse Rendering We included more qualitative cases to demonstrate the ability of Uni-renderer to perform inverse. Best viewed in color.

the original DDIM latent sequences at inference with a high
CFG value. It then performs prompt-to-prompt for image
editing. For optimization, the steps are set to 300, and the
prompt pairs used are in the format "a object name" to "a
metallic/rough object name." For inference, the steps are set
to 50, and the CFG value is set to 7.5.

O75 Subias et. al [3] The method takes an image as input
o76 along with a scalar as input to perform relative material
o77 editing for glossiness and metallic. The method requires

an input mask for localizing the object in the image. We generated the mask using the provided format input scripts. We input the image with the required transformation and tested it with a scalar of 1.0.

1.5. More Visual Samples

In this section, we will include more qualitative samples to083further support the effectiveness of our method. We also084include some cases to demonstrate the capacity of our frame-085



Figure 3. More Qualitative on Accurate attributes editing We included more qualitative cases to demonstrate the ability of Uni-renderer to perform rendering. The leftmost are reference images, and we have provided both increasing metallic and roughness. Best viewed in color.

086 work to perform smooth rendering.

087 References

- [1] Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow image editing instructions, 2023. 1
- (2) Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and
 (3) Daniel Cohen-Or. Null-text inversion for editing real images
 (4) using guided diffusion models, 2022. 1
- [3] J. Daniel Subias and Manuel Lagunas. In-the-wild material
 appearance editing using perceptual attributes, 2023. 1, 2