

---

# Night-to-Day Translation via Illumination Degradation Disentanglement

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Night-to-Day translation (Night2Day) aims to achieve day-like vision for nighttime  
2 scenes. However, processing night images with complex degradations remains a  
3 significant challenge under unpaired conditions. Previous methods that uniformly  
4 mitigate these degradations have proven inadequate in simultaneously restoring  
5 daytime domain information and preserving underlying semantics. In this paper,  
6 we propose **N2D3** (Night-to-Day via **D**egradation **D**isentangle**m**ent) to identify  
7 different degradation patterns in nighttime images. Specifically, our method com-  
8 prises a degradation disentanglement module and a degradation-aware contrastive  
9 learning module. Firstly, we extract physical priors from a photometric model  
10 based on Kubelka-Munk theory. Then, guided by these physical priors, we design a  
11 disentanglement module to discriminate among different illumination degradation  
12 regions. Finally, we introduce the degradation-aware contrastive learning strategy  
13 to preserve semantic consistency across distinct degradation regions. Our method  
14 is evaluated on two public datasets, **demonstrating a significant improvement of**  
15 **5.4 FID on BDD100K and 10.3 FID on Alderley.**

## 16 1 Introduction

17 Nighttime images often suffer from severe information loss, posing significant challenges to both  
18 human visual recognition and computer vision tasks including detection, segmentation, *etc.* [14].  
19 In contrast, daylight images exhibit rich content and intricate details. Achieving day-like nighttime  
20 vision remains a primary objective in nighttime perception, sparking numerous pioneering works [30].  
21 Night-to-Day image translation (Night2Day) offers a comprehensive solution to achieve day-like  
22 vision at night. The primary goal is to transform images from nighttime to daytime while maintaining  
23 their underlying semantic structure. However, achieving this goal is challenging. It requires to process  
24 complex degraded images using unpaired data, which raises additional difficulties compared to other  
25 image translation tasks.

26 Recently, explorations have been made in Night2Day. Early approaches, such as ToDayGAN,  
27 demonstrated the effectiveness of cycle-consistent learning in maintaining semantic structure [1].  
28 Subsequent methods incorporated auxiliary structure regularization techniques, including perceptual  
29 loss and uncertainty regularization, to better preserve the original structure [33, 18]. Furthermore,  
30 some methods utilized daytime images with nearby GPS locations to aid in coarse structure regular-  
31 ization [26]. However, these methods often neglect the complex degradations at nighttime, applying  
32 structure regularization uniformly and resulting in severe artifacts. To address this issue, more recent  
33 approaches adopt auxiliary human annotations to maintain semantic consistency, such as segmenta-  
34 tion maps and bounding boxes [16, 22]. Despite their potential, these methods are labor-intensive  
35 and challenging, especially since many nighttime scenes are beyond human cognition.

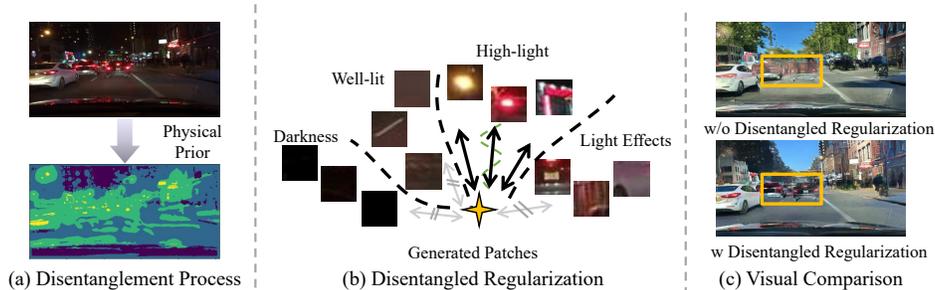


Figure 1: Illustration of our motivation. (a) The disentanglement process leverages physical priors. (b) The image patches are restored individually for each degradation type. (c) The proposed Disentangled Regularization improves the overall performance.

36 The critical limitation of the aforementioned methods is the disregard for complex degraded regions.  
 37 Specifically, different regions in nighttime images possess varying characteristics, such as extreme  
 38 darkness, well-lit regions, light effects, *etc.* Treating all these degraded regions equally could adversely  
 39 impact the results. As illustrated in Figure 1, our key insight emphasizes that nighttime images suffer  
 40 from various degradations, necessitating customizing restoration for different degradation types.  
 41 Intuitively, we manage to disentangle nighttime images into patches according to the recognized  
 42 degradation type and learn individual restoration patterns for them to enhance the overall performance.

43 Motivated by this point, we propose N2D3 (Night to Day via Degradation Disentanglement), which  
 44 utilizes Generative Adversarial Networks (GANs) to bridge the domain gap between nighttime and  
 45 daytime in a degradation-aware manner, as illustrated in Figure 2. There are two modules in N2D3,  
 46 including physical-informed degradation disentanglement and degradation-aware contrastive learning,  
 47 which are employed to preserve the semantic structure of nighttime images. In the disentanglement  
 48 of nighttime degradation, a photometric model tailored to nighttime scenes is conducted to extract  
 49 physical priors. Subsequently, the illuminance and physical priors are integrated to disentangle  
 50 regions into darkness, well-lit, high-light, and light effects. Building on this, degradation-aware  
 51 contrastive learning is designed to constrain the similarity of the source and generated images in  
 52 different regions. It comprises disentanglement-guided sampling and reweighting strategies. The  
 53 sampling strategy mines valuable anchors and hard negative examples, while the reweighting process  
 54 assigns their weights. They enhance vanilla contrastive learning by prioritizing valuable patches with  
 55 appropriate attention. Ultimately, our method yields highly faithful results that are visually pleasing  
 56 and beneficial for downstream vision tasks including keypoint matching and semantic segmentation.

57 Our contributions are summarized as follows:

- 58 (1) We propose the N2D3 translation method based on the illumination degradation disentanglement  
 59 module, which enables degradation-aware restoration of nighttime images.
- 60 (2) We present a novel degradation-aware contrastive learning module to preserve the semantic  
 61 structure of generated results. The core design incorporates disentanglement-guided sampling and  
 62 reweighting strategies, which greatly enhance the performance of vanilla contrastive learning.
- 63 (3) Experimental results on two public datasets underscore the significance of considering distinct  
 64 degradation types in nighttime scenes. Our method achieves state-of-the-art performance in visual  
 65 effects and downstream tasks.

## 66 2 Related Work

67 **Unpaired Image-to-Image Translation.** Unpaired image-to-image translation addresses the chal-  
 68 lenge of lacking paired data, providing an effective self-supervised learning strategy. To overcome the  
 69 efficiency limitations of traditional cycle-consistency learning, Park *et al.*, first introduces contrastive  
 70 learning to this domain, achieving efficient one-sided learning[20]. Following this work, several stud-  
 71 ies have improved the contrastive learning by generating hard negative examples [24], re-weighting  
 72 positive-negative pairs [31], and selecting key samples [9]. Furthermore, other constraints, such as  
 73 density [27] and path length [28], have been explored in unpaired image translation. However, all  
 74 these works neglect physical priors in the nighttime, leading to suboptimal results in Night2Day.

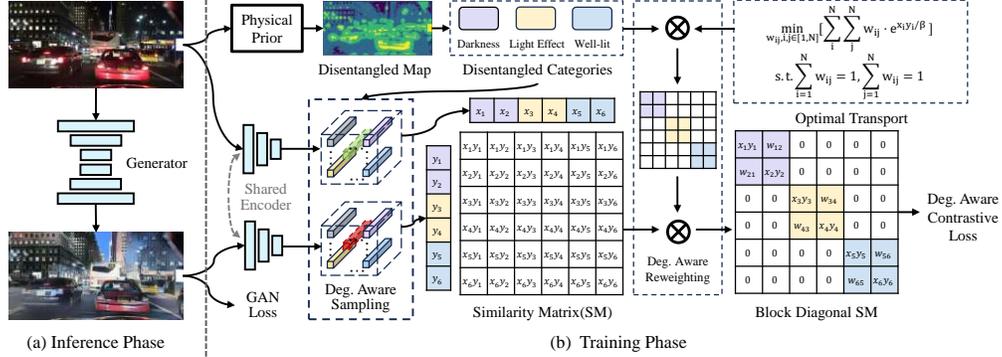


Figure 2: The overall architecture of the proposed N2D3 method. The training phase contains the physical prior informed degradation disentanglement module and degradation-aware contrastive learning module. They are utilized to optimize the ResNet-based generator which is the main part in the inference phase.

75 **Nighttime Domain Translation.** Domain translation techniques have been applied to address adverse  
 76 nighttime conditions. An early contribution is made by Anoosheh *et al.*, which demonstrates the  
 77 effectiveness of cycle-consistent learning in Night2Day[1]. Following this, many works incorporate  
 78 different modules into cycle-consistent learning to enhance structural modeling capabilities. Zheng *et*  
 79 *al.* incorporate a fork-shaped encoder to enhance visual perceptual quality[33]. AUGAN employs  
 80 uncertainty estimation to mine useful features in nighttime images[18]. Fan *et al.* explore inter-  
 81 frequency relation knowledge to streamline the Night2Day process[5]. Xia *et al.* utilize nearby GPS  
 82 locations to form paired night and daytime images, providing weak supervision[26]. Some other  
 83 studies incorporate human annotations to impose structural constraints, overlooking the practical  
 84 difficulty of acquiring such annotations at nighttime with multiple degradations [11][16] [22]. To  
 85 address the concerns of the aforementioned methods, the proposed N2D3 explores patch-wise  
 86 contrastive learning with physical guidance, so as to achieve degradation-aware Night2Day. N2D3 is  
 87 free of human annotations and offers comprehensive structural modeling to provide faithful translation  
 88 results.

### 89 3 Methods

90 Given nighttime image  $\mathbf{I}_N \in \mathcal{N}$  and daytime image  $\mathbf{I}_D \in \mathcal{D}$ , the goal of Night2Day is to translate  
 91 images from nighttime to daytime while preserving content semantic consistency. This involves the  
 92 construction of a mapping function  $\mathcal{F}$  with parameters  $\theta$ , which can be formulated as  $\mathcal{F}_\theta : \mathbf{I}_N \rightarrow \mathbf{I}_D$ .  
 93 Our method N2D3 is illustrated in Figure 2. To train a generator for Night2Day, we employ GANs as  
 94 the overall learning framework to bridge the domain gap between nighttime and daytime. Our core  
 95 design, consisting of the degradation disentanglement module and the degradation-aware contrastive  
 96 learning module, aims to preserve the structure from the source images and suppress artifacts.

97 In this section, we first introduce physical priors in the nighttime environment, and then describe  
 98 the degradation disentanglement module and the degradation-aware contrastive learning module,  
 99 respectively.

#### 100 3.1 Physical Priors for Nighttime Environment

101 The illumination degradations at night are primarily categorized as darkness, well-lit regions, high-  
 102 light regions, and light effects. As shown in Figure 3, well-lit represents the diffused reflectance under  
 103 normal light, while the light effects denote phenomena such as flare, glow, and specular reflections.  
 104 Intuitively, these regions can be disentangled through the analysis of illumination distribution. Among  
 105 these degradation types, darkness and high-light are directly correlated with illuminance and can be  
 106 effectively disentangled through illumination estimation.

107 As a common practice, we estimate the illuminance map  $L$  by utilizing the maximum RGB channel  
 108 of image  $\mathbf{I}_N$  as  $L = \max_{c \in R, G, B} \mathbf{I}_N^c$ . Then k-nearest neighbors [4] is employed to acquire three  
 109 clusters representing darkness, well-lit, and high-light regions. These clusters are aggregated as  
 110 masks  $M_d, M_n, M_h$ . However, the challenge arises with light effects that are mainly related to

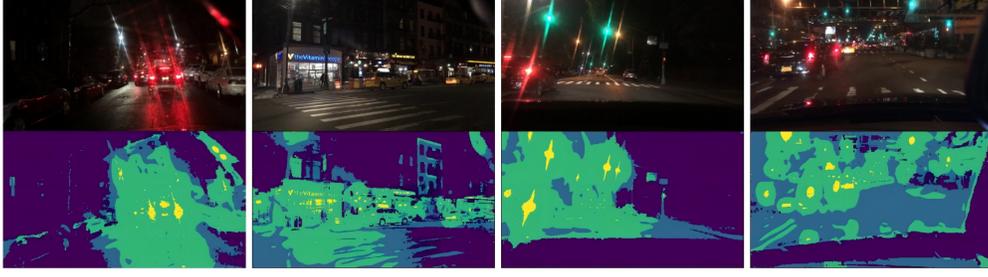


Figure 3: The first row displays nighttime images, while the second row shows the corresponding degradation disentanglement results. The color progression from **blue**, **light blue**, **green** to **yellow** corresponds to the following regions: darkness, well-lit, light effects, and high-light, respectively.

111 the illumination. Light effects regions tend to intertwine with well-lit regions when using only the  
 112 illumination map, as they often share similar illumination densities. To disentangle light effects from  
 113 well-lit regions, we need to introduce additional physical priors.

114 To extract the physical priors for disentangling light effects, we develop a photometric model derived  
 115 from Kubelka-Munk theory [17]. This model characterizes the spectrum of light  $E$  reflected from an  
 116 object as follows:

$$E(\lambda, x) = e(\lambda, x)(1 - \rho_f(x))^2 R_\infty(\lambda, x) + e(\lambda, x)\rho_f(x), \quad (1)$$

117 here  $x$  represents the horizontal component for analysis, while the analysis of the vertical component  
 118  $y$  is the same as the horizontal component.  $\lambda$  corresponds to the wavelength of light.  $e(\lambda, x)$  signifies  
 119 the spectrum, representing the illumination density and color.  $\rho_f$  stands for the Fresnel reflectance  
 120 coefficient.  $R_\infty$  is the material reflectivity function, formulated as follows at a specific location  
 121  $x = x_0$ :

$$R(\lambda) = a(\lambda) - \sqrt{a(\lambda)^2 - 1}, a(\lambda) = 1 + \frac{k(\lambda)}{s(\lambda)}, \quad (2)$$

122 where  $k(\lambda)$  and  $s(\lambda)$  denote the absorption and scattering coefficients, respectively. This formulation  
 123 implies that for any local pixels, the material reflectivity is determined if the material is given.  
 124 Assuming  $C$  is the material distribution function, which describes the material type varying across  
 125 locations, the material reflectivity  $R_\infty$  can be formulated as:

$$R_\infty(\lambda, x) = R(\lambda)C(x). \quad (3)$$

126 Since the mixture of light effects and well-lit regions has been obtained previously, the core of  
 127 disentangling light effects from well-lit regions lies in separating the illumination  $e(\lambda, x)$  and  
 128 reflectance components  $R(\lambda)C(x)$ . Note that the Fresnel reflectance coefficient  $\rho_f(x)$  approaches 0 in  
 129 reflectance-dominating well-lit regions, while  $\rho_f(x)$  approaches 1 in illumination-dominating light  
 130 effects regions. According to Equation (1), the photometric model for the mixture of light effects and  
 131 well-lit regions is formulated as:

$$E(\lambda, x) = \begin{cases} e(\lambda, x), & \text{if } x \notin \Omega \\ e(\lambda, x)R(\lambda)C(x), & \text{if } x \in \Omega \end{cases}, \quad (4)$$

132 where  $\Omega$  denotes the reflectance-dominating well-lit regions.

133 Subsequently, we observe that the following color invariant response to the regions with high color  
 134 saturation, which is suitable to extract the illumination:

$$N_{\lambda^m x^n} = \frac{\partial^{m+n-1}}{\partial \lambda^{m-1} \partial x^n} \left\{ \frac{1}{E(\lambda, x)} \frac{\partial E(\lambda, x)}{\partial \lambda} \right\}, \quad (5)$$

135 This invariant has the following characteristics:

$$\begin{aligned} N_{\lambda^m x^n} &= \frac{\partial^{m+n-2}}{\partial \lambda^{m-1} \partial x^{n-1}} \frac{\partial}{\partial x} \left\{ \frac{1}{E(\lambda, x)} \frac{\partial E(\lambda, x)}{\partial \lambda} \right\} \\ &= \frac{\partial^{m+n-2}}{\partial \lambda^{m-1} \partial x^{n-1}} \frac{\partial}{\partial x} \left\{ \frac{1}{e(\lambda, x)} \frac{\partial e(\lambda, x)}{\partial \lambda} + \frac{1}{R(\lambda)C(x)} \frac{\partial R(\lambda)C(x)}{\partial \lambda} \right\} \\ &= \frac{\partial^{m+n-1}}{\partial \lambda^{m-1} \partial x^n} \left\{ \frac{1}{e(\lambda, x)} \frac{\partial e(\lambda, x)}{\partial \lambda} \right\}. \end{aligned} \quad (6)$$

136 Equation (5) to Equation (6) demonstrate that the invariant  $N_{\lambda^m x^n}$  captures the features only related  
 137 to illumination  $e(\lambda, x)$ . Consequently, we assert that  $N_{\lambda^m x^n}$  functions as a light effects detector  
 138 because light effects are mainly related to the illumination. It allows us to design the illumination  
 139 disentanglement module based on this physical prior.

### 140 3.2 Degradation Disentanglement Module

141 In this subsection, we will elucidate how to incorporate the invariant for extracting light effects into  
 142 the disentanglement in computation. As common practice, the following second and third-order  
 143 components, both horizontally and vertically, are taken into account in the practical calculation of the  
 144 final invariant, which is denoted as  $N$ :

$$N = \sqrt{N_{\lambda x}^2 + N_{\lambda \lambda x}^2 + N_{\lambda y}^2 + N_{\lambda \lambda y}^2}. \quad (7)$$

145 here  $N_{\lambda x}$  and  $N_{\lambda \lambda x}$  can be computed through  $E(\lambda, x)$  by simplifying Equation (5). The calculation  
 146 of  $N_{\lambda y}$  and  $N_{\lambda \lambda y}$  are the same. Specifically,

$$N_{\lambda x} = \frac{E_{\lambda x} E - E_{\lambda} E_x}{E^2}, N_{\lambda \lambda x} = \frac{E_{\lambda \lambda x} E^2 - E_{\lambda \lambda} E_x E - 2E_{\lambda x} E_{\lambda} E + 2E_{\lambda}^2 E_x}{E^3}, \quad (8)$$

147 where  $E_x$  and  $E_{\lambda}$  denote the partial derivatives of  $x$  and  $\lambda$ .

148 To compute each component in the invariant  $N$ , we develop a computation scheme starting with the  
 149 estimation of  $E$  and its partial derivatives  $E_{\lambda}$  and  $E_{\lambda \lambda}$  using the Gaussian color model:

$$\begin{bmatrix} E(x, y) \\ E_{\lambda}(x, y) \\ E_{\lambda \lambda}(x, y) \end{bmatrix} = \begin{bmatrix} 0.06, & 0.63, & 0.27 \\ 0.3, & 0.04, & -0.35 \\ 0.34, & -0.6, & 0.17 \end{bmatrix} \begin{bmatrix} R(x, y) \\ G(x, y) \\ B(x, y) \end{bmatrix}, \quad (9)$$

150 where  $x, y$  are pixel locations of the image. Then, the spatial derivatives  $E_x$  and  $E_y$  are calculated by  
 151 convolving  $E$  with Gaussian derivative kernel  $g$  and standard deviation  $\sigma$ :

$$E_x(x, y, \sigma) = \sum_{t \in \mathbf{Z}} E(t, y) \frac{\partial g(x - t, \sigma)}{\partial x}, \quad (10)$$

152 where  $t$  denotes the index of the horizontal component  $x$  and  $\mathbf{Z}$  represents set of integers. The spatial  
 153 derivatives for  $E_{\lambda x}$  and  $E_{\lambda \lambda x}$  are obtained by applying Equation (10) to  $E_{\lambda}$  and  $E_{\lambda \lambda}$ . Then invariant  
 154  $N$  can be obtained following Equation (8) and Equation (7).

155 To extract the light effects, ReLU and normalization functions are first applied to filter out minor  
 156 disturbances. Then, by filtering invariant  $N$  with the well-lit mask  $M_n$ , we obtain the light effects  
 157 from the well-lit regions. The operations above can be formulated as:

$$M_{le} = \text{ReLU}\left(\frac{N - \mu(N)}{\sigma(N)}\right) \odot M_n, \quad (11)$$

158 while the well-lit mask are refined:  $M_n \leftarrow M_n - M_{le}$ .

159 With the initial disentanglement in Section 3.1, we obtain the final disentanglement:  $M_d, M_n, M_h$   
 160 and  $M_{le}$ . All the masks are stacked to obtain the disentanglement map. Through the employment of  
 161 the aforementioned techniques and processes, we successfully achieve the disentanglement of various  
 162 degradation regions.

### 163 3.3 Degradation-Aware Contrastive Learning

164 For unpaired image translation, contrastive learning has validated its effectiveness for the preservation  
 165 of content. It targets to maximize the mutual information between patches in the same spatial location  
 166 from the generated image and the source image as below:

$$\ell(v, v^+, v^-) = -\log \frac{\exp(v \cdot v^+ / \tau)}{\exp(v \cdot v^+ / \tau) + \sum_{n=1}^Q \exp(v \cdot v_n^- / \tau)}, \quad (12)$$

167  $v$  is the anchor that denotes the patch from the generated image. The positive example  $v^+$  corresponds  
 168 to the source image patch with the same location as the anchor  $v$ . The negative examples  $v^-$  represent

169 patches with locations distinct from that of the anchor  $v$ .  $Q$  denotes the total number of negative  
 170 examples. In our work, the key insight of degradation-aware contrastive learning lies in two folds: (1)  
 171 How to sample the anchor, positive, and negative examples. (2) How to manage the focus on different  
 172 negative examples.

173 **Degradation-Aware Sampling.** In this paper, N2D3 selects the anchor, positive, and negative patches  
 174 under the guidance of the disentanglement results. Initially, based on the disentanglement mask  
 175 obtained in the Section 3.2, we compute the patch count for different degradation types, denoting as  
 176  $K_s$ ,  $s \in [1, 4]$ . Then, within each degradation region, the anchors  $v$  are randomly selected from the  
 177 patches of generated daytime images  $I_{\mathcal{N} \rightarrow \mathcal{D}}$ . The positive examples  $v^+$  are sampled from the same  
 178 locations with the anchors in the source nighttime images  $I_{\mathcal{N}}$ , and the negative examples  $v^-$  are  
 179 randomly selected from other locations of  $I_{\mathcal{N}}$ . For each anchor, there is one corresponding positive  
 180 example and  $K_s$  negative examples. Subsequently, the sample set with the same degradation type  
 181 will be assigned weights and the contrastive loss will be computed in the following steps.

182 **Degradation-Aware Reweighting.** Despite the careful selection of anchor, positive, and negative  
 183 examples, the importance of anchor-negative pairs still differs within the same degradation. A known  
 184 principle of designing contrastive learning is that the hard anchor-negative pairs (*i.e.*, the pairs with  
 185 high similarity) should assign higher attention. Thus, weighted contrastive learning can be formulated  
 186 as:

$$\ell(v, v^+, v^-, w_n) = -\log \frac{\exp(v \cdot v^+ / \tau)}{\exp(v \cdot v^+ / \tau) + \sum_{n=1}^Q w_n \exp(v \cdot v_n^- / \tau)}, \quad (13)$$

187  $w_n$  denotes the weight of the  $n$ -th anchor-negative pairs.

188 The contrastive objective is depicted in the *Similarity Matrix* in Figure 2. The patches in different  
 189 regions are obviously easy examples. We suppress their weights to 0, which transforms the similarity  
 190 matrix into a blocked diagonal matrix with  $\text{diag}(A_1, \dots, A_4)$ . Within each degradation matrix  
 191  $A_s$ ,  $s \in [1, 4]$ , a soft reweighting strategy is implemented. Specifically, for each anchor-negative  
 192 pair, we apply optimal transport to yield an optimal transport plan, serving as a reweighting matrix  
 193 associated with the disentangled results. It can adaptively optimize and avoid manual design. The  
 194 reweight matrix for each degradation type is formulated as:

$$\begin{aligned} \min_{w_{ij}, i, j \in [1, K_s]} & \left[ \sum_{i=1}^{K_s} \sum_{j=1, i \neq j}^{K_s} w_{ij} \cdot \exp(v_i \cdot v_j^- / \tau) \right], \\ & \sum_{i=1}^{K_s} w_{ij} = 1, \sum_{j=1}^{K_s} w_{ij} = 1, i, j \in [1, K_s], \end{aligned} \quad (14)$$

195 The aforementioned operations transform the contrastive objective to the *Block Diagonal Similarity*  
 196 *Matrix* depicted in Figure 2. As a common practice, our degradation-aware contrastive loss is applied  
 197 to the  $S$  layers of the CNN feature extractor, formulated as:

$$\mathcal{L}_{DegNCE}(\mathcal{F}) = \sum_{l=1}^S \ell(v, v^+, v^-, w_n). \quad (15)$$

### 198 3.4 Other Regularizations

199 As a common practice, GANs are employed to bridge the domain gap between daytime and nighttime.  
 200 The adversarial loss is formulated as:

$$\begin{aligned} \mathcal{L}_{adv}(\mathcal{F}) &= \|D(\mathbf{I}_{\mathcal{N} \rightarrow \mathcal{D}}) - 1\|_2^2, \\ \mathcal{L}_{adv}(D) &= \|D(\mathbf{I}_{\mathcal{D}}) - 1\|_2^2 + \|D(\mathbf{I}_{\mathcal{N} \rightarrow \mathcal{D}})\|_2^2, \end{aligned} \quad (16)$$

201 where  $D$  denotes the discriminator network. The final loss function is formatted as :

$$\begin{aligned} \mathcal{L}(\mathcal{F}) &= \mathcal{L}_{adv}(\mathcal{F}) + \mathcal{L}_{DegNCE}(\mathcal{F}), \\ \mathcal{L}(D) &= \mathcal{L}_{adv}(D). \end{aligned} \quad (17)$$

## 202 4 Experiments

### 203 4.1 Experimental Settings

204 **Datasets.** Experiments are conducted on the two public datasets BDD100K [29] and Alderley [19].  
 205 **Alderley** dataset consists of images captured along the same route twice: once on a sunny day and  
 206 another time during a stormy rainy night. The nighttime images in this dataset are often blurry due to  
 207 the rainy conditions, which makes Night2Day challenging. **BDD100K** dataset is a large-scale high-  
 208 resolution autonomous driving dataset. It comprises 100,000 video clips under various conditions.  
 209 For each video, a keyframe is selected and meticulously annotated with details. We reorganized this  
 210 dataset based on its annotations, resulting in 27,971 night images for training and 3,929 night images  
 211 for evaluation.

212 **Evaluation Metric.** Following common practice, we utilize the *Fréchet Inception Distance* (FID)  
 213 scores [7] to assess whether the generated images align with the target distribution. This assessment  
 214 helps determine if a model effectively transforms images from the night domain to the day domain.  
 215 Additionally, we seek to understand the extent to which the generated daytime images maintain  
 216 structural consistency compared to the original inputs. To measure this, we employ SIFT scores,  
 217 mIoU scores and LPIPS distance [32].

218 **DownStream Vision Task.** Two downstream tasks are conducted. In the Alderley dataset, GPS  
 219 annotations indicate the locations of two images, one in the nighttime and the other in the daytime,  
 220 as the same. We calculate the number of SIFT-detected key points between the generated daytime  
 221 images and their corresponding daytime images to measure if the two images represent the same  
 222 location. The BDD100K dataset includes 329 night images with semantic annotations. We employ  
 223 Deeplabv3 pretrained on the Cityscapes dataset as the semantic segmentation model [2], then perform  
 224 inference on our generated daytime images without any additional training and compute the mIoU  
 225 (mean Intersection over Union).

Table 1: The quantitative results on Alderley and BDD100k. ↓ means lower result is better. ↑ means higher is better.

| Dataset           |             | Alderley    |              |              | BDD100k     |              |              |
|-------------------|-------------|-------------|--------------|--------------|-------------|--------------|--------------|
| Methods           |             | FID↓        | LPIPS↓       | SIFT↑        | FID↓        | LPIPS↓       | mIoU↑        |
| Original          | Conf./Jour. | 210         | -            | 3.12         | 101         | -            | 15.63        |
| CycleGAN[34]      | ICCV 2017   | 167         | 0.706        | 3.36         | 51.7        | 0.477        | 13.42        |
| StarGAN[3]        | CVPR 2018   | 117         | -            | 3.28         | 68.3        | -            | -            |
| ToDayGAN[1]       | ICRA 2019   | 104         | 0.770        | 4.14         | 43.8        | 0.577        | 16.77        |
| UGATIT[15]        | ICLR 2020   | 170         | -            | 2.51         | 72.2        | -            | -            |
| CUT[20]           | ECCV 2020   | 64.7        | 0.707        | 6.78         | 55.5        | 0.583        | 9.30         |
| ForkGAN[33]       | ECCV 2020   | 61.2        | 0.759        | 12.1         | 37.6        | 0.581        | 11.81        |
| AUGAN[18]         | BMVC 2021   | 65.2        | -            | -            | 38.6        | -            | -            |
| MoNCE[31]         | CVPR 2022   | 72.7        | 0.737        | 6.35         | 40.2        | 0.502        | 17.21        |
| Decent[27]        | NIPS 2022   | 76.5        | 0.768        | 6.31         | 40.3        | 0.582        | 10.49        |
| Santa[28]         | CVPR 2023   | 67.1        | 0.757        | 6.93         | 36.9        | 0.559        | 11.03        |
| N2D-LPNet[5]      | CVPR 2023   | -           | -            | -            | 69.1        | -            | -            |
| EnlightenGAN [13] | TIP 2021    | 209.8       | -            | 2.00         | 103.5       | -            | 16.10        |
| Zero-DCE [6]      | TPAMI 2022  | 246.4       | -            | 4.34         | 90.5        | -            | 15.90        |
| DeLight [21]      | ECCV 2022   | 222.9       | -            | 3.07         | 113.8       | -            | 14.48        |
| LLformer [23]     | AAAI 2023   | 275.6       | -            | 7.62         | 123.1       | -            | 15.28        |
| WCDM [12]         | ToG 2023    | 239.6       | -            | 7.10         | 124.3       | -            | 16.32        |
| GSAD [8]          | NIPS 2023   | 214.7       | -            | 6.29         | 116.0       | -            | 15.76        |
| N2D3(Ours)        | -           | <b>50.9</b> | <b>0.650</b> | <b>16.62</b> | <b>31.5</b> | <b>0.466</b> | <b>21.58</b> |

### 226 4.2 Results on Alderley

227 We first apply Night2Day on the Alderley dataset, a challenging collection of nighttime images  
 228 captured on rainy nights. In Figure 4, we present a visual comparison of the results. CycleGAN [34]  
 229 and CUT [20] manage to preserve the general structural information of the entire image but often  
 230 lose many fine details. ToDayGAN [1], ForkGAN [33], Decent [27], and Santa [28] tend to miss  
 231 important elements such as cars in their results.

232 In Table 1, thirteen translation methods and three enhancement methods are compared, considering  
 233 both visual effects and keypoint matching metrics. Our method showcases **an improvement of 10.3**

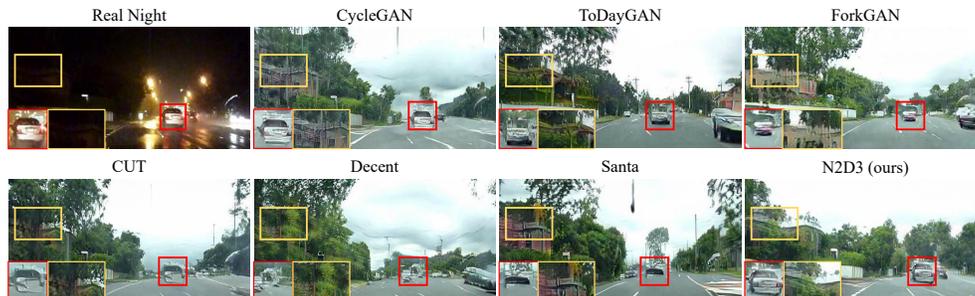


Figure 4: The qualitative comparison results on the Alderley dataset.

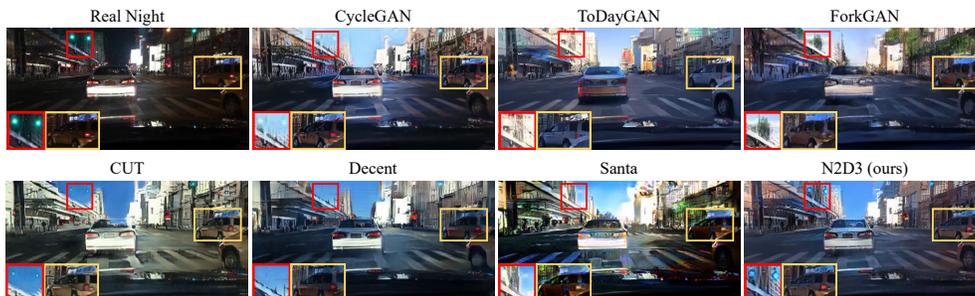


Figure 5: The qualitative comparison results on the BDD100K dataset.

234 **in FID scores and 4.52 in SIFT scores** compared to the previous state-of-the-art. This suggests that  
 235 N2D3 successfully achieves photorealistic daytime image generation, underscoring its potential for  
 236 robotic localization applications. The qualitative comparison results are demonstrated in Figure 4. In  
 237 conclusion, N2D3 achieves top scores in both FID and LPIPS metrics, demonstrating its superiority  
 238 in the Night2Day task. N2D3 excels in generating photorealistic daytime images while effectively  
 239 preserving structures, even in challenging scenarios such as rainy nights in the Alderley.

### 240 4.3 Results on BDD100K

241 We conducted experiments on a larger-scale dataset, BDD100K, focusing on more general night  
 242 scenes. The qualitative results can be found in Figure 5. CycleGAN, ToDayGAN, and CUT succeed  
 243 in preserving the structure in well-lit regions. ForkGAN, Santa, and Decent demonstrate poor  
 244 performance in such challenging scenes. Regrettably, none of them excel in handling light effects and  
 245 exhibit weak performance in maintaining global structures. With a customized design specifically  
 246 addressing light effects, our method successfully preserves the structure in all regions.

247 The quantitative results are presented in Table 1. As the scale of the dataset increases, all the  
 248 compared methods show an improvement in their performance. Notably, N2D3 demonstrates the best  
 249 performance with **a significant improvement of 5.4 in FID scores**, showcasing its ability to handle a  
 250 broader range of nighttime scenes and establishing itself as the most advanced method in this domain.

251 We also investigate the potential of Night2Day in enhancing downstream vision tasks in nighttime  
 252 environments using the BDD100K dataset. The quantitative results are summarized in Table 1.  
 253 The enhancement methods demonstrate a slight improvement in segmentation results, while some  
 254 image-to-image translation methods have a negative impact on performance. N2D3 exhibits the best  
 255 performance in enhancing nighttime semantic segmentation with **a remarkable improvement of**  
 256 **5.95 in mIoU** compared to inferring the segmentation model directly on nighttime images.

257 In conclusion, N2D3 achieves top scores in both FID and LPIPS metrics, establishing itself as the  
 258 most advanced method for the Night2Day task. It excels in generating photorealistic daytime images  
 259 while preserving local and global structures. Moreover, the substantial improvement in nighttime  
 260 semantic segmentation highlights its benefits for downstream tasks and its potential for wide-ranging  
 261 applications.

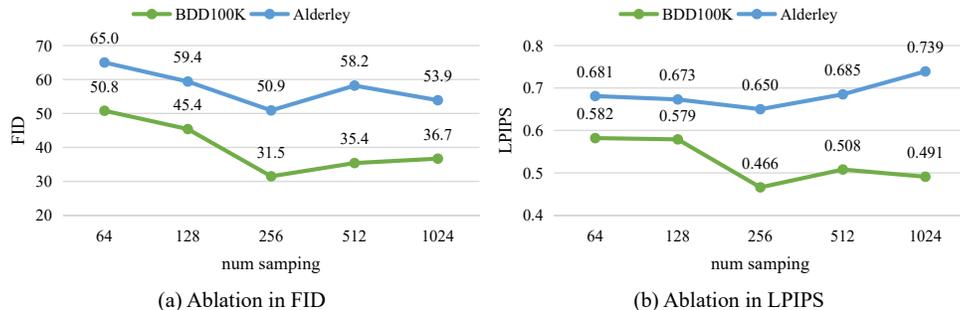


Figure 6: The quantitative results of ablation on the number of patches of the degradation-aware sampling.

Table 2: The quantitative results of ablation on the main component of degradation-aware contrastive learning. (a) denotes the degradation-aware sampling, and (b) denotes the degradation-aware reweighting.  $L$  and  $N$  denotes the invariant types.

| Main Component |     | BDD100K |       | Alderley |       |       | Invariant Type |   | BDD100K |       | Alderley |       |       |
|----------------|-----|---------|-------|----------|-------|-------|----------------|---|---------|-------|----------|-------|-------|
| (a)            | (b) | FID     | LPIPS | FID      | LPIPS | SIFT  | L              | N | FID     | LPIPS | FID      | LPIPS | SIFT  |
| ✗              | ✗   | 55.5    | 0.583 | 64.7     | 0.707 | 6.78  | ✗              | ✗ | 55.5    | 0.583 | 64.7     | 0.707 | 6.78  |
| ✓              | ✗   | 36.9    | 0.495 | 56.6     | 0.698 | 16.52 | ✓              | ✗ | 49.1    | 0.592 | 62.9     | 0.726 | 9.83  |
| ✓              | ✓   | 31.5    | 0.466 | 50.9     | 0.650 | 16.62 | ✓              | ✓ | 31.5    | 0.466 | 50.9     | 0.650 | 16.62 |

#### 262 4.4 Ablation Study

263 **Ablation on the main component of degradation-aware contrastive learning.** The core design of  
 264 the degradation-aware contrastive learning module relies on two main components: (a) degradation-  
 265 aware sampling, and (b) degradation-aware reweighting. As shown in Table 2, when degradation-  
 266 aware sampling is exclusively activated, there is a noticeable decrease in FID on both datasets  
 267 compared to the baseline (no components activated). Notably, the combination of degradation-aware  
 268 sampling and reweighting achieves the lowest FID on both BDD100K and Alderley, indicating the  
 269 effectiveness of degradation-aware sampling in conjunction with degradation-aware reweighting.

270 **Ablation on the number of patches in the degradation-aware sampling.** To explore the impact  
 271 of the number of sampling patches in our method, we conduct an ablation study on the number of  
 272 sampling patches with settings of 64, 128, 256, 512, and 1024 for degradation-aware sampling. The  
 273 FID and LPIPS scores are evaluated, as shown in Figure 6. The optimal performance is achieved with  
 274 256 patches, and increasing the number of sampling patches beyond this point leads to a degradation  
 275 in performance.

276 **Ablation on the type of the invariant in disentanglement.** To explore different invariants for  
 277 obtaining degradation-disentangled prototypes, we conduct an ablation study on the type of invariant.  
 278 As shown in Table 2, when  $L$  is enabled, the FID decreases from 55.5 to 49.1 on BDD100K and  
 279 from 64.7 to 62.9 on Alderley. This suggests that incorporating illuminance maps helps in reducing  
 280 the perceptual gap between generated and source nighttime images. When  $N$  is activated, there  
 281 is a consistent improvement in FID on both datasets, indicating that considering physical priors  
 282 invariant contributes to more realistic image generation. The combination of both illuminance map  
 283 and physical prior invariant results in the lowest FID on both datasets, showcasing the complementary  
 284 nature of these degradation types in improving contrastive learning.

## 285 5 Conclusion

286 This paper introduces a novel solution for the Night2Day image translation task, focusing on trans-  
 287 lating nighttime images to their corresponding daytime counterparts while preserving semantic  
 288 consistency. To achieve this objective, the proposed method begins by disentangling the degradation  
 289 presented in nighttime images, which is the key insight of our method. To achieve this, we contribute  
 290 a degradation disentanglement module and a degradation-aware contrastive learning module. Our  
 291 method outperforms the existing state-of-the-art, which shows the effectiveness of N2D3 and the  
 292 superiority of the insight to disentangle the degradation.

293 **References**

- 294 [1] Asha Anoosheh, Torsten Sattler, Radu Timofte, Marc Pollefeys, and Luc Van Gool. Night-to-day  
295 image translation for retrieval-based localization. In *2019 International Conference on Robotics  
296 and Automation (ICRA)*, pages 5958–5964. IEEE, 2019.
- 297 [2] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous  
298 convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.
- 299 [3] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo.  
300 Stargan: Unified generative adversarial networks for multi-domain image-to-image translation.  
301 In *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, pages 8789–8797,  
302 2018.
- 303 [4] T. Cover and P. Hart. Nearest neighbor pattern classification. *IEEE Transactions on Information  
304 Theory*, 13(1):21–27, 1967.
- 305 [5] Zhentao Fan, Xianhao Wu, Xiang Chen, and Yufeng Li. Learning to see in nighttime driving  
306 scenes with inter-frequency priors. In *Proceedings of the IEEE/CVF Conference on Computer  
307 Vision and Pattern Recognition*, pages 4217–4224, 2023.
- 308 [6] Chunle Guo, Chongyi Li, Jichang Guo, Chen Change Loy, Junhui Hou, Sam Kwong, and  
309 Runmin Cong. Zero-reference deep curve estimation for low-light image enhancement. In *Proc.  
310 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1780–1789, 2020.
- 311 [7] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.  
312 Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in  
313 Neural Information Processing Systems*, 30, 2017.
- 314 [8] Jinhui Hou, Zhiyu Zhu, Junhui Hou, Hui Liu, Huanqiang Zeng, and Hui Yuan. Global structure-  
315 aware diffusion process for low-light image enhancement. *Advances in Neural Information  
316 Processing Systems*, 36, 2024.
- 317 [9] Xueqi Hu, Xinyue Zhou, Qiusheng Huang, Zhengyi Shi, Li Sun, and Qingli Li. Qs-attn: Query-  
318 selected attention for contrastive learning in i2i translation. In *Proceedings of the IEEE/CVF  
319 Conference on Computer Vision and Pattern Recognition*, pages 18291–18300, 2022.
- 320 [10] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei Efros. Image-to-image translation with  
321 conditional adversarial networks. In *Proc. IEEE Conference on Computer Vision and Pattern  
322 Recognition*, pages 1125–1134, 2017.
- 323 [11] Somi Jeong, Youngjung Kim, Eungbean Lee, and Kwanghoon Sohn. Memory-guided unsuper-  
324 vised image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer  
325 Vision and Pattern Recognition*, pages 6558–6567, 2021.
- 326 [12] Hai Jiang, Ao Luo, Haoqiang Fan, Songchen Han, and Shuaicheng Liu. Low-light image  
327 enhancement with wavelet-based diffusion models. *ACM Transactions on Graphics (TOG)*,  
328 42(6):1–14, 2023.
- 329 [13] Yifan Jiang, Xinyu Gong, Ding Liu, Yu Cheng, Chen Fang, Xiaohui Shen, Jianchao Yang, Pan  
330 Zhou, and Zhangyang Wang. Enlightengan: Deep light enhancement without paired supervision.  
331 *IEEE Transactions on Image Processing*, 30:2340–2349, 2021.
- 332 [14] Mikhail Kennerley, Jian-Gang Wang, Bharadwaj Veeravalli, and Robby T Tan. 2pcnet: Two-  
333 phase consistency training for day-to-night unsupervised domain adaptive object detection. In  
334 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
335 11484–11493, 2023.
- 336 [15] Junho Kim, Minjae Kim, Hyeonwoo Kang, and Kwanghee Lee. U-gat-it: Unsupervised  
337 generative attentional networks with adaptive layer-instance normalization for image-to-image  
338 translation. *arXiv preprint arXiv:1907.10830*, 2019.

- 339 [16] Soohyun Kim, Jongbeom Baek, Jihye Park, Gyeongnyeon Kim, and Seungryong Kim. In-  
340 staformer: Instance-aware image-to-image translation with transformer. In *Proceedings of*  
341 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18321–18331,  
342 2022.
- 343 [17] Paul Kubelka. Ein beitrage zur optik der farbanstriche (contribution to the optic of paint).  
344 *Zeitschrift fur technische Physik*, 12:593–601, 1931.
- 345 [18] Jeong-gi Kwak, Youngsaeng Jin, Yuanming Li, Dongsik Yoon, Donghyeon Kim, and Hanseok  
346 Ko. Adverse weather image translation with asymmetric and uncertainty-aware gan. *arXiv*  
347 *preprint arXiv:2112.04283*, 2021.
- 348 [19] Michael J. Milford and Gordon. F. Wyeth. Seqslam: Visual route-based navigation for sunny  
349 summer days and stormy winter nights. In *2012 IEEE International Conference on Robotics*  
350 *and Automation*, pages 1643–1649, 2012.
- 351 [20] Taesung Park, Alexei Efros, Richard Zhang, and Jun-Yan Zhu. Contrastive learning for unpaired  
352 image-to-image translation. In *European Conference on Computer Vision*, pages 319–345,  
353 2020.
- 354 [21] Aashish Sharma and Robby T Tan. Nighttime visibility enhancement by increasing the dynamic  
355 range and suppression of light effects. In *Proceedings of the IEEE/CVF Conference on Computer*  
356 *Vision and Pattern Recognition*, pages 11977–11986, 2021.
- 357 [22] Seokbeom Song, Suhyeon Lee, Hongje Seong, Kyoungwon Min, and Euntai Kim. Shunit: Style  
358 harmonization for unpaired image-to-image translation. *Proceedings of the AAAI Conference*  
359 *on Artificial Intelligence*, 37(2):2292–2302, Jun. 2023.
- 360 [23] Tao Wang, Kaihao Zhang, Tianrun Shen, Wenhan Luo, Bjorn Stenger, and Tong Lu. Ultra-  
361 high-definition low-light image enhancement: A benchmark and transformer-based method. In  
362 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 2654–2662,  
363 2023.
- 364 [24] Weilun Wang, Wengang Zhou, Jianmin Bao, Dong Chen, and Houqiang Li. Instance-wise hard  
365 negative example generation for contrastive learning in unpaired image-to-image translation.  
366 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14020–  
367 14029, 2021.
- 368 [25] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment:  
369 from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–  
370 612, 2004.
- 371 [26] Youya Xia, Josephine Monica, Wei-Lun Chao, Bharath Hariharan, Kilian Q Weinberger, and  
372 Mark Campbell. Image-to-image translation for autonomous driving from coarsely-aligned  
373 image pairs. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages  
374 7756–7762. IEEE, 2023.
- 375 [27] Shaoan Xie, Qirong Ho, and Kun Zhang. Unsupervised image-to-image translation with density  
376 changing regularization. In *Advances in Neural Information Processing Systems*, 2022.
- 377 [28] Shaoan Xie, Yanwu Xu, Mingming Gong, and Kun Zhang. Unpaired image-to-image translation  
378 with shortest path regularization. In *Proceedings of the IEEE/CVF Conference on Computer*  
379 *Vision and Pattern Recognition (CVPR)*, pages 10177–10187, June 2023.
- 380 [29] Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht  
381 Madhavan, and Trevor Darrell. Bdd100k: A diverse driving dataset for heterogeneous mul-  
382 titask learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*  
383 *recognition*, pages 2636–2645, 2020.
- 384 [30] Zhenjie Yu, Shuang Li, Yirui Shen, Chi Harold Liu, and Shuigen Wang. On the difficulty of  
385 unpaired infrared-to-visible video translation: Fine-grained content-rich patches transfer. In  
386 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,  
387 pages 1631–1640, June 2023.

- 388 [31] Fangneng Zhan, Jiahui Zhang, Yingchen Yu, Rongliang Wu, and Shijian Lu. Modulated contrast  
389 for versatile image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
390 *and Pattern Recognition (CVPR)*, pages 18280–18290, June 2022.
- 391 [32] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unrea-  
392 sonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE*  
393 *conference on computer vision and pattern recognition*, pages 586–595, 2018.
- 394 [33] Ziqiang Zheng, Yang Wu, Xinran Han, and Jianbo Shi. Forkgan: Seeing into the rainy night. In  
395 *European conference on computer vision*, pages 155–170. Springer, 2020.
- 396 [34] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei Efros. Unpaired image-to-image transla-  
397 tion using cycle-consistent adversarial networks. In *Proc. IEEE International Conference on*  
398 *Computer Vision*, pages 2223–2232, 2017.

## 399 A Overview

400 This supplementary material is organized as follows. Appendix B provides additional details about  
 401 the proof that the invariant  $N_{\lambda^m x^n}$  is exclusively related to the illumination. Appendix C outlines the  
 402 limitations and failure case of N2D3. Appendix D illustrates the implementation details, including  
 403 N2D3 and other methods used in the experiments. Appendix E presents additional visualization  
 404 results.

## 405 B More Proof Details

406 We provide a detailed proof process to demonstrate **how the invariant  $N_{\lambda^m x^n}$  is exclusively related  
 407 to the illumination and can function as the light effect detector**. First, consider the following  
 408 equations, corresponding to Equation (5) in the main paper:

$$\begin{aligned} N_{\lambda^m x^n} &= \frac{\partial^{m+n-2}}{\partial \lambda^{m-1} \partial x^{n-1}} \frac{\partial}{\partial x} \left\{ \frac{1}{E(\lambda, x)} \frac{\partial E(\lambda, x)}{\partial \lambda} \right\} \\ &= \frac{\partial^{m+n-2}}{\partial \lambda^{m-1} \partial x^{n-1}} \frac{\partial}{\partial x} \left\{ \frac{1}{e(\lambda, x)} \frac{\partial e(\lambda, x)}{\partial \lambda} + \frac{1}{R(\lambda)C(x)} \frac{\partial R(\lambda)C(x)}{\partial \lambda} \right\}, \end{aligned} \quad (18)$$

409 by applying the additivity of linear differential operators, the first term represents the invariants only  
 410 related to the illumination. The second term can be simplified by applying the chain rule as follows:

$$\begin{aligned} &\frac{\partial}{\partial x} \left\{ \frac{1}{R(\lambda)C(x)} \frac{\partial R(\lambda)C(x)}{\partial \lambda} \right\} \\ &= \frac{1}{R(\lambda)^2 C(x)^2} \left( \frac{\partial^2 \{R(\lambda)C(x)\}}{\partial \lambda \partial x} \cdot R(\lambda)C(x) - \frac{\partial \{R(\lambda)C(x)\}}{\partial \lambda} \cdot \frac{\partial \{R(\lambda)C(x)\}}{\partial x} \right) \\ &= \frac{1}{R(\lambda)^2 C(x)^2} \left( \frac{\partial R(\lambda)}{\partial \lambda} \frac{\partial C(x)}{\partial x} \cdot R(\lambda)C(x) - \frac{\partial R(\lambda)}{\partial \lambda} C(x) \cdot R(\lambda) \frac{\partial C(x)}{\partial x} \right) = 0. \end{aligned} \quad (19)$$

411 Finally, we conclude that the invariant  $N_{\lambda^m x^n}$  is **exclusively related to the illumination** and can be  
 412 formulated as follows:

$$\begin{aligned} N_{\lambda^m x^n} &= \frac{\partial^{m+n-2}}{\partial \lambda^{m-1} \partial x^{n-1}} \frac{\partial}{\partial x} \left\{ \frac{1}{E(\lambda, x)} \frac{\partial E(\lambda, x)}{\partial \lambda} \right\} \\ &= \frac{\partial^{m+n-1}}{\partial \lambda^{m-1} \partial x^n} \left\{ \frac{1}{e(\lambda, x)} \frac{\partial e(\lambda, x)}{\partial \lambda} \right\}. \end{aligned} \quad (20)$$



Figure 7: Failure Cases of N2D3: Our method struggles to handle various other types of degradation.

## 413 C Limitations and Failure Case

414 Despite the superior performance of N2D3 in Night2Day, it still exhibits certain limitations. On the  
 415 one hand, this work focuses solely on addressing light degradation, while nighttime environments  
 416 encompass various other types of degradation, including blur caused by rain, motion, and other

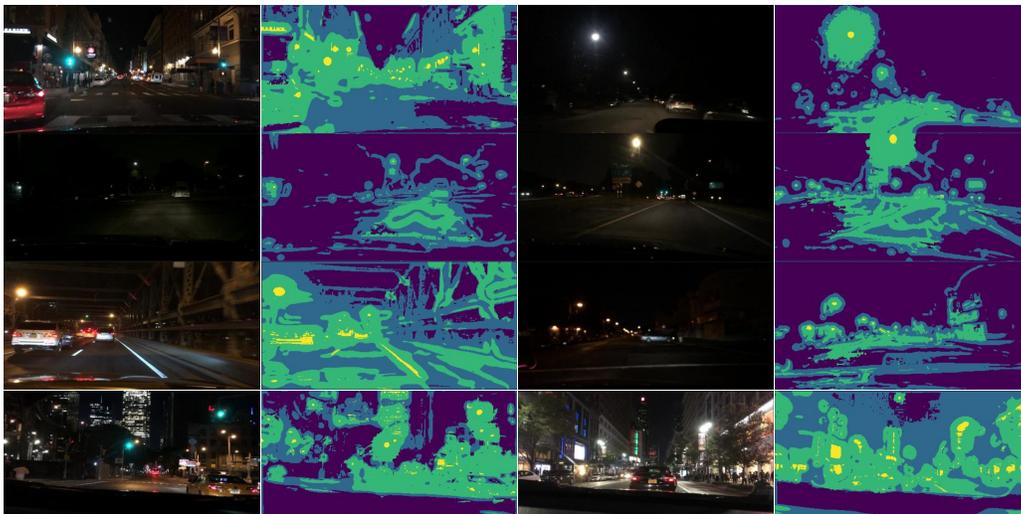


Figure 8: More disentanglement results. The first and third rows display nighttime images, while the second and fourth rows show the corresponding degradation disentanglement results. The color progression from blue, light blue, green to yellow corresponds to the following regions: darkness, well-lit, light effects, and high-light.



Figure 9: Qualitative comparison ablation results.

417 factors. Our method currently struggles to handle these situations effectively. On the other hand, the  
 418 limitations of visible imaging in night vision arise from the scarcity of photos captured in low-light  
 419 conditions, as illustrated by the failure cases presented in Figure 7. Future advancements in night  
 420 vision will likely incorporate additional modalities, such as infrared images, radar, and other sensor  
 421 data, to overcome these challenges and improve performance.

## 422 D Implementation Details

423 **Training Details.** We adopt the *resnet\_9blocks*, a ResNet-based model with nine residual blocks, as  
 424 the backbone for generator  $G$ . Additionally, we utilize the patch-wise discriminator  $D$  following  
 425 PatchGAN[10]. To conduct degradation-aware contrastive learning on multiple layers, we extract  
 426 features from 5 layers of the generator  $G$  encoder, as done in [20]. These layers include RGB pixels,  
 427 the first and second downsampling convolution, and the first and fifth residual block. For the features  
 428 of each layer, we apply a 2-layer MLP to acquire final 256-dimensional features. These features are  
 429 then utilized in our degradation-aware contrastive learning.

430 All the comparison methods are reproduced using their released source code with default settings.  
 431 Training procedures are consistent across all methods. All models are trained using the Adaptive  
 432 Moment Estimation optimizer with an initial learning rate of  $10^{-4}$ , a momentum of 0.9, and weight  
 433 decay of  $10^{-4}$ . For the BDD100K dataset, training consists of 10 epochs with the initial learning  
 434 rate, followed by another 10 epochs with a decreased learning rate using the polynomial annealing  
 435 procedure with a power of 0.9. On the Alderley dataset, given the limited training data compared  
 436 to BDD100K, we extend the training to 20 epochs with the initial learning rate and an additional



Figure 10: More qualitative comparison results on the Alderley dataset.

437 20 epochs with the decayed learning rate. All the experiments are run on a single A100 GPU with  
438 80GB of memory. Training our method with a smaller patch size and batch size on a device with less  
439 memory is feasible.

440 **Evaluation Details.** In the evaluation, we compute the *Fréchet Inception Distance* (FID) [7],  
441 Structural Similarity Index (SSIM) [25], and Learned Perceptual Image Patch Similarity (LPIPS)  
442 [32] scores on  $256 \times 512$  images. Partial FID scores are provided by ForkGAN [33], and all SSIM  
443 and LPIPS scores are reproduced by us.

444 Semantic segmentation evaluation are conducted as follows. First, we use Deeplabv3 pretrained  
445 on the Cityscapes dataset as the semantic segmentation model [2]. The model is provided by  
446 <https://github.com/open-mmlab/msegmentation> with an R-18-D8 backbone and trained at  
447 a resolution of  $512 \times 1024$ . Second, we perform  $512 \times 1024$  Night2Day translation to obtain the  
448 generation results. Finally, we infer the semantic segmentation on the generated daytime images.

## 449 E More Visualization Results

450 **More Ablation Visualization Results.** We provide ablation visualization results on both Alderley  
451 and BDD100K in Figure 9. The complete method is presented along with ablation studies on the  
452 invariant  $N$  and without degradation-aware reweighting. All the modules contribute to improving the  
453 ability to maintain semantic consistency.

454 **More Disentanglement Results.** We provide additional disentanglement results in Figure 8. Our  
455 disentanglement methods offer a comprehensive representation of different illumination degradation  
456 types in various nighttime scenes.

457 **More Qualitative Comparison.** We present more qualitative comparisons in Figure 10 and Figure 11  
458 alongside other methods. Our method demonstrates visually pleasing results under various nighttime  
459 conditions.

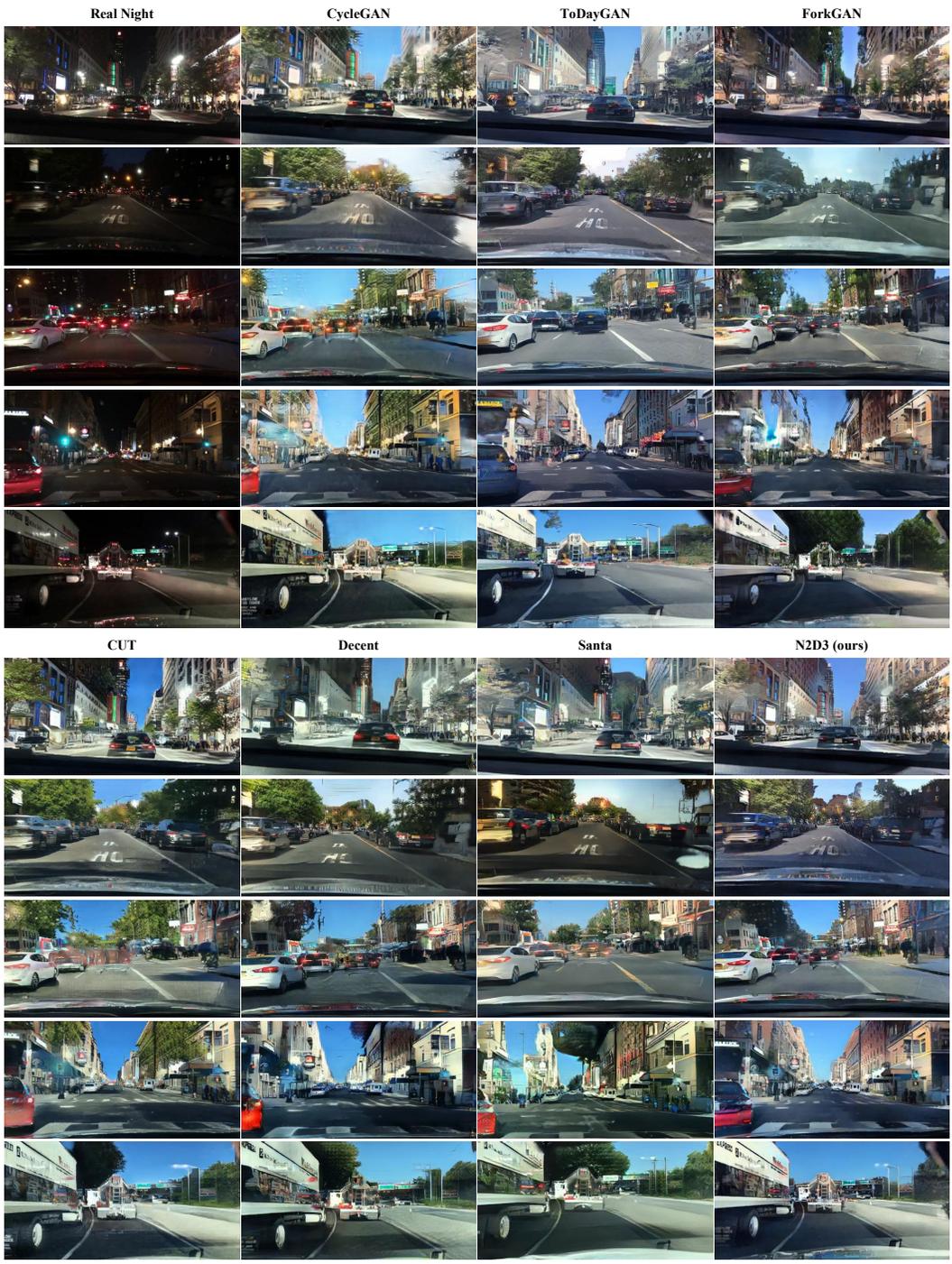


Figure 11: More qualitative comparison results on the BDD100K dataset.

460 **NeurIPS Paper Checklist**

461 **1. Claims**

462 Question: Do the main claims made in the abstract and introduction accurately reflect the  
463 paper's contributions and scope?

464 Answer: [Yes]

465 Justification: We claim our main contribution as N2D3, which achieves SOTA performance  
466 by bridging the domain gap between nighttime and daytime in a degradation-aware manner.

467 Guidelines:

- 468 • The answer NA means that the abstract and introduction do not include the claims  
469 made in the paper.
- 470 • The abstract and/or introduction should clearly state the claims made, including the  
471 contributions made in the paper and important assumptions and limitations. A No or  
472 NA answer to this question will not be perceived well by the reviewers.
- 473 • The claims made should match theoretical and experimental results, and reflect how  
474 much the results can be expected to generalize to other settings.
- 475 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
476 are not attained by the paper.

477 **2. Limitations**

478 Question: Does the paper discuss the limitations of the work performed by the authors?

479 Answer: [Yes]

480 Justification: We discuss our limitation in degradations beyond light and low-light image  
481 scarcity in the appendix.

482 Guidelines:

- 483 • The answer NA means that the paper has no limitation while the answer No means that  
484 the paper has limitations, but those are not discussed in the paper.
- 485 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 486 • The paper should point out any strong assumptions and how robust the results are to  
487 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
488 model well-specification, asymptotic approximations only holding locally). The authors  
489 should reflect on how these assumptions might be violated in practice and what the  
490 implications would be.
- 491 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
492 only tested on a few datasets or with a few runs. In general, empirical results often  
493 depend on implicit assumptions, which should be articulated.
- 494 • The authors should reflect on the factors that influence the performance of the approach.  
495 For example, a facial recognition algorithm may perform poorly when image resolution  
496 is low or images are taken in low lighting. Or a speech-to-text system might not be  
497 used reliably to provide closed captions for online lectures because it fails to handle  
498 technical jargon.
- 499 • The authors should discuss the computational efficiency of the proposed algorithms  
500 and how they scale with dataset size.
- 501 • If applicable, the authors should discuss possible limitations of their approach to  
502 address problems of privacy and fairness.
- 503 • While the authors might fear that complete honesty about limitations might be used by  
504 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
505 limitations that aren't acknowledged in the paper. The authors should use their best  
506 judgment and recognize that individual actions in favor of transparency play an impor-  
507 tant role in developing norms that preserve the integrity of the community. Reviewers  
508 will be specifically instructed to not penalize honesty concerning limitations.

509 **3. Theory Assumptions and Proofs**

510 Question: For each theoretical result, does the paper provide the full set of assumptions and  
511 a complete (and correct) proof?

512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565

Answer: [Yes]

Justification: We provide the full set of assumptions and complete proofs in both Section 3.1 and Appendix B .

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: All the information needed to reproduce the main experimental results is included in the Section 3 and Appendix D.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

566 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
567 tions to faithfully reproduce the main experimental results, as described in supplemental  
568 material?

569 Answer: [No]

570 Justification: Code will be released latter.

571 Guidelines:

- 572 • The answer NA means that paper does not include experiments requiring code.
- 573 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
574 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 575 • While we encourage the release of code and data, we understand that this might not be  
576 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
577 including code, unless this is central to the contribution (e.g., for a new open-source  
578 benchmark).
- 579 • The instructions should contain the exact command and environment needed to run to  
580 reproduce the results. See the NeurIPS code and data submission guidelines ([https:  
581 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 582 • The authors should provide instructions on data access and preparation, including how  
583 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 584 • The authors should provide scripts to reproduce all experimental results for the new  
585 proposed method and baselines. If only a subset of experiments are reproducible, they  
586 should state which ones are omitted from the script and why.
- 587 • At submission time, to preserve anonymity, the authors should release anonymized  
588 versions (if applicable).
- 589 • Providing as much information as possible in supplemental material (appended to the  
590 paper) is recommended, but including URLs to data and code is permitted.

## 591 6. Experimental Setting/Details

592 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
593 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
594 results?

595 Answer: [Yes]

596 Justification: The training details and dataset information are provided in Section 4.

597 Guidelines:

- 598 • The answer NA means that the paper does not include experiments.
- 599 • The experimental setting should be presented in the core of the paper to a level of detail  
600 that is necessary to appreciate the results and make sense of them.
- 601 • The full details can be provided either with the code, in appendix, or as supplemental  
602 material.

## 603 7. Experiment Statistical Significance

604 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
605 information about the statistical significance of the experiments?

606 Answer: [No]

607 Justification: Error bars are not reported because it would be too computationally expensive.  
608 We report our results using a fixed random seed.

609 Guidelines:

- 610 • The answer NA means that the paper does not include experiments.
- 611 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
612 dence intervals, or statistical significance tests, at least for the experiments that support  
613 the main claims of the paper.
- 614 • The factors of variability that the error bars are capturing should be clearly stated (for  
615 example, train/test split, initialization, random drawing of some parameter, or overall  
616 run with given experimental conditions).

- 617 • The method for calculating the error bars should be explained (closed form formula,  
618 call to a library function, bootstrap, etc.)
- 619 • The assumptions made should be given (e.g., Normally distributed errors).
- 620 • It should be clear whether the error bar is the standard deviation or the standard error  
621 of the mean.
- 622 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
623 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
624 of Normality of errors is not verified.
- 625 • For asymmetric distributions, the authors should be careful not to show in tables or  
626 figures symmetric error bars that would yield results that are out of range (e.g. negative  
627 error rates).
- 628 • If error bars are reported in tables or plots, The authors should explain in the text how  
629 they were calculated and reference the corresponding figures or tables in the text.

## 630 8. Experiments Compute Resources

631 Question: For each experiment, does the paper provide sufficient information on the com-  
632 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
633 the experiments?

634 Answer: [Yes]

635 Justification: We report the compute resources in Appendix D.

636 Guidelines:

- 637 • The answer NA means that the paper does not include experiments.
- 638 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
639 or cloud provider, including relevant memory and storage.
- 640 • The paper should provide the amount of compute required for each of the individual  
641 experimental runs as well as estimate the total compute.
- 642 • The paper should disclose whether the full research project required more compute  
643 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
644 didn't make it into the paper).

## 645 9. Code Of Ethics

646 Question: Does the research conducted in the paper conform, in every respect, with the  
647 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

648 Answer: [Yes]

649 Justification: The research conducted in this paper conforms, in every respect, with the  
650 NeurIPS Code of Ethics.

651 Guidelines:

- 652 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 653 • If the authors answer No, they should explain the special circumstances that require a  
654 deviation from the Code of Ethics.
- 655 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
656 eration due to laws or regulations in their jurisdiction).

## 657 10. Broader Impacts

658 Question: Does the paper discuss both potential positive societal impacts and negative  
659 societal impacts of the work performed?

660 Answer: [Yes]

661 Justification: The societal impacts are discussed in the manuscript and appendix.

662 Guidelines:

- 663 • The answer NA means that there is no societal impact of the work performed.
- 664 • If the authors answer NA or No, they should explain why their work has no societal  
665 impact or why the paper does not address societal impact.

- 666
- 667
- 668
- 669
- 670
- 671
- 672
- 673
- 674
- 675
- 676
- 677
- 678
- 679
- 680
- 681
- 682
- 683
- 684
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
  - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
  - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 685 11. Safeguards

686 Question: Does the paper describe safeguards that have been put in place for responsible  
687 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
688 image generators, or scraped datasets)?

689 Answer: [NA]

690 Justification: Our model does not have such risks, and all the datasets used in the experiments  
691 are open-source benchmarks in this field.

692 Guidelines:

- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701
- 702
- The answer NA means that the paper poses no such risks.
  - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
  - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
  - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 703 12. Licenses for existing assets

704 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
705 the paper, properly credited and are the license and terms of use explicitly mentioned and  
706 properly respected?

707 Answer: [Yes]

708 Justification: The code and data are properly credited, and the license and terms of use are  
709 explicitly mentioned and properly documented.

710 Guidelines:

- 711
- 712
- 713
- 714
- 715
- 716
- 717
- The answer NA means that the paper does not use existing assets.
  - The authors should cite the original paper that produced the code package or dataset.
  - The authors should state which version of the asset is used and, if possible, include a URL.
  - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
  - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- 718
- 719
- 720
- 721
- 722
- 723
- 724
- 725
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
  - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
  - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 726 13. **New Assets**

727 Question: Are new assets introduced in the paper well documented and is the documentation  
728 provided alongside the assets?

729 Answer: [Yes]

730 Justification: The code introduced in the paper is well-documented, and the documentation  
731 is provided alongside it.

732 Guidelines:

- 733
- 734
- 735
- 736
- 737
- 738
- 739
- 740
- The answer NA means that the paper does not release new assets.
  - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
  - The paper should discuss whether and how consent was obtained from people whose asset is used.
  - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 741 14. **Crowdsourcing and Research with Human Subjects**

742 Question: For crowdsourcing experiments and research with human subjects, does the paper  
743 include the full text of instructions given to participants and screenshots, if applicable, as  
744 well as details about compensation (if any)?

745 Answer: [NA]

746 Justification: The paper does not involve crowdsourcing nor research with human subjects.

747 Guidelines:

- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
  - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

### 756 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 757 Subjects**

758 Question: Does the paper describe potential risks incurred by study participants, whether  
759 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
760 approvals (or an equivalent approval/review based on the requirements of your country or  
761 institution) were obtained?

762 Answer: [NA]

763 Justification: The paper does not involve crowdsourcing nor research with human subjects.

764 Guidelines:

- 765
- 766
- 767
- 768
- 769
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

770  
771  
772  
773  
774

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.