

GlowGAN: Unsupervised Learning of HDR Images from LDR Images in the Wild

Chao Wang¹, Ana Serrano², Xingang Pan¹, Bin Chen¹, Hans-Peter Seidel¹
 Christian Theobalt¹, Karol Myszkowski¹, Thomas Leimkühler¹

¹ Max Planck Institute for Informatics, Germany, ² Universidad de Zaragoza, Spain



Figure 1. Typical image datasets have a low dynamic range (LDR), resulting in over- and underexposed pixels (a). We devise a high-dynamic-range (HDR) generator trained only on in-the-wild LDR data. Our HDR samples (b, tonemapped for display) exhibit details at all brightness levels. Our model can be used for inverse tone mapping to recover details in large-scale saturated regions (c).

Abstract

Most in-the-wild images are stored in Low Dynamic Range (LDR) form, serving as a partial observation of the High Dynamic Range (HDR) visual world. Despite limited dynamic range, these LDR images are often captured with different exposures, implicitly containing information about the underlying HDR image distribution. Inspired by this intuition, in this work we present, to the best of our knowledge, the first method for learning a generative model of HDR images from in-the-wild LDR image collections in a fully unsupervised manner. The key idea is to train a generative adversarial network (GAN) to generate HDR images which, when projected to LDR under various exposures, are indistinguishable from real LDR images. The projection from HDR to LDR is achieved via a camera model that captures the stochasticity in exposure and camera response function. Experiments show that our method GlowGAN can synthesize photorealistic HDR images in many challenging cases such as landscapes, lightning, or windows, where previous

supervised generative models produce overexposed images. We further demonstrate the new application of unsupervised inverse tone mapping (ITM) enabled by GlowGAN. Our ITM method does not need HDR images or paired multi-exposure images for training, yet it reconstructs more plausible information for overexposed regions than state-of-the-art supervised learning models trained on such data.

1. Introduction

High Dynamic Range (HDR) images [55] are capable of capturing and displaying much richer appearance information than Low Dynamic Range (LDR) images, thus playing an important role in image representation and visualization. The most popular method to acquire HDR images is multiple exposure blending, which requires capturing a set of LDR images of the same scene with different exposures [15, 50, 57]. However, this is time and effort intensive and only suitable for static scenes. Due to this limitation, existing HDR image

datasets only cover limited scene categories and have much fewer images than LDR datasets. Thus, supervised learning methods [16, 17, 26, 39, 41, 43, 45, 60, 69, 71] that reconstruct an HDR image from an LDR image are constrained by the HDR datasets and cannot extend to cases where no HDR training data is available, e.g., lightning.

While HDR images are hard to collect, it is much easier to collect a large number of LDR images from the Internet. This motivates us to investigate a new unsupervised learning problem: *Can we learn to reconstruct HDR images from in-the-wild LDR images?* The LDR images do not need to depict the same scene, it is enough if they contain a roughly similar class of scenes (e.g., landscapes) with various exposures. This weak multi-exposure assumption is often naturally satisfied for in-the-wild LDR datasets as images can come from different camera parameters or different adjustments of the auto-exposure mode. This problem is challenging as only one exposure is available for each scene, therefore, a way to merge the multi-exposure information spread across different scenes is required.

In this work, we address this challenge via *GlowGAN*, which, to our knowledge, is the first method to learn an HDR generative model from in-the-wild LDR image collections in a fully unsupervised manner. GlowGAN uses adversarial training of an HDR generator with a discriminator that operates merely in LDR. Specifically, the generator produces an HDR image, which is projected to LDR via a camera model and is then sent to the discriminator as a fake image for adversarial training. The camera model consists of multiplying the HDR sample with an exposure value, clipping the dynamic range, and applying a camera response function (CRF). Importantly, during training, we use a randomly sampled exposure from a prior Gaussian distribution when projecting HDR to LDR. This requires the generated HDR images to be realistic under any possible exposure, thus only valid HDR samples would satisfy this “multi-exposure constraint”. Furthermore, we also model the stochasticity in the non-linear camera response by randomly sampling CRFs according to a well-established parametric distribution [16, 19]. During the inference process, we can disable the camera model so that the generator produces HDR imagery directly.

We conduct extensive experiments on several datasets collected from the Internet, including landscapes, windows, lightning, fireplaces, and fireworks. By training on these LDR images, GlowGAN successfully learns to generate high-quality HDR images that capture rich appearance information from dark to very bright regions. These details can be presented on HDR displays, or via suitable tone mapping to create appealing imagery. In contrast, previous unconditional GANs tend to miss information in over- or under-exposed regions.

By modelling a distribution of HDR images, GlowGAN paves the way for new applications such as unsupervised

inverse tone mapping (ITM). ITM aims to reconstruct an HDR image from a single-exposure LDR input, where a key challenge is to restore the flat-white overexposed regions [16, 60]. We can use a pre-trained GlowGAN as a prior and apply GAN inversion to optimize latent code and exposure, making the model generate the corresponding HDR image for the input LDR image. An exciting result is that our method can, without using any HDR imagery or paired multi-exposure data, reconstruct starkly more plausible information for large overexposed regions than other supervised learning approaches trained on such data. Furthermore, the HDR samples generated by GlowGAN can be used as versatile environment maps in rendering. Our contributions are summarized as follows:

- We are the first to present unsupervised learning of HDR images from in-the-wild LDR images. This gets rid of the reliance on ground truth HDR images that are much harder to collect.
- To achieve this, we propose a novel GlowGAN, which bridges HDR space and LDR space via a camera model. GlowGAN can synthesize diverse high-quality images with a much higher dynamic range than vanilla GANs, opening up new avenues for getting cheap abundant HDR data.
- Using GlowGAN as a prior, we design an unsupervised inverse tone mapping method (ITM), which reconstructs large overexposed regions significantly better than the state-of-the-art fully-supervised approaches.

The supplementary material is provided in [this link](#). Our code and pretrained models will be released.

2. Related Work

2.1. High Dynamic Range Imaging

The real world has a vast dynamic range. Therefore, HDR imaging [55] is crucial for creating and manipulating immersive viewing experiences. While HDR capture is cumbersome and HDR displays are not yet commonplace in most environments [63], the representation of visual information free from the limitations of typical LDR encodings has evolved significantly in recent years [67]. Working with HDR content is crucial for rendering [14] and has been shown to be beneficial for 3D scene reconstruction: An HDR representation can naturally handle multi-exposure [25, 27] and raw data [49]. Reconstructing HDR in conjunction with an explicit tone mapping module can compensate for poorly calibrated cameras [59].

A particular interest has evolved around the conversion between HDR and LDR content. Tone mapping, the transformation from HDR to LDR with as little information loss as possible, is a mature field with well-understood trade-offs [55, 65]. In contrast, inverse tone mapping (ITM) [5, 56],

the recovery of HDR content from LDR imagery, remains a challenging inverse problem. It typically involves multiple steps, including linearization, dynamic range expansion, reconstruction of over- and underexposed regions, artifact reduction, and color correction [3]. Among these, the reconstruction of saturated pixels is considered the most challenging [16, 60], as it requires the hallucination of content [66].

Early ITM works considered expansion curves using either global [2, 38, 46, 47] or content-driven operators [4–6, 48, 56, 56], without explicitly reconstructing saturated regions. More recent learning-based solutions can be categorized into two streams: A neural network either predicts the HDR image directly [12, 16, 43, 45, 60, 69, 71], or it predicts multiple LDR images with different exposures [17, 26, 39–41], which are subsequently merged into an HDR image [15, 50, 57]. ITM and exposure fusion can be combined with adversarial training [40, 52, 71]. Also, the extension to video ITM has been explored extensively, leveraging inter-frame consistency [7, 21, 22, 28, 35, 36].

All learning-based techniques discussed above rely on supervision from paired LDR–HDR training data. This constitutes a fundamental problem: HDR image data is hard to obtain and therefore naturally scarce. Further, most HDR capture techniques require static content, which significantly restricts the applicability of learning-based methods to arbitrary scenes. In contrast, we are the first to train an HDR image generator *unsupervised* from an LDR dataset. We believe this is an important step towards solving two main problems in HDR imaging: First, our generator can synthesize an abundance of HDR samples, alleviating the scarcity of HDR content. Second, our system allows to perform ITM with a significant improvement in the reconstruction of overexposed regions.

2.2. Lossy Generative Adversarial Networks

Generative Adversarial Networks (GANs) [18] are very successful in modelling distributions of images with high visual fidelity. The StyleGAN family [31–34] marks the current state of the art, scaling to high resolutions, while the recent extension StyleGAN-XL [61] allows for unprecedented diversity in the generated content. To this date, (unconditional) GANs operate in LDR, since high-quality HDR data at the scale required for successful training is difficult to obtain. We propose a simple modification to the GAN training pipeline, which allows to train an HDR generator from readily available LDR data only.

AmbientGAN [10] has demonstrated that it is feasible to train a generative model from lossy measurements, i.e., a GAN can be trained from degraded samples, as long as the stochastic properties of the degradation are known. This concept has been used to learn a generator for clean images from noisy data [30], or for all-in-focus images from data that contains shallow depth of field [29]. Most prominently, the

idea has been applied to learn 3D generators from 2D images by explicitly modeling the projection from 3D to 2D using a distribution of extrinsic camera parameters [11, 23, 42, 51, 62, 64]. We follow the paradigm of injecting a degradation model into the GAN training pipeline, by devising a novel model of the distribution of processing steps in a digital camera, converting HDR radiance into LDR pixel intensities.

3. Method

We train a GAN [18] to capture the distribution of HDR images in a domain (e.g., landscapes) by combining it with a stochastic camera model that transforms the generated HDR images into their LDR counterparts. The discriminator is only fed LDR images, which allows the system to be trained on an easily accessible in-the-wild LDR image dataset. Our camera model can be inserted into any GAN model to yield HDR outputs as long as the LDR training dataset consists of images exhibiting different exposures across samples. We consider this a mild assumption, which in particular in-the-wild photo datasets easily satisfy. See Fig. 2 for an overview of our system. In the following, we describe our pipeline in detail (Sec. 3.1), before turning to our main application of unsupervised inverse tone mapping (Sec. 3.2).

3.1. GlowGAN

We seek to capture the unknown true distribution of HDR images p_{HDR} from samples of the distribution of LDR images p_{LDR} . Similar to the standard GAN setup, we achieve this by training a generator G which turns a random latent vector $\mathbf{z} \in \mathbb{R}^k$ into an HDR sample $\mathbf{r} \in \mathbb{R}_{\geq 0}^{H \times W \times 3}$, an RGB image with $H \times W$ pixels and no upper restriction on the value range. To train G , we inject a camera model $C \in \mathbb{R}_{\geq 0}^{H \times W \times 3} \rightarrow [0, 1]^{H \times W \times 3}$ into the adversarial training pipeline, turning the unbounded HDR image into an LDR image with values between 0 and 1. C captures the distribution p_{cam} of pixel-wise image transformations typically applied in a digital camera to convert incoming radiance values to final pixel intensities, including varying exposures, clipping, and varying non-linearities arising from the camera response function (CRF). The result of this process is an LDR image $\mathbf{l} = C(\mathbf{r}) = C(G(\mathbf{z}))$. The discriminator D is tasked with differentiating the fake samples \mathbf{l} from samples from the distribution of true LDR images p_{LDR} . Since the samples \mathbf{r} undergo stochastic projections from HDR to LDR, G is forced to produce valid HDR images, resulting in the distribution p_{HDR}^G of generated HDR images approaching the true distribution p_{HDR} [10].

As our generator and discriminator backbone, we choose the current state-of-the-art model StyleGAN-XL [61]. This model has been demonstrated to yield excellent image quality on diverse datasets. The generator model consists of two stages (left block in Fig. 2): First, a mapping network M_{θ_M} in the form of an MLP with parameters θ_M turns the

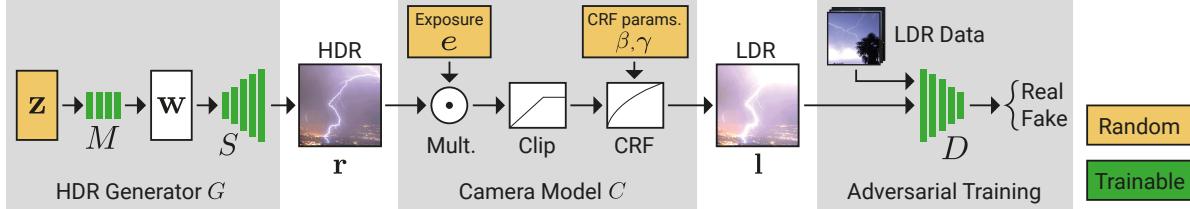


Figure 2. Overview of GlowGAN. The generator generates an HDR image \mathbf{r} from a random noise \mathbf{z} . Then a camera model C projects \mathbf{r} to an LDR image \mathbf{l} with a random exposure and CRF. The model is trained merely on in-the-wild LDR images in an adversarial manner.

initial random vector $\mathbf{z} \in \mathbb{R}^k$ into a more disentangled latent feature representation $\mathbf{w} \in W \subset \mathbb{R}^k$. Second, \mathbf{w} is fed into a synthesis CNN S_{θ_S} with parameters θ_S to yield the final HDR output $\mathbf{r} = S_{\theta_S}(M_{\theta_M}(\mathbf{z}))$. Notice that the particular choice of the generator and discriminator networks is orthogonal to our approach. Except for the camera model introduced in the next paragraphs, we use StyleGAN-XL without any modifications in architecture or training hyperparameters. We verify this in Sec. 4.2.

At the core of our method is the stochastic camera model C (central block in Fig. 2) which projects the HDR image \mathbf{r} onto an LDR counterpart \mathbf{l} . It is designed to model the distribution of typical processing steps in a digital camera:

$$C(\mathbf{r}) = \text{CRF}_{\beta, \gamma} \left(\min(2^{\frac{e}{2}} \cdot \mathbf{r}, 1) \right). \quad (1)$$

In the first step, we multiply each pixel of \mathbf{r} with a single global exposure value. The exposure is parameterized by the random variable e , capturing the exposure distribution of typical cameras arising from the combined effect of aperture, shutter speed, and sensor sensitivity (ISO). We do not have access to the true distribution of exposure parameters e , as images from in-the-wild image collections frequently do not have EXIF headers that would contain this information. Since the exposure is a combined effect of multiple factors (aperture, shutter, ISO, time of day, etc.), we choose to model e using a normal distribution, i.e., $e \sim \mathcal{N}(0, \sigma_e^2)$. The exposure variance σ_e^2 is the only hyper-parameter in our system and is analyzed in Sec. 4.2.

After applying a random exposure, the min operator in Eq. 1 clips large radiance values to 1, effectively flattening all definition in overexposed regions. Notice that this loss of information is precisely the reason why C is not invertible for individual images: By selecting an exposure, we only observe a bracket of radiance values. In contrast, our method seeks to invert C over the distribution of images and camera models [10].

The final component of our model in Eq. 1 is the non-linear sensor response. The CRF describes the mapping

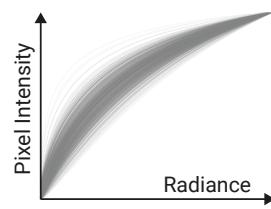


Figure 3. CRF Distribution.

from radiance values arriving at the sensor to pixel intensities stored in the final image. We follow the established distribution of Eilertsen et al. [16]:

$$\text{CRF}_{\beta, \gamma}(x) = \frac{(1 + \beta)x^\gamma}{\beta + x^\gamma},$$

with $\beta \sim \mathcal{N}(0.6, 0.1)$ and $\gamma \sim \mathcal{N}(0.9, 0.1)$ as obtained from the analysis of a large dataset [19]. We visualize the distribution of CRFs arising from this model in Fig. 3.

With the full stochastic camera model in place, our system can be trained in an adversarial fashion from scratch without further modifications to the standard GAN pipeline to yield an HDR generator trained only on an LDR image dataset.

3.2. Unsupervised Inverse Tone Mapping

In addition to producing HDR image samples, a trained GlowGAN can be used to perform unsupervised inverse tone mapping (ITM). While state-of-the-art ITM approaches typically rely on supervision from LDR–HDR image pairs, our method allows for recovering HDR images from their LDR counterparts without HDR data using GAN inversion [68, 72].

To obtain high-quality ITM results and to facilitate multimodality, we choose a per-image optimization-based approach to perform the inversion [1, 13]. Given an LDR image $\hat{\mathbf{l}}$, we consider the following optimization objective:

$$[e^*, \mathbf{w}^*, \theta_S^*] = \underset{e, \mathbf{w}, \theta_S}{\operatorname{argmin}} \Phi \left(C(S_{\theta_S}(\mathbf{w})), \hat{\mathbf{l}} \right). \quad (2)$$

Here, we jointly optimize over exposure e and the latent code \mathbf{w} while fine-tuning the synthesis network parameters θ_S to obtain a faithful match between the target LDR image $\hat{\mathbf{l}}$ and its reconstruction using our pipeline. Φ denotes the discrepancy measure between the two images. Using the Adam optimizer [37] with standard parameters, we proceed in two stages: In the first stage, we exclude the generator weights θ_S from the optimization and measure image discrepancy Φ using the LPIPS perceptual distance [70]. In the second stage, we only optimize (fine-tune) θ_S using the pivotal tuning technique of Roich et al. [58], with Φ being the sum of a pixel-wise ℓ_2 loss and the LPIPS perceptual distance. Following most previous work, we relax \mathbf{w} to explore the extended latent space W^+ [1]. We did not find it necessary



Figure 4. (a) Dynamic range comparison between our GlowGAN (top) and a vanilla GAN (bottom) when varying the exposure values (EV) of the image. (b) HDR samples generated by our GlowGAN for the five tested datasets (please see Fig. 1 for more samples). Our method generates high quality images with much higher dynamic range and without overexposure.

to optimize over the CRF parameters β, γ for high-quality results and consequently fix them to their mean values. More details on the optimization can be found in the supplemental. Upon completion of the optimization, we obtain the HDR version of $\hat{\mathbf{l}}$ via $\mathbf{r}^* = S_{\theta_S}(\mathbf{w}^*)$.

Following the lossy projection in Eq. 1, the mapping from LDR to HDR images is not unique: Overexposed regions in the LDR image can be explained by many different HDR images. Our system allows capturing this multi-modality by running the optimization of Eq. 2 multiple times with different parameter initializations for \mathbf{w} and e . Specifically, we initialize each optimization run with $\mathbf{w} = M_{\theta_M}(\mathbf{z})$, where \mathbf{z} is a normally distributed random vector. This allows us to obtain multiple plausible HDR solutions, which almost exclusively differ in the overexposed regions.

Obtaining pixel-accurate GAN inversion results is challenging [1, 9]. Fortunately, in most cases, we are only interested in hallucinating content in the saturated image regions, while well-exposed pixels can be re-used after linearization. Therefore, optionally, we diminish potential distortions arising from the inversion by blending the linearized original LDR image $\hat{\mathbf{l}}$ with the reconstructed HDR result \mathbf{r}^* [16, 60] as follows:

$$\mathbf{r}_{\text{blend}}^* = e^* \cdot (\mathbf{m} \odot \mathbf{r}^*) + (1 - \mathbf{m}) \odot \text{CRF}^{-1}(\hat{\mathbf{l}}).$$

Here, \odot denotes the Hadamard product and \mathbf{m} is a soft mask, indicating saturated pixels in $\hat{\mathbf{l}}$, which we compute for each pixel i following Eilertsen et al. [16]:

$$\mathbf{m}_i = \frac{\max \left(0, \max_c \hat{\mathbf{l}}_{i,c} - \tau \right)}{1 - \tau},$$

where $\hat{\mathbf{l}}_{i,c}$ denotes the LDR image with pixel index i and color channel c . We set the threshold $\tau = 0.97$ in all our experiments, resulting in a short ramp towards saturation.

4. Experiments

We have conducted the following experiments to demonstrate the effectiveness of our approach. Sec. 4.1 shows the

generated HDR data and compares it with an LDR equivalent; Sec. 4.2 explores the influences of the exposure distribution, model backbone, and sampled camera response curve; and Sec. 4.3 presents the results of our unsupervised ITM. We use the tone mapper of Mantiuk et al. [44] to display HDR content in this paper. We refer readers to the supplemental for more results and full HDR data.

Implementation details. We implement our method on top of the official StyleGAN-XL [61] implementation under the PyTorch [54] environment. Training takes six days with four RTX 8000 GPUs for 256×256 resolution dataset – roughly the same time as training a vanilla StyleGAN-XL model.

Datasets. We collect five different datasets which naturally contain scenes with high dynamic range from the Internet: Landscapes (~ 7700 images), Lightning (~ 7000 images), Windows (~ 4200 images), Fireplaces (~ 2600 images), and Fireworks (~ 5600 images). We randomly crop and resize each image to the target resolution. Please refer to Fig. 1 and Fig. 4 for examples of generated images from models trained on these datasets. We collect our datasets from several websites: Flickr, Pexels, Instagram, and 500PX. We will make the datasets available upon request. Further, we will make all source code and pre-trained models publicly available upon publication.

4.1. Generation of HDR Images

We show in Fig. 4 a variety of samples generated with GlowGAN and a comparison with a vanilla StyleGAN-XL. Generated samples from the vanilla GAN often bear overexposed regions similar to those in the LDR training images. In contrast, samples from GlowGAN preserve detailed appearance information even for bright objects such as the sun, as they have more extensive dynamic ranges than those from the vanilla GAN. To show the difference in dynamic range, we plot the image brightness histogram of the two models in Fig. 5, where each histogram is computed from 500 randomly sampled images. It can be seen that the histogram of the vanilla GAN is cut off at 1, while GlowGAN clearly

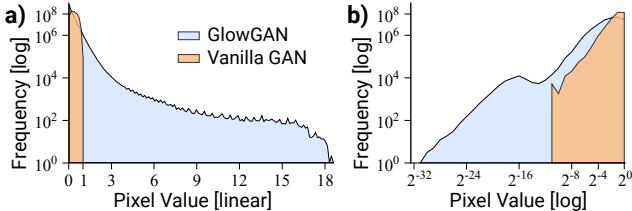


Figure 5. Pixel values from our model and a vanilla GAN. The linear-scale histogram (a) shows that we significantly extend the dynamic range for high pixel values, while the log-scale histogram (b, truncated to only show values up to 1) demonstrates an extension for low pixel values as well.

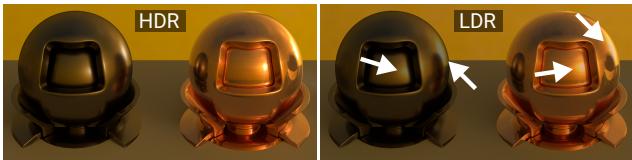


Figure 6. Image-based lighting using a generated HDR environment map (left) vs. a plain LDR equivalent (right). Arrows mark regions where the differences are most pronounced. The HDR illumination leads to high-contrast and high-frequency reflections in the specular materials (right object) and produces a more realistic bloom effect in the glossy material (left object).

avoids pixel intensity clamping (Fig. 5a) and has a much wider histogram, also in the dark regions (Fig. 5b). These results satisfy our expectation of learning HDR information from LDR data. Moreover, we see that GlowGAN can synthesize diverse images that do not exist in the real world. It thus opens up a new avenue for getting cheap abundant HDR data which can be used in several applications, e.g., for creating environment maps for image-based lighting (IBL) [14], as showcased in Fig. 6. Additionally, GlowGAN can interpolate between two environment maps, achieving a smooth transition effect, as demonstrated in the supplemental.

4.2. Ablation Study

Exposure Distribution. Our method assumes the exposure e follows a Gaussian distribution with variance σ_e^2 . Here we study how σ_e^2 impacts the generated image quality and the dynamic range. We employ the commonly used FID [24] and KID [8] scores to evaluate image quality on the Landscapes dataset. As we do not have ground truth HDR images, the scores are computed between the generated LDR images (i.e., output of the camera model) and the LDR training images. We also compute the dynamic range (DR) for each generated HDR image r as $DR = \log_2(r_{\max}/r_{\min})$, where r_{\max} and r_{\min} are the max and min values of the image, respectively. We report the median and the 90th percentile of DR computed over 50k images, referred to as DR50 and DR90, respectively. From Table 1, we can observe a trade-off between image quality and the dynamic range, i.e., increasing

Table 1. Effects of model and σ_e^2 on quality and dynamic range.

Model	σ_e^2	FID \downarrow	KID($\times 10^4$) \downarrow	DR50	DR90
SG-XL ¹	–	3.48	2.69	8.0	8.0
Ours	1.0	3.61	3.07	15.4	20.2
Ours	3.0	3.87	4.04	16.2	20.7
Ours	5.0	4.00	4.78	16.5	20.8
SG2 ²	–	9.2	23.37	8.0	8.0
Ours w/ SG2 ³	1.0	10.02	28.41	16.3	22.9

¹ Refers to a vanilla StyleGAN-XL model.

² Refers to a vanilla StyleGAN2-ADA model.

³ Refers to our approach with a StyleGAN2-ADA backbone.

Table 2. Comparing quality with fixed and stochastic CRFs.

Dataset	CRF	FID \downarrow	KID($\times 10^4$) \downarrow
Landscapes	Fixed	3.89	3.80
	Stochastic	3.61	3.07
Lightning	Fixed	3.40	4.77
	Stochastic	3.29	4.57

σ_e^2 leads to a higher dynamic range (with diminishing returns for high σ_e^2) but slightly worse FID and KID scores. To understand the positive correlation between σ_e^2 and DR, suppose that r has a low dynamic range, then with a very small or large exposure e (which is more likely to happen for large σ_e^2), it would produce an out-of-distribution LDR image 1 that is overly dark or bright. In other words, only a valid high dynamic range r can yield realistic LDR images when processed with different exposures. On the other hand, as σ_e^2 increases, the camera model interferes more with the image generation process, which may increase the training difficulty as the Gaussian distribution is only an approximation to the underlying exposure distribution. In practice, users can choose a suitable σ_e^2 depending on their goal. In most of our experiments, we use $\sigma_e^2 = 1$ as it already removes overexposure while featuring good quality and it also obtains better scores in the inverse tone mapping application. We provide more results on the effect of σ_e^2 in the supplemental. From Table 1 we can further see that a vanilla StyleGAN-XL exhibits slightly higher image quality, but a substantially smaller (log-scale) dynamic range.

Generator backbone. We further test our method based on StyleGAN2-ADA [31]. As Table 1 shows, our method also successfully synthesizes images with high dynamic range, albeit the final image quality directly depends on the baseline generative model.

Stochastic CRF. We further study the effects of the stochastic CRF sampling process in our model. Table 2 compares our stochastic CRF sampling with a fixed CRF, using $\beta = 0.6$ and $\gamma = 0.9$. Modelling the stochasticity leads to a clear improvement in FID and KID scores. This is because the in-the-wild LDR images used for training are captured with different cameras with diverse CRFs, which can be better modelled via stochastic CRF sampling.

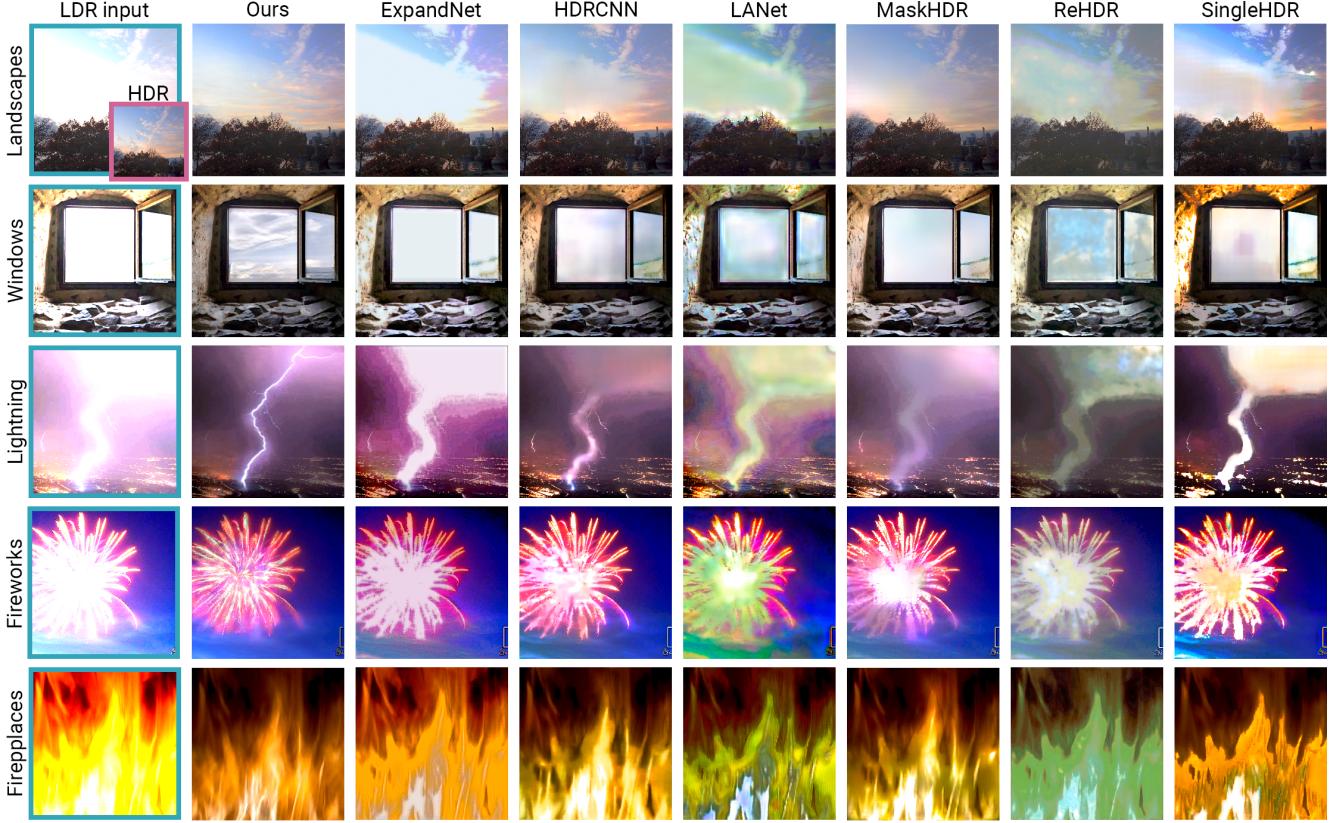


Figure 7. Results for the Inverse Tone Mapping application and comparisons to six state-of-the-art methods for different datasets. Given a single LDR image as input, our method can produce an HDR image that contains plausible and realistic content in the previously overexposed regions, while previous methods tend to produce blurred results or even noticeable artifacts in such regions. In the first row (Landscapes dataset), it can be seen that recovering the ground truth content present in the original HDR scene is not possible in fully saturated regions, however, our method is able to produce plausible results that are consistent with the scene.

Table 3. Evaluation of the Inverse Tone Mapping application. We achieve on-par quality with the best-performing supervised methods in reference metrics while outperforming all previous approaches in the non-reference metric that evaluates the overall naturalness and quality of the image. Note that reference metrics are not best suited for this evaluation, since the strength of our method lies in the reconstruction of missing overexposed regions, and therefore it is expected that the hallucinated content does not match that of the original HDR scene.

Method	Unsupervised	Reference			Non-Reference
		HDR-VDP3 \uparrow	PU21-VSI \uparrow	PU21-PSNR \uparrow	PU21-PIQE \downarrow
HDRCNN [16]	✗	7.42 \pm 1.02	0.961 \pm 0.034	32.1 \pm 5.1	36.1 \pm 5.3
MaskHDR [60]	✗	7.60 \pm 0.93	0.962 \pm 0.032	32.4 \pm 5.1	33.3 \pm 6.5
SingleHDR [43]	✗	7.01 \pm 1.17	0.956 \pm 0.031	30.1 \pm 4.5	40.3 \pm 6.2
ExpandNet [45]	✗	6.66 \pm 1.61	0.957 \pm 0.033	30.7 \pm 4.2	43.2 \pm 6.6
ReHDR [40]	✗	7.06 \pm 1.31	0.953 \pm 0.035	30.3 \pm 4.2	39.7 \pm 4.7
LANet [69]	✗	6.94 \pm 0.98	0.956 \pm 0.031	29.0 \pm 3.6	40.6 \pm 6.5
Ours	✓	7.44 \pm 0.94	0.961 \pm 0.032	31.8 \pm 4.4	31.8 \pm 5.1

4.3. Application: Inverse Tone Mapping

A potential application of our approach is unsupervised inverse tone mapping (ITM). We show both quantitatively and through a user study that our method outperforms previous approaches in hallucinating content in large overexposed regions, effectively recovering HDR content from a single LDR image.

Objective comparisons. We compare our unsupervised approach to six state-of-the-art fully-supervised ITM methods, which we abbreviate for simplicity as HDRCNN [16], MaskHDR [60], SingleHDR [43], ExpandNet [45], ReHDR [39], and LANet [69]. Following the work of Hanji et al. on quality assessment of single image HDR reconstruction methods [20], we select their three recommended



Figure 8. Our method can generate different but plausible HDR images from a single LDR input with large overexposed regions.

full-reference metrics (PU21-PSNR, PU21-VSI, and HDR-VDP3) as well as their recommended non-reference metric (PU21-PIQE). As test set for the reference metrics we use as HDR ground truth a set of 62 images collected from existing datasets [20, 53] and generate the corresponding LDR input images following the pipeline proposed by Eilertsen et al. [16]. For the non-reference metric, we use an extended set of 100 LDR images obtained from the Internet which we use directly as input with 15% to 45% of the pixels saturated. For fairness in these comparisons, both sets are composed of landscape images, since fully-supervised approaches are typically trained with datasets mainly containing this type of content. We show in Table 3 the results of these metrics and in Fig. 7 visual comparisons for our five datasets. Note that, since our method focuses on hallucinating content in completely saturated regions, it is highly unlikely that this content fully matches that of the original ground truth image, therefore full-reference metrics are not well suited for assessing the quality of our reconstructions. Nevertheless, our unsupervised method is still on par with previous fully-supervised approaches for the reference metrics, while it excels in the non-reference metric, showing that our hallucinated content is more plausible in terms of naturalness. Additionally, previous methods can generate only one potential reconstruction given an input LDR image, while our approach allows the generation of multiple results with different but plausible semantic information, as Fig. 8 shows.

User Study. Since one of the main strengths of our work is the capability to hallucinate plausible content in overexposed regions, reference metrics are unsuitable for comparisons. Additionally, the image diversity of available ground-truth HDR datasets is limited. Therefore, we perform a subjective study in order to further assess the quality of our generated results for the inverse tone mapping application. We include 20 scenes (four for each of our five datasets). For each scene, HDR results obtained with each of the seven methods were shown on a single screen, and participants were asked to rank the seven images from 1 (most preferred) to 7 (least preferred). The presentation order of the scenes and methods was randomized. The images were displayed in an HDR display Dell UP3221Q (3840×2160 resolution) in a standard

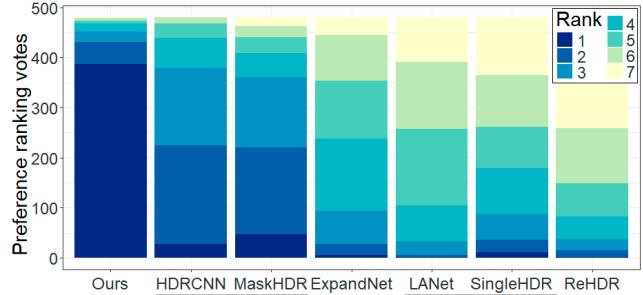


Figure 9. Preference rankings for the seven inverse tone mapping methods aggregated across participants and scenes. Different colors indicate the rankings (from 1 to 7). Methods marked in the same set (gray underline) are statistically indistinguishable, while all others present statistically significant differences in their distributions.

office room with natural illumination, and participants sat at a distance of 0.5 meters from the display. A total of 24 participants (38% female, aged 22 to 37 years old with normal or corrected-to-normal vision) participated in the study. We show in Fig. 9 the preference rankings for each method, aggregated for all scenes and participants. To analyze the results, we use pairwise Kruskal-Wallis tests adjusted by Bonferroni correction for multiple comparisons since the rankings do not follow a normal distribution. Results reveal that our method was ranked significantly higher than all others ($p < 0.001$), and it was selected as the top performing method in over 80% of the trials.

5. Conclusion & Discussion

We have introduced GlowGAN, a novel paradigm for learning HDR imagery from LDR data. Our method is orthogonal to other advances in generative adversarial learning and can be easily incorporated into any GAN-based pipeline. A trained GlowGAN acts as a strong prior, producing starkly more plausible inverse tone mapping results than previous approaches. Our inverse tone mapping method builds on GAN inversion via optimization, which can sometimes result in low-quality images, especially for high-frequency content (Fig. 10) – a problem that is orthogonal to our approach.



Figure 10. Failure case.

We heavily rely on training data with a diverse exposure distribution. While this assumption is oftentimes naturally satisfied for in-the-wild photo datasets, the dynamic range we can obtain is tightly linked to the exposure variance in the dataset. We hope that our approach inspires future work on learning rich models from casually captured images.

References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2StyleGAN: How to embed images into the StyleGAN latent space? In *ICCV*, pages 4432–4441, 2019. 4, 5
- [2] Ahmet Oğuz Akyüz, Roland Fleming, Bernhard E Riecke, Erik Reinhard, and Heinrich H Bülthoff. Do HDR displays support LDR content? A psychophysical evaluation. *ACM TOG*, 26(3):38–es, 2007. 3
- [3] Francesco Banterle, Alessandro Artusi, Kurt Debattista, and Alan Chalmers. *Advanced high dynamic range imaging*. AK Peters/CRC Press, 2017. 3
- [4] Francesco Banterle, Patrick Ledda, Kurt Debattista, Marina Bloj, Alessandro Artusi, and Alan Chalmers. A psychophysical evaluation of inverse tone mapping techniques. In *Computer Graphics Forum*, volume 28, pages 13–25. Wiley Online Library, 2009. 3
- [5] Francesco Banterle, Patrick Ledda, Kurt Debattista, and Alan Chalmers. Inverse tone mapping. In *Proceedings of the 4th International Conference on Computer Graphics and Interactive Techniques in Australasia and Southeast Asia*, pages 349–356, 2006. 2, 3
- [6] Francesco Banterle, Patrick Ledda, Kurt Debattista, and Alan Chalmers. Expanding low dynamic range videos for high dynamic range applications. In *Proceedings of the 24th Spring Conference on Computer Graphics*, pages 33–41, 2008. 3
- [7] Francesco Banterle, Demetris Marnerides, Kurt Debattista, and Thomas Bashford-Rogers. Unsupervised HDR imaging: What can be learned from a single 8-bit video? *arXiv preprint arXiv:2202.05522*, 2022. 3
- [8] Mikołaj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying MMD GANs. *arXiv preprint arXiv:1801.01401*, 2018. 6
- [9] Yochai Blau and Tomer Michaeli. The perception-distortion tradeoff. In *CVPR*, pages 6228–6237, 2018. 5
- [10] Ashish Bora, Eric Price, and Alexandros G Dimakis. AmbientGAN: Generative models from lossy measurements. In *ICLR*, 2018. 3, 4
- [11] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-GAN: Periodic implicit generative adversarial networks for 3D-aware image synthesis. In *CVPR*, pages 5799–5809, 2021. 3
- [12] Zhaoxi Chen, Guangcong Wang, and Ziwei Liu. Text2light: Zero-shot text-driven HDR panorama generation. *ACM TOG*, 41(6):1–16, 2022. 3
- [13] Antonia Creswell and Anil Anthony Bharath. Inverting the generator of a generative adversarial network. *IEEE Transactions on Neural Networks and Learning Systems*, 30(7):1967–1974, 2018. 4
- [14] Paul Debevec. Image-based lighting. *IEEE Computer Graphics and Applications*, 22(02):26–34, 2002. 2, 6
- [15] Paul E. Debevec and Jitendra Malik. Recovering high dynamic range radiance maps from photographs. In *Proc. ACM SIGGRAPH*, page 369–378, 1997. 1, 3
- [16] Gabriel Eilertsen, Joel Kronander, Gyorgy Denes, Rafał K Mantiuk, and Jonas Unger. HDR image reconstruction from a single exposure using deep cnns. *ACM TOG*, 36(6):1–15, 2017. 2, 3, 4, 5, 7, 8
- [17] Yuki Endo, Yoshihiro Kanamori, and Jun Mitani. Deep reverse tone mapping. *ACM TOG*, 36(6):177–1, 2017. 2, 3
- [18] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, pages 2672–2680, 2014. 3
- [19] Michael D Grossberg and Shree K Nayar. What is the space of camera response functions? In *CVPR*, volume 2, pages II–602, 2003. 2, 4
- [20] Param Hanji, Rafal Mantiuk, Gabriel Eilertsen, Saghi Hajisharif, and Jonas Unger. Comparison of single image HDR reconstruction methods—the caveats of quality assessment. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–8, 2022. 7, 8
- [21] Gang He, Shaoyi Long, Li Xu, Chang Wu, Jinjia Zhou, Ming Sun, Xing Wen, and Yurong Dai. Global priors guided modulation network for joint super-resolution and inverse tone-mapping. *arXiv preprint arXiv:2208.06885*, 2022. 3
- [22] Gang He, Kepeng Xu, Li Xu, Chang Wu, Ming Sun, Xing Wen, and Yu-Wing Tai. SDRTV-to-HDRTV via hierarchical dynamic context feature mapping. *arXiv preprint arXiv:2207.00319*, 2022. 3
- [23] Philipp Henzler, Niloy J Mitra, and Tobias Ritschel. Escaping Plato’s cave: 3D shape from adversarial rendering. In *ICCV*, pages 9984–9993, 2019. 3
- [24] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. *NeurIPS*, 30, 2017. 6
- [25] Xin Huang, Qi Zhang, Ying Feng, Hongdong Li, Xuan Wang, and Qing Wang. HDR-NeRF: High dynamic range neural radiance fields. In *CVPR*, pages 18398–18408, 2022. 2
- [26] So Yeon Jo, Siyeong Lee, Namhyun Ahn, and Suk-Ju Kang. Deep arbitrary HDRI: Inverse tone mapping with controllable exposure changes. *IEEE Transactions on Multimedia*, 2021. 2, 3
- [27] Kim Jun-Seong, Kim Yu-Ji, Moon Ye-Bin, and Tae-Hyun Oh. HDR-Plenoxels: Self-calibrating high dynamic range radiance fields. *arXiv preprint arXiv:2208.06787*, 2022. 2
- [28] Nima Khademi Kalantari and Ravi Ramamoorthi. Deep HDR video from sequences with alternating exposures. In *Computer Graphics Forum*, volume 38, pages 193–205. Wiley Online Library, 2019. 3
- [29] Takuhiro Kaneko. Unsupervised learning of depth and depth-of-field effect from natural images with aperture rendering generative adversarial networks. In *CVPR*, pages 15679–15688, 2021. 3
- [30] Takuhiro Kaneko and Tatsuya Harada. Noise robust generative adversarial networks. In *CVPR*, pages 8404–8414, 2020. 3
- [31] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. *NeurIPS*, 33:12104–12114, 2020. 3, 6
- [32] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *NeurIPS*, 34:852–863, 2021. 3

- [33] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, pages 4401–4410, 2019. 3
- [34] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *CVPR*, pages 8110–8119, 2020. 3
- [35] Soo Ye Kim, Jihyong Oh, and Munchurl Kim. Deep SR-ITM: Joint learning of super-resolution and inverse tone-mapping for 4k UHD HDR applications. In *ICCV*, pages 3116–3125, 2019. 3
- [36] Soo Ye Kim, Jihyong Oh, and Munchurl Kim. JSI-GAN: GAN-based joint super-resolution and inverse tone-mapping with pixel-wise task-specific filters for UHD HDR video. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 11287–11295, 2020. 3
- [37] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 4
- [38] Hayden Landis. Production-ready global illumination. *Siggraph course notes*, 16(2002):11, 2002. 3
- [39] Siyeong Lee, Gwon Hwan An, and Suk-Ju Kang. Deep chain HDRI: Reconstructing a high dynamic range image from a single low dynamic range image. *IEEE Access*, 6:49913–49924, 2018. 2, 3, 7
- [40] Siyeong Lee, Gwon Hwan An, and Suk-Ju Kang. Deep recursive HDRI: Inverse tone mapping using generative adversarial networks. In *ECCV*, pages 596–611, 2018. 3, 7
- [41] Siyeong Lee, So Yeon Jo, Gwon Hwan An, and Suk-Ju Kang. Learning to generate multi-exposure stacks with cycle consistency for high dynamic range imaging. *IEEE Transactions on Multimedia*, 23:2561–2574, 2020. 2, 3
- [42] Yiyi Liao, Katja Schwarz, Lars Mescheder, and Andreas Geiger. Towards unsupervised learning of generative models for 3D controllable image synthesis. In *CVPR*, pages 5871–5880, 2020. 3
- [43] Yu-Lun Liu, Wei-Sheng Lai, Yu-Sheng Chen, Yi-Lung Kao, Ming-Hsuan Yang, Yung-Yu Chuang, and Jia-Bin Huang. Single-image HDR reconstruction by learning to reverse the camera pipeline. In *CVPR*, pages 1651–1660, 2020. 2, 3, 7
- [44] Rafal Mantiuk, Scott Daly, and Louis Kerofsky. Display adaptive tone mapping. In *ACM SIGGRAPH 2008 papers*, pages 1–10. 2008. 5
- [45] Demetris Marnerides, Thomas Bashford-Rogers, Jonathan Hatchett, and Kurt Debattista. Expandnet: A deep convolutional neural network for high dynamic range expansion from low dynamic range content. In *Computer Graphics Forum*, volume 37, pages 37–49. Wiley, 2018. 2, 3, 7
- [46] Belen Masia, Sandra Agustin, Roland W Fleming, Olga Sorkine, and Diego Gutierrez. Evaluation of reverse tone mapping through varying exposure conditions. In *ACM SIGGRAPH Asia 2009*, pages 1–8. 2009. 3
- [47] Belen Masia, Ana Serrano, and Diego Gutierrez. Dynamic range expansion based on image statistics. *Multimedia Tools and Applications*, 76(1):631–648, 2017. 3
- [48] Laurence Meylan, Scott Daly, and Sabine Süsstrunk. The reproduction of specular highlights on high dynamic range displays. In *Color and Imaging Conference*, volume 2006, pages 333–338, 2006. 3
- [49] Ben Mildenhall, Peter Hedman, Ricardo Martin-Brualla, Pratul P Srinivasan, and Jonathan T Barron. NeRF in the dark: High dynamic range view synthesis from noisy raw images. In *CVPR*, pages 16190–16199, 2022. 2
- [50] Tomoo Mitsunaga and Shree K. Nayar. Radiometric self calibration. *Proceedings. 1999 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (Cat. No PR00149)*, 1:374–380 Vol. 1, 1999. 1, 3
- [51] Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yong-Liang Yang. HoloGAN: Unsupervised learning of 3d representations from natural images. In *ICCV*, pages 7588–7597, 2019. 3
- [52] Yuzhen Niu, Jianbin Wu, Wenxi Liu, Wenzhong Guo, and Rynson WH Lau. HDR-GAN: HDR image reconstruction from multi-exposed LDR images with large motions. *IEEE Transactions on Image Processing*, 30:3885–3896, 2021. 3
- [53] Karen Panetta, Landry Kezebou, Victor Oludare, Sos Agaian, and Zehua Xia. Tmo-net: A parameter-free tone mapping operator using generative adversarial network, and performance benchmarking on large scale HDR dataset. *IEEE Access*, 9:39500–39517, 2021. 8
- [54] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *NeurIPS*, 32, 2019. 5
- [55] Erik Reinhard, Wolfgang Heidrich, Paul Debevec, Sumanta Pattanaik, Greg Ward, and Karol Myszkowski. *High dynamic range imaging: Acquisition, display, and image-based lighting*. Morgan Kaufmann, 2010. 1, 2
- [56] Allan G Rempel, Matthew Trentacoste, Helge Seetzen, H David Young, Wolfgang Heidrich, Lorne Whitehead, and Greg Ward. LDR2HDR: On-the-fly reverse tone mapping of legacy video and photographs. *ACM TOG*, 26(3):39–es, 2007. 2, 3
- [57] Mark A Robertson, Sean Borman, and Robert L Stevenson. Estimation-theoretic approach to dynamic range enhancement using multiple exposures. *Journal of electronic imaging*, 12(2):219–228, 2003. 1, 3
- [58] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based editing of real images. *arXiv preprint arXiv:2106.05744*, 2021. 4
- [59] Darius Rückert, Linus Franke, and Marc Stamminger. ADOP: Approximate differentiable one-pixel point rendering. *ACM TOG*, 41(4):1–14, 2022. 2
- [60] Marcel Santana Santos, Tsang Ing Ren, and Nima Khademi Kalantari. Single image HDR reconstruction using a CNN with masked features and perceptual loss. *arXiv preprint arXiv:2005.07335*, 2020. 2, 3, 5, 7
- [61] Axel Sauer, Katja Schwarz, and Andreas Geiger. StyleGAN-XL: Scaling StyleGAN to large diverse datasets. In *ACM SIGGRAPH*, pages 1–10, 2022. 3, 5
- [62] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. GRAF: Generative radiance fields for 3D-aware image synthesis. *NeurIPS*, 33:20154–20166, 2020. 3

- [63] Helge Seetzen, Wolfgang Heidrich, Wolfgang Stuerzlinger, Greg Ward, Lorne Whitehead, Matthew Trentacoste, Abhijeet Ghosh, and Andrejs Vorozcovs. High dynamic range display systems. In *ACM SIGGRAPH*, pages 760–768. 2004. 2
- [64] Attila Szabó, Givi Meishvili, and Paolo Favaro. Unsupervised generative 3D shape learning from natural images. *arXiv preprint arXiv:1910.00287*, 2019. 3
- [65] Chao Wang, Bin Chen, Hans-Peter Seidel, Karol Myszkowski, and Ana Serrano. Learning a self-supervised tone mapping operator via feature contrast masking loss. In *Computer Graphics Forum*, volume 41, pages 71–84. Wiley Online Library, 2022. 2
- [66] Lvdi Wang, Li-Yi Wei, Kun Zhou, Baining Guo, and Heung-Yeung Shum. High dynamic range image hallucination. *Rendering Techniques*, 321:326, 2007. 3
- [67] Lin Wang and Kuk-Jin Yoon. Deep learning for HDR imaging: State-of-the-art and future trends. *IEEE PAMI*, 2021. 2
- [68] Weihao Xia, Yulun Zhang, Yujiu Yang, Jing-Hao Xue, Bolei Zhou, and Ming-Hsuan Yang. GAN inversion: A survey. *IEEE PAMI*, 2022. 4
- [69] Hanning Yu, Wentao Liu, Chengjiang Long, Bo Dong, Qin Zou, and Chunxia Xiao. Luminance attentive networks for HDR image and panorama reconstruction. In *Computer Graphics Forum*, volume 40, pages 181–192. Wiley Online Library, 2021. 2, 3, 7
- [70] 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. 4
- [71] Yang Zhang and TO Aydin. Deep HDR estimation with generative detail reconstruction. In *Computer Graphics Forum*, volume 40, pages 179–190. Wiley Online Library, 2021. 2, 3
- [72] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A Efros. Generative visual manipulation on the natural image manifold. In *ECCV*, pages 597–613. Springer, 2016. 4