

SINGLE IMAGE LDR TO HDR CONVERSION USING CONDITIONAL DIFFUSION

Dwip Dalal*

Gautam Vashishtha*

Prajwal Singh[†]

Shanmuganathan Raman[‡]

Computer Vision, Imaging and Graphics Lab
Indian Institute of Technology Gandhinagar, India
{dwip.dalal, gautam.pv, singh_prajwal, shanmuga}@iitgn.ac.in

ABSTRACT

Digital imaging aims to replicate realistic scenes, but Low Dynamic Range (LDR) cameras cannot represent the wide dynamic range of real scenes, resulting in under-/overexposed images. This paper presents a deep learning-based approach for recovering intricate details from shadows and highlights while reconstructing High Dynamic Range (HDR) images. We formulate the problem as an image-to-image (I2I) translation task and propose a conditional Denoising Diffusion Probabilistic Model (DDPM) based framework using classifier-free guidance. We incorporate a deep CNN-based autoencoder in our proposed framework to enhance the quality of the latent representation of the input LDR image used for conditioning. Moreover, we introduce a new loss function for LDR-HDR translation tasks, termed Exposure Loss. This loss helps direct gradients in the opposite direction of the saturation, further improving the results' quality. By conducting comprehensive quantitative and qualitative experiments, we have effectively demonstrated the proficiency of our proposed method. The results indicate that a simple conditional diffusion-based method can replace the complex camera pipeline-based architectures.

Index Terms— Diffusion Model, Autoencoder, High Dynamic Range Imaging, Computational Photography

1. INTRODUCTION

High dynamic range (HDR) imaging is a promising technique for improving the viewing experience of digital content by capturing real-world lighting and details. However, low dynamic range (LDR) cameras are unable to capture the vast dynamic range in real-world scenes. A workaround for this is to capture several LDR images taken at various exposures and concatenate them to get HDR images, but it frequently results in ghosting artifacts, especially in scenes that are dynamic. To overcome these limitations, deep convolutional neural networks (CNNs) have been used to develop single-image HDR reconstruction techniques [1], [2], [3]. These techniques address the problems with LDR-to-HDR mapping, which is difficult because HDR pixels (32-bit floating point) have much more variation than LDR pixels (8-bit unsigned integers).

Research on LDR to HDR translation has recently received intense attention. Endo et al. [4] proposed the deep-learning-based approach for fully automatic inference using deep convolutional neural networks. They adopt a bracketed approach by inferring from a sequence of k LDR images of different exposures and then reconstructing an HDR image by merging the LDR images [5]. The work

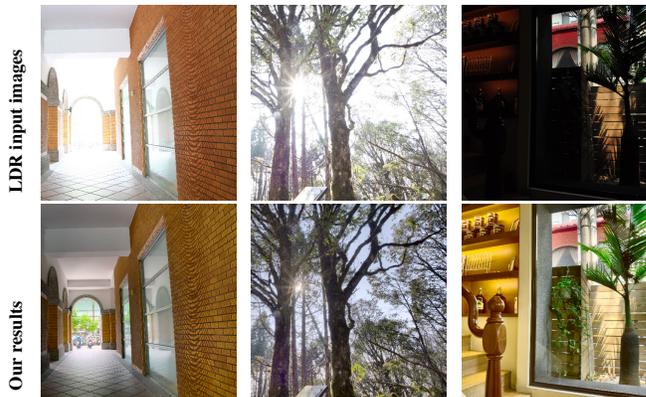


Table 1: HDR reconstruction from a single LDR image. We condition the DDPM and guide the gradient to recover missing details from shadows and highlights.

[3] used a full CNN design in the form of a hybrid dynamic range autoencoder that transformed the LDR input image using the encoder network to generate a compact feature representation and generated the output HDR image by passing it to the HDR decoder operating in the log domain.

Liu et al. [1] approached to tackle this problem by using the domain knowledge of LDR image formation pipeline to decompose the reconstruction into three sub-tasks of i) dequantization, ii) linearization, and iii) hallucination. They developed networks for each of these tasks using CNNs at the core of each architecture. The work in [2] approaches the problem in a similar fashion but utilizes a condition GAN-based framework.

In this paper, we introduce a novel method for reconstructing high-quality HDR images from a single LDR image without the requirement of an explicit inverse camera pipeline [1] [2]. Our approach relies upon the utilization of a conditional classifier-free diffusion architecture [6]. Besides the fundamental structure of the diffusion architecture, the model includes an encoder network that generates a latent representation of the LDR input image, which is used to condition the output for the generation of the subsequent HDR images. We propose a novel methodology of employing a CNN-based decoder network to enhance the obtained latent representation, yielding superior conditioning over input LDR images. The proposed loss function employed in this work integrates several terms, including reconstruction loss for the autoencoder architecture, multiscale training loss [7], and perceptual loss [8]. Additionally, we introduce a novel loss function named "Exposure Loss," which helps in achieving optimal balance in the exposure of the reconstructed images. By prioritizing this metric, our approach significantly improves the quality of the resulting output images. Incorporating these

*Equal contribution and [†]This work is supported by Prime Minister Research Fellowship (PMRF-2122-2557) and [‡] Jibaben Patel Chair.

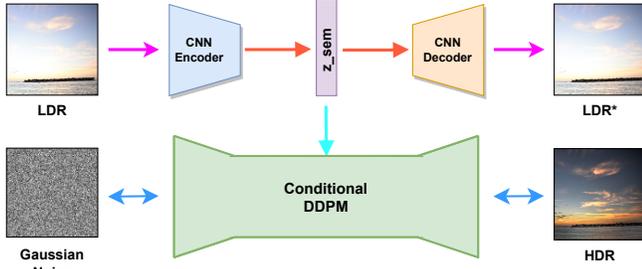


Fig. 1: The figure illustrates conditional diffusion architecture. In it, z_{sem} represents the latent representation of the input LDR image, and LDR^* is the output reconstruction of the input LDR image.

terms in the loss function helps achieve high-quality HDR image reconstruction from a single LDR image. The effectiveness of the proposed approach is evaluated through a series of experiments, which demonstrate its superiority in terms of both reconstruction quality and convergence rate compared to existing state-of-the-art methods.

2. BACKGROUND

The family of generative models that includes score-based generative models and diffusion-based models (DPMs) has been shown effective in modeling the target distribution through a denoising process with varying noise levels. These models can accurately transform an arbitrary Gaussian noise map of the prior distribution, $\mathcal{N}(\mathbf{0}, \mathbf{I})$, into a clear image sample after multiple denoising passes. To achieve this, Ho et al. [9] proposed a method for learning a function, $\epsilon_\theta(x_t, t)$, that predicts the corresponding noise of a noisy image, x_t , using a UNet architecture [10]. The model is trained using a loss function, $|\epsilon_\theta(x_t, t) - \epsilon|$, where ϵ is the noise added to x_0 producing x_t . The present formulation represents a simplified adaptation of the variational lower bound for the marginal log-likelihood, which has gained widespread adoption within the research community [11] [12] [13] [14].

2.1. Forward Diffusion

The forward diffusion adds Gaussian noise to a given HDR image in a series of T steps. The initial image is sampled from the training data distribution $q(x)$, and a variance schedule $\beta_t \in (0, 1)$ controls the noise step sizes. Specifically, at each step, the noisy version of the image \mathbf{x}_0 is generated by $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$ [15], resulting in a sequence of samples x_1, x_2, \dots, x_T . Due to Gaussian diffusion, we can produce the noisy version of the image \mathbf{x}_0 at any timestep t as

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\sqrt{\alpha_t} \mathbf{x}_0, (1 - \alpha_t) \mathbf{I}); \alpha_t = \prod_{s=1}^t (1 - \beta_s) \quad (1)$$

In most cases, α -cosine schedule is used among the various possible choices for selecting the variance scheduler. However, for higher resolutions images, Hoogeboom et al. [7] showed that the noise added by α -cosine schedule is not enough. Hence, here we use a modified noise scheduler introduced by [7]. The noise schedule is adjusted to hold the signal-to-noise ratio (SNR) constant at 64×64 resolution scale.

$$\log \text{SNR}_{\text{shift } 64}^{256 \times 256}(t) = -2 \log \tan(\pi t / 2) + 2 \log(64 / 256) \quad (2)$$

2.2. Reverse Diffusion

Upon successfully adding sequential noise to an input HDR image, our focus shifts to the inverse process, namely, the distribution $p(\mathbf{x}_{t-1} | \mathbf{x}_t)$. We utilize our AttenRecResUNet Fig. 2 to model this distribution. If the time gap between $t - 1$ and t is negligible (i.e., $T = \infty$), the probability function takes a Gaussian form, expressed as $\mathcal{N}(\mu_\theta(\mathbf{x}_t, t), \sigma_t)$ [9]. Various techniques can be employed to model this distribution, including the standard epsilon loss $\epsilon_\theta(x_t, t)$ [9]. It is vital to note that since $T = \infty$ is an unrealistic assumption in practice, DPMs can only offer approximations.

2.3. Conditional Diffusion

Similar to other generative frameworks, diffusion models can be made to sample conditionally given some variable of interest $p_\theta(x_0 | y)$ like a class label or a sentence description. Particularly in our case of generation of HDR images given a single LDR input, we want our output that is generated by starting from Gaussian noise to be conditioned on the input LDR image. [11] show that guiding the diffusion process using a separate classifier can help. In this setup, we take a pre-trained classifier to guide the reverse diffusion process. Specifically, we push process in the direction of the gradient of the target label probability. The downside of this approach is the reliance of another guiding network. An alternative approach proposed by [6] eliminates this reliance by using special training of the diffusion model itself to guide the sampling.

$$\hat{\epsilon}_\theta(x_t, t, y) = \epsilon_\theta(x_t, t, \phi) + s \cdot \epsilon_\theta(x_t, t, y) - \epsilon_\theta(x_t, t, \phi) \quad (3)$$

During training, the conditioning label, denoted by y , can be assigned a null label with a certain probability. At the inference stage, an artificial shift towards the conditional direction is applied to the reconstructed samples using a parameter termed the guidance scale (s) to distance them from the null label and thus enhance the effect of conditioning. This approach has been shown to yield superior sample quality based on human evaluation compared to classifier guidance, as reported in [16].

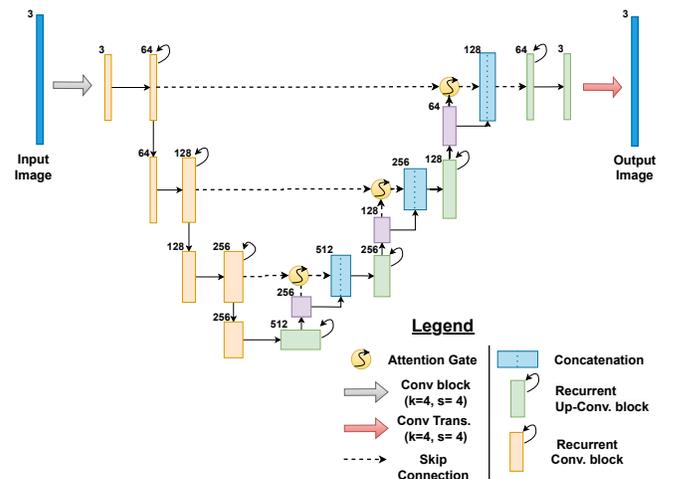


Fig. 2: Architecture of Attention Residual UNet, used for estimating Gaussian noise in image per time-step t . Downsampling is applied with max-pooling = 2, and upsampling with scale = 2

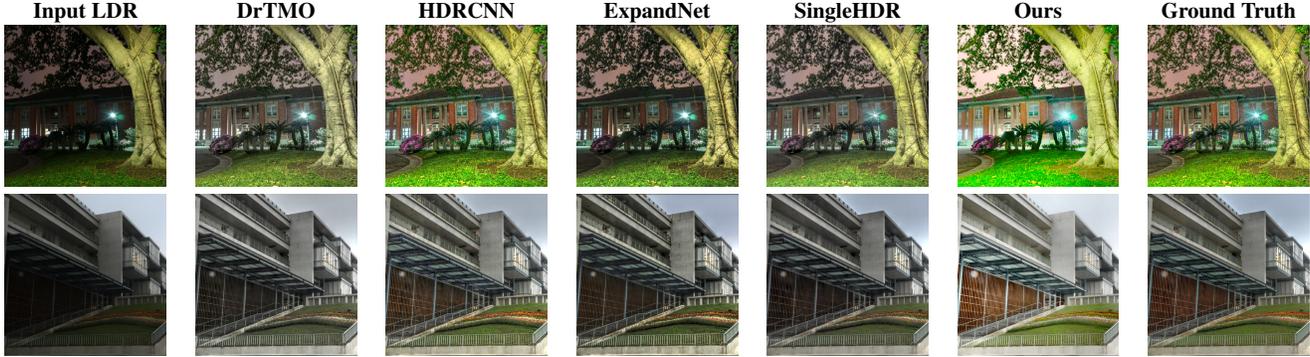


Table 2: Qualitative comparison of HDR image generated from our proposed approach with the previous state-of-the-art methods.

Method	HDR-EYE			HDR-REAL		
	PSNR	SSIM	VDP 2.2	PSNR	SSIM	VDP 2.2(1)
DrTMO [4]	9.28 ± 2.98	0.69 ± 0.15	48.33 ± 5.16	5.54 ± 3.55	0.36 ± 0.26	43.64 ± 6.51
HDRCNN [3]	16.12 ± 3.77	0.74 ± 0.12	50.75 ± 5.75	13.34 ± 7.68	0.53 ± 0.29	47.20 ± 7.77
ExpandNet [17]	17.12 ± 4.27	0.79 ± 0.13	51.09 ± 5.87	12.84 ± 6.90	0.50 ± 0.29	46.70 ± 7.90
SingleHDR [1]	15.47 ± 6.65	0.71 ± 0.19	53.15 ± 5.91	19.07 ± 7.16	0.64 ± 0.28	50.13 ± 7.74
Ours	16.97 ± 4.94	0.81 ± 0.14	52.29 ± 5.82	15.76 ± 6.68	0.66 ± 0.27	49.85 ± 7.27

Table 3: Quantitative evaluation of our approach. The best-performing method is highlighted **red** and second-best performing method is indicated by **blue** for each metric.

2.4. Diffusion model for higher resolution

The diffusion model struggles to converge when dealing with higher resolution images (256×256) [7]. To overcome this limitation, we have incorporated three methodologies, (1) modified noise schedules, (2) multiscale loss function and (3) architecture scaling, proposed in [7]. Our proposed model can handle sufficiently large under-/over exposed regions even with relatively fewer artifacts without any explicit inverse camera pipeline and is able to compete with the state-of-the-art models [1] [3] [18] [4] in this domain with a much more streamlined architecture.

3. METHODOLOGY

In this study, we have proposed a framework for a single image-based LDR-HDR reconstruction using a probabilistic diffusion model [9]. The framework consists of an autoencoder and conditional DDPM, as shown in Fig.2. The autoencoder model is used to generate encoding for the LDR image, which is later added to conditional DDPM. In the forward pass of the ddpm method, we first encode HDR images to isotropic gaussian noise, and the task of the backward pass is to reconstruct the HDR image using sampled gaussian noise by conditioning it on the encoded LDR image. We utilize four loss functions to ensure the convergence of the framework during training.

3.1. Multiscale training Loss

For high-resolution images, the unweighted loss on ϵ_t introduced by [9] fails to converge due to the domination of high-frequency details [7]. Hence, here we used multiscale training loss [7], which comprises the weighted sum of losses of multiple resolutions.

$$\tilde{L}_{\theta}^{256 \times 256}(\mathbf{x}) = \sum_{s \in \{32, 64, 128, 256\}} \frac{1}{s} L_{\theta}^{s \times s}(\mathbf{x}) \quad (4)$$

where $L_{\theta}^{s \times s}$ denotes,

$$L_{\theta}^{d \times d}(\mathbf{x}) = \frac{1}{d^2} \mathbb{E}_{\epsilon, t} \|D^{d \times d}[\epsilon] - D^{d \times d}[\hat{\epsilon}_{\theta}(\alpha_t \mathbf{x} + \sigma_t \epsilon, t)]\|_2^2 \quad (5)$$

and $D^{d \times d}$ downsamples to $d \times d$ resolution. Here d will take values belonging to $\{32, 64, 128, 256\}$

3.2. Reconstruction Loss

The decoder network is trained using reconstruction loss, computed on the image generated by the decoder using the compact latent space derived from the low dynamic range (LDR) image. This loss is formally defined as follows:

$$L_{rec} = \|\xi_{LDR} - \tilde{\xi}\|^2 \quad (6)$$

Here, ξ_{LDR} represents the ground truth LDR image, and $\tilde{\xi}$ denotes the decoded output of the autoencoder network. The reconstruction loss is designed to ensure that the autoencoder creates a more meaningful semantic representation of the input LDR image, which is then used to condition our AttenRecResUNet (Fig. 2) architecture in the reverse diffusion process.

3.3. Perceptual Loss

The Learned Perceptual Image Patch Similarity (LPIPS) [8] metric is utilized to assess the perceptual similarity of two images, and has been demonstrated to align with human perception. Here, the LPIPS score is employed to expedite model convergence by ensuring that the higher-level semantics of the predicted noisy image align perceptually with those of the original noisy image at each time step.

3.4. Exposure Loss

In most cases, we want the exposure of the output HDR image to be inverse with respect to the input LDR image, i.e., if an LDR image

is over-exposed, we would want to push the gradients of our model in the opposite direction, making sure that the output HDR image has balanced exposures. To be able to achieve this, we introduce a new loss defined as Exposure loss L_{exp} that helps in achieving the opposite gradient flow.

$$L_{exp} = -\zeta * \left(\left| \frac{\sum_p x_t}{\sum_p \max(x_t, 1)} - \frac{\sum_p x_{ldr}}{\sum_p \max(x_{ldr}, 1)} \right| \right) \quad (7)$$

Here \sum_p denotes pixel-wise summation, x_t denotes the normalized predicted HDR image, and x_{ldr} denotes the normalized input LDR image. The negation ensures that we minimize the penalty for the loss when there is a contrast between the exposures of the input LDR and the predicted HDR image. Scaling factor ζ was introduced to balance the impact(magnitude) of L_{exp} with the other terms, ensuring a fair contribution from each loss term during optimization. The final loss function that we use to train the proposed frames is as follows:

$$\mathcal{L}_T = L_{MSTL} + L_{rec} + L_{lips} + L_{exp} \quad (8)$$

4. EXPERIMENTS AND RESULTS

We evaluate our proposed method by conducting a comparison against several single-image-based HDR reconstruction approaches, including DrTMO [4], HDRCNN [3], ExpandNet [17], and Single-HDR [1]. We performed a comparison on two publicly available datasets HDR-Eye and HDR-REAL (test split). There are 1838 LDR-HDR image pairs in the HDR-REAL dataset test set and 46 LDR-HDR image pairs in HDR-Eye. The direct inference was conducted on publicly available pre-trained models of these approaches to obtain output HDR images.

The qualitative comparison is illustrated in Table 2, and for the quantitative evaluation, we used three different metrics: HDR-VDP 2.2 [19], PSNR, SSIM [20] metrics. Table 3 shows the scores obtained on HDR-Eye and HDR-Real Dataset (Test split) for the three metrics. Table 2 & 3 showcases that our model performs commendably when compared to the state-of-the-art models.

5. ABLATIONS

5.1. Autoencoder Diffusion

We introduce a novel deep CNN based decoder network that focuses on enhancing the quality of the latent representation formed from our encoder network that is used in the conditioning of our AttenRecResUNet (Fig. 2) during the reverse diffusion process. We validate the results of our Autoencoder Diffusion Model with the following variants:

- **Autoencoder Diffusion** - This is the complete autoencoder architecture used as a subpart of our proposed network.
- **Encoder Diffusion** - In this model, we remove the Decoder module from our autoencoder Diffusion architecture.

Table 4 & 5 qualitatively and quantitatively shows that adding the decoder block in our model causes the output HDR image to become significantly sharper compared to its variant.

5.2. Loss function

We also performed ablation studies on the importance of exposure loss. In Case 1 (Exposure P.), we included L_{exp} in the total loss computation, which was defined as total Loss = $L_{MSTL} + L_{rec} +$

$L_{lips} + L_{exp}$. In Case 2 (Exposure A.), we excluded L_{exp} and kept all other loss components intact. Table 4 & 5 qualitatively and quantitatively demonstrates that the addition of the Exposure Loss improves the exposure balance of the output HDR as compared to its variant.

Method	PSNR	SSIM	VDP 2.2
Vanilla Conditional Diffusion	16.36 ± 4.55	0.74 ± 0.11	50.93 ± 5.68
Autoencoder Diffusion	16.71 ± 4.69	0.79 ± 0.19	51.67 ± 5.79
Autoencoder Diffusion + Total Loss	16.97 ± 4.94	0.81 ± 0.14	52.29 ± 5.82

Table 4: Evaluation of our approach on HDR-Eye dataset

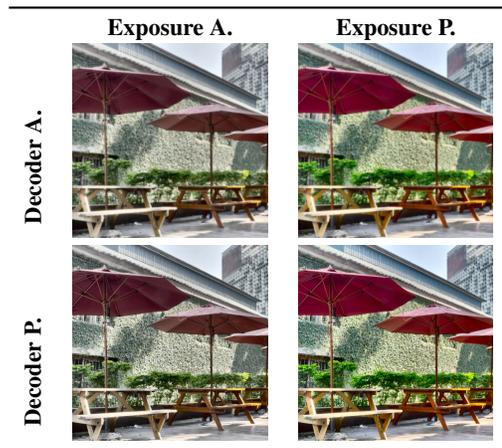


Table 5: Images plotted in the form of a confusion matrix to demonstrate the improvement of output on the addition of decoder block and Exposure loss. Here A. means absent, and P. means present.

6. LIMITATION AND FUTURE WORKS

With Exposure loss and decoder network in our architecture, we were able to substantially increase the speed of convergence as well as improve the quality of the generated HDR samples. However, in areas with lesser artifacts in the saturated regions, the quality of the output HDR images can be further improved by incorporating an inverse camera pipeline as done by [1] and [2]. By incorporating masking on under and over-exposed regions to treat them separately, the quality of the model in more stringent conditions can be improved.

7. CONCLUSION

The proposed framework reconstruct an HDR image from a single LDR image and is based on conditional diffusion with an autoencoder. We show that adding noise in HDR images during the forward pass of diffusion can help the network reconstruct lost information in LDR images during the backward pass. We also proposed a novel loss, named Exposure loss, that focused on explicitly directing the gradients in the direction opposite to the saturation, thus enhancing the quality of the results. Additionally, to ensure faster convergence and learn the semantics of the artifacts in the saturated regions quicker, we used perceptual loss. We show the effectiveness of the proposed framework over the previous state-of-the-art methods through multiple experiments and ablation studies.

8. REFERENCES

- [1] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1651–1660.
- [2] Prarabdh Raipurkar, Rohil Pal, and Shanmuganathan Raman, “Hdr-cgan: single ldr to hdr image translation using conditional gan,” in *Proceedings of the Twelfth Indian Conference on Computer Vision, Graphics and Image Processing*, 2021, pp. 1–9.
- [3] Gabriel Eilertsen, Joel Kronander, Gyorgy Denes, Rafał K Mantiuk, and Jonas Unger, “Hdr image reconstruction from a single exposure using deep cnns,” *ACM transactions on graphics (TOG)*, vol. 36, no. 6, pp. 1–15, 2017.
- [4] Yuki Endo, Yoshihiro Kanamori, and Jun Mitani, “Deep reverse tone mapping,” *ACM Trans. Graph.*, vol. 36, no. 6, pp. 177–1, 2017.
- [5] Paul E Debevec and Jitendra Malik, “Recovering high dynamic range radiance maps from photographs,” in *ACM SIGGRAPH 2008 classes*, pp. 1–10. 2008.
- [6] Jonathan Ho and Tim Salimans, “Classifier-free diffusion guidance,” *arXiv preprint arXiv:2207.12598*, 2022.
- [7] Emiel Hoogeboom, Jonathan Heek, and Tim Salimans, “simple diffusion: End-to-end diffusion for high resolution images,” *arXiv preprint arXiv:2301.11093*, 2023.
- [8] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang, “The unreasonable effectiveness of deep features as a perceptual metric,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 586–595.
- [9] Jonathan Ho, Ajay Jain, and Pieter Abbeel, “Denoising diffusion probabilistic models,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020.
- [10] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III 18*. Springer, 2015, pp. 234–241.
- [11] Prafulla Dhariwal and Alexander Nichol, “Diffusion models beat gans on image synthesis,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 8780–8794, 2021.
- [12] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al., “Imagen video: High definition video generation with diffusion models,” *arXiv preprint arXiv:2210.02303*, 2022.
- [13] Alexander Quinn Nichol and Prafulla Dhariwal, “Improved denoising diffusion probabilistic models,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 8162–8171.
- [14] Jiaming Song, Chenlin Meng, and Stefano Ermon, “Denoising diffusion implicit models,” *arXiv preprint arXiv:2010.02502*, 2020.
- [15] Konpat Preechakul, Nattanat Chatthee, Suttisak Wizadwongsa, and Supasorn Suwajanakorn, “Diffusion autoencoders: Toward a meaningful and decodable representation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10619–10629.
- [16] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen, “Glide: Towards photorealistic image generation and editing with text-guided diffusion models,” *arXiv preprint arXiv:2112.10741*, 2021.
- [17] 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*. Wiley Online Library, 2018, vol. 37, pp. 37–49.
- [18] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol, “Extracting and composing robust features with denoising autoencoders,” in *Proceedings of the 25th international conference on Machine learning*, 2008, pp. 1096–1103.
- [19] Rafał Mantiuk, Kil Joong Kim, Allan G Rempel, and Wolfgang Heidrich, “Hdr-vdp-2: A calibrated visual metric for visibility and quality predictions in all luminance conditions,” *ACM Transactions on graphics (TOG)*, vol. 30, no. 4, pp. 1–14, 2011.
- [20] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli, “Image quality assessment: from error visibility to structural similarity,” *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.