

SEMANTIC-PRESERVING IMAGE CODING BASED ON CONDITIONAL DIFFUSION MODELS

Francesco Pezone^{1,4}, Osman Musa³, Giuseppe Caire⁴, Sergio Barbarossa²

¹DIAG Department, ²DIET Department, Sapienza University of Rome

³BIFOLD, ⁴Communications and Information Theory Group, Technische Universität Berlin

E-mail: {francesco.pezone, sergio.barbarossa}@uniroma1.it, {osman.musa, caire}@tu-berlin.de

ABSTRACT

Semantic communication, rather than on a bit-by-bit recovery of the transmitted messages, focuses on the meaning and the goal of the communication itself. In this paper, we propose a novel semantic image coding scheme that preserves the semantic content of an image, while ensuring a good trade-off between coding rate and image quality. The proposed Semantic-Preserving Image Coding based on Conditional Diffusion Models (SPIC) transmitter encodes a Semantic Segmentation Map (SSM) and a low-resolution version of the image to be transmitted. The receiver then reconstructs a high-resolution image using a Denoising Diffusion Probabilistic Models (DDPM) doubly conditioned to the SSM and the low-resolution image. As shown by the numerical examples, compared to state-of-the-art (SOTA) approaches, the proposed SPIC exhibits a better balance between the conventional rate-distortion trade-off and the preservation of semantically-relevant features. Code available at <https://github.com/frapez1/SPIC>

Index Terms— Semantic communications, image segmentation, denoising diffusion probabilistic models, super-resolution diffusion models .

1. INTRODUCTION AND RELATED WORK

Semantic communications is lately receiving great attention because of its potential to improve the efficiency of communication systems, focusing directly on the semantics (meaning) of the transmitted messages rather than on recovering the bits used to represent the transmitted images [1, 2]. In semantic communication, there is no semantic error at the receiver if the reconstructed message is semantically equivalent to the transmitted one, even if the representations of the transmitted and recovered images do not coincide at the bit level. For example, in the transmission of an image captured by a web camera in an autonomous car, we might require that the reconstructed image should retain as accurately as possible the ability to detect semantically relevant objects, e.g. pedestrians, vehicles, traffic lights, etc., while providing contextually a sufficiently good trade-off between quality of the reconstructed image and the number of transmitted bits. This is just an example of combining semantics (identification of a class of meaningful objects) and the goal of communication (image reconstruction and the ability to segment the image properly at the receiver side). Several works have considered “semantic

coding” as joint source-channel coding with semantic side information [3]. However, since legacy protocols and network architectures are standardized according to the separation principle (OSI layers), and link layer control mechanisms do not pass erroneous data packets to the upper layers, we consider image compression at the “application layer” and do not consider transmission errors. The fusion of semantic segmentation with image reconstruction techniques has surfaced as a potent strategy to improve the quality of image reconstruction [4–6]. By attributing semantically-relevant labels to each pixel, SSMs represent a fundamental tool to encapsulate semantics within the image representation and can then play a key role in semantic communications.

Classical image compression techniques, such as JPEG, BPG or JPEG2000, target to achieve the best trade-off between compression ratio and image artifacts. However, this may come at the expense of semantic retention. Moreover, classical approaches can efficiently compress images without considering that some objects might be more relevant than others. The problem becomes even more relevant when some objects of interest have a small size, comparable to the patches used for compression. An example might be a pedestrian crossing the street in the distance. By applying a classical compression algorithm like JPEG2000, since the pedestrian size might occupy just a few 8×8 patches, a possible distortion in the reconstructed image might involve a small degradation of the overall image quality, but a big loss in the ability to recover crucial information like detecting the pedestrian.

Our goal in this work is to design compression methods able to balance high compression ratios and image quality while preserving the semantic information present in the original image. Exploiting semantic information to guide the image reconstruction process, ensuring the preservation of crucial details of the original image, has already been considered. Recently, prominent approaches for semantic-guided image reconstruction have been built using generative models like Generative Adversarial Networks (GAN) [7]. For example, Isola et al. [5] unveiled the pix2pix model, a conditional GAN that uses a SSM as input and outputs an image that preserves the same semantic information as the original one. Despite its visual allure, this method often overlooks the original image, resulting in a substantially different image, given the fact that the reconstruction is created starting only from the SSM. Wang et al. [6] tackled the problem of reconstructing an image not close enough to the original by suggesting a conditional GAN-centric model that integrates both the SSM and features derived from the original image. This combination ensures better retention of the original content in the reconstructed images. Yet, these techniques often sidestep the pivotal aspect of efficient image compression since they are designed solely for guaranteeing an image that uses as low as possible Bits Per Pixel (BPP), while optimizing a metric that does not distinguish

Caire’s and Musa’s work was partially funded by the German Ministry for Education and Research as BIFOLD – Berlin Institute for the Foundations of Learning and Data. Barbarossa’s work was funded by the Italian National Recovery and Resilience Plan (NRRP) of NextGenerationEU, partnership on “Telecommunications of the Future” (PE00000001 - program RESTART and Huawei Technology France SASU, under agreement N.TC20220919044

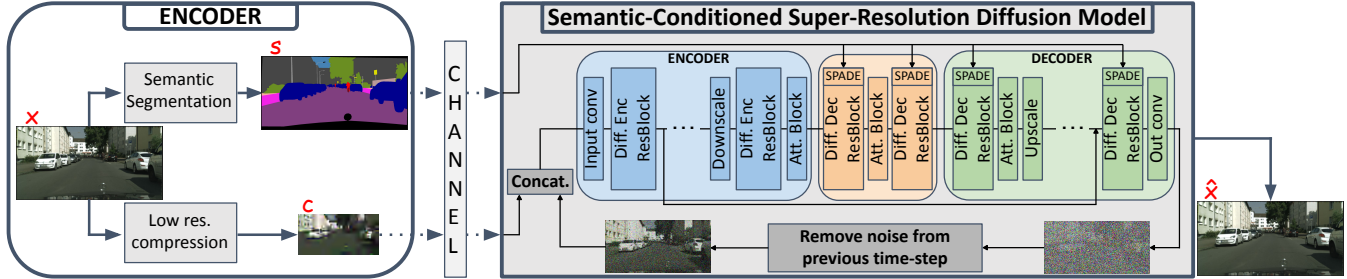


Fig. 1: Overview of the SPIC Architecture. The diagram illustrates our novel approach, combining a Semantic Segmentation Map (s) and a coarse low-resolution image (c), both compressed with classical out-of-the-shelf algorithms for efficient encoding. The reconstruction employs the proposed Semantic-Conditioned Super-Resolution Diffusion Model, leveraging both s and c for high-fidelity semantic-relevant image recovery even at low BPP.

different regions of the image, like PSNR. Instead, we would like to tackle the problem of high-quality image reconstruction and ensure that the method can be used as a valid alternative to classical image compression algorithms, which is essential for real-world applications with bandwidth and storage limitations.

Recently, DDPM [8], a class of generative models that match a data distribution by learning to reverse a gradual multi-step noising process, has exhibited incredible results in image synthesis [9–11]. The authors of [12] improved the works of Isola and Wang et al. introducing a DDPM model that conditions the image generation to its semantic map, hinging on the previous work [13]. The results obtained in these works are promising, but the regenerated images are again obtained considering the SSM solely without taking into account the coarse image.

In this paper, we propose an innovative semantic image communication scheme where the transmitter encodes the SSM losslessly, together with a low-resolution version of the image itself. The receiver uses the proposed Semantic-Conditioned Super-Resolution Diffusion Model (SCSRDM) to regenerate the full-resolution image. While slightly suboptimal with respect to conventional approaches, in terms of the overall rate-distortion curve, the proposed method enables a much better reconstruction and positioning of the semantically relevant objects. The scheme is similar to what proposed in [4], but with some important differences: i) our approach uses a DDPM, as opposed to [4] that uses a GAN, because diffusion models are known for having better image synthesis capabilities [9]; ii) differently from [4], we do not transmit the residual error between the input and the reconstruction, to make our transmitter much simpler to implement and to limit the transmission rate; iii) instead of using a single end-to-end architecture that learns, jointly the Semantic Segmentation Map (s) and the compressed low-resolution image (c), as in [4], we use a modular approach that computes them separately. This simplifies the method considerably, enabling a separate control of the segmentation and compression tasks, using SOTA task-specific algorithms for the SSM generation, e.g. INTERN-2.5 [14], and employing classical compression algorithms, e.g. BPG [15] and FLIF [16], to compress the coarse image and the SSM. From the computational and explainability points of view, the proposed modular approach is more efficient. Exploiting off-the-shelf SOTA components, rather than training a much bigger DNN for the joint approach, allows a model with fewer parameters to train and total control over s and c . Moreover, the modular approach allows the framework to be improved easily, for example, by implementing a new SOTA model for semantic segmentation and replacing only the semantic block without the need to retrain the whole model.

2. PROPOSED METHOD

In this section, we introduce the encoding and decoding parts of our semantic-preserving image coding based on conditional diffusion model.

2.1. Encoder

As illustrated in Figure 1, the encoder is composed of two separate blocks that extract s and c from the input image x .

For the segmentation part, in this work, we used the INTERN-2.5 model [14], known to have high performance in terms of semantic segmentation, but of course, as discussed before, other choices are possible. After generating the SSM, we compress it with a lossless encoder since we assign high priority to the accurate reconstruction of the SSM at the receiver. More specifically, we applied the lossless compression technique FLIF, ensuring efficient encoding with an average of 0.112 BPP. As far as the generation of the coarse image is concerned, we adopted different approaches. The first and simplest one is the average down-scaling operator that shrinks the image dimensions from 256×512 to 64×128 . Based on top of the down-scaled version, to further compress the coarse image before transmission, we employed the BPG compression algorithm on the down-scaled image.

2.2. Semantically-Conditioned Super-Resolution Diffusion Model Decoder

The Decoder takes the received SSM s and the coarse image c and sends them to the SCSRDM, depicted in the right side of Figure 1, whose goal is to reconstruct a Super-Resolution (SR) image, i.e. an image with the same dimension as the original one and similar (or even better) resolution.

At the model’s core lies a U-net structure [17], encompassing three different substructures: an encoder, a central bottleneck, and a decoder. As with every DDPM, synthesizing a singular image necessitates passing through the same U-net multiple times. At each iteration, the model inputs the previous iteration’s output and the conditioning variables to predict the noise to be removed from the input image at that time step. Because of this iterative approach, DDPMs can progressively refine the image, starting from white noise. Because of the random nature of this process, it is necessary to direct the denoising process to avoid a purely random generation disconnected from the original image. Different approaches [10, 18, 19] can be employed to avoid this purely random generation. The two



Fig. 2: (a) Resulting image compressed with the BPG algorithm at 0.176 BPP. (b) Reconstructed image employing our approach at 0.166 BPP. (c) Detail of the image compressed with the BPG. (d) Detail of the image reconstructed with our approach.

main methods are guidance and conditioning; in this work, we adopt a conditioning approach, as it allows us to implement the dual conditioning process, which lies at the core of the proposed SCSRDM, in an efficient manner. As stated before, the idea builds on the very foundation of DDPM and, more specifically, on the concept of SR Diffusion Models [20], properly modified to guarantee the dual conditioning used in our strategy.

Differently from a classical DDPM, our SR Diffusion Model, during the learning phase, instead of starting the denoising process from white noise alone, concatenates the noise with a coarse image expanded to its original size. This conditioning recurs at every step during training, ensuring that the model is consistently driven from the coarse image. Specifically, during training, the model starts with a tensor of dimensions $6 \times 256 \times 512$, with the initial three (colour) channels representing white noise and the subsequent three channels containing the expanded coarse image. Throughout the training, the model undergoes 1000 iterations to operate the transition from the white noise of the first three channels to the reconstructed image and leave the coarse image conditioning unchanged. During inference, only 20 iterations are executed to save time (and energy), using as input always a tensor of size $6 \times 256 \times 512$, but this time substituting the noise of the first three channels with the coarse image itself. This ensures that the starting point is closely aligned with the original image. The second conditioning is the one on the SSM. To do so, we adapted the SPADE technique [13] to our model. As shown in Fig. 1, the conditioning occurs at every ResBlock layer [21] of the bottleneck and decoder subnetworks of the U-net, as also depicted in [12]. In essence, the proposed SCSRDM introduces a contextual diffusion strategy, conditioned on dual inputs, achieving superior SR outcomes, properly steered by the SSM.

3. NUMERICAL RESULTS

In this section, we delve into a comprehensive presentation of the results and advantages associated with the proposed SPIC. It is important to clarify that while the images utilized for these comparative analyses are sourced from the validation folder of the Cityscapes dataset, none of these images were employed during the training or validation phase.

As mentioned before, a paramount advantage of our model is its capability to retain semantic information while able to provide a good trade-off between the overall image quality and compression rate. Several existing compression algorithms and SR models often reconstruct visually pleasing images. However, a closer inspection reveals a significant drawback: the degradation of semantic content, particularly evident as the size of the semantically relevant objects within the image diminishes. For larger foreground objects, most available approaches are able to detect and generate the correct semantic segmentation correctly. However, as the object size shrinks, conventional models falter, failing to accurately process the image and evaluate a precise SSM. This aspect can be grasped by looking at Figure 2: on the left, we see the image reconstructed after compression with a BPG algorithm (a) and its zoom on the center part (c); on the right, we observe the image reconstructed using our approach ((b), and the corresponding detail (d). At first glance, even because of the little advantage in BPP (0.176 vs 0.166), the image on the left looks clearer and more detailed than the one on the right. But, as soon as we zoom in, the story is completely different because our reconstruction clearly shows a person on a bicycle on the right and a pedestrian in the distance, which are not at all clear in the left image. As a further example, in Figure 3, we compare the reconstruction capabilities of our model with the SOTA SR model introduced in [22].

Both models are evaluated on their ability to amplify the image size by a factor of four, transitioning from 128×64 to 512×256 pixels, without any further source compression. The distinguishing difference between the two approaches is that our model is conditioning the reconstruction on the SSM. Looking at Fig. 3, which reports the zoom on the central part of the reconstructed images, we can see that the resolution of our model is better and, more specifically, the three pedestrians between the two cars and the road signs are clearly visible in our case, while they are only barely observable using the SOTA SR model proposed in [22].

To compare the performance of our model with available alternatives, in terms of semantic segmentation retention, we used as a performance metric the mean Intersection over Union (mIoU), a number that quantifies the degree of overlap between the ground truth and the predicted regions corresponding to the objects of interest. More specifically, given two boxes s_1^i and s_2^i , with $i = 1, \dots, n_c$, computed over the ground truth and the reconstructed image, where n_c is the number of classes of objects of interest, the $\text{mIoU}(s_1, s_2)$ is defined as follows:

$$\text{mIoU}(s_1, s_2) = \frac{1}{n_c} \sum_{i=1}^{n_c} \text{IoU}(s_1^i, s_2^i) = \frac{1}{n_c} \sum_{i=1}^{n_c} \frac{|s_1^i \cap s_2^i|}{|s_1^i \cup s_2^i|}$$

In the given semantically-preserving coding scheme, the quality of the reconstructed image cannot be assessed by using conventional metrics, like PSNR, which focus only on a pixel-by-pixel reconstruction, and then fail to capture the semantic content. For this reason, since the reconstructed images are also sensitive to various types of distortion, such as blurriness, noise, and artifacts, we assess the difference between the original and reconstructed images in terms of the Fréchet Inception Distance (FID) [23], a widely used metric in computer vision, which compares the features maps extracted by an Inception-v3 model [24], and is expressed as follows:

$$\text{FID} = \|\mu_{f(x)} - \mu_{f(y)}\|^2 + \text{Tr}(\Sigma_{f(x)} + \Sigma_{f(y)} - 2(\Sigma_{f(x)}\Sigma_{f(y)})^{1/2})$$

with $f(\cdot)$ the output of the *pool3* layer of the Inception-v3, and μ and Σ are the mean and covariance matrix of the 2048 feature vectors.



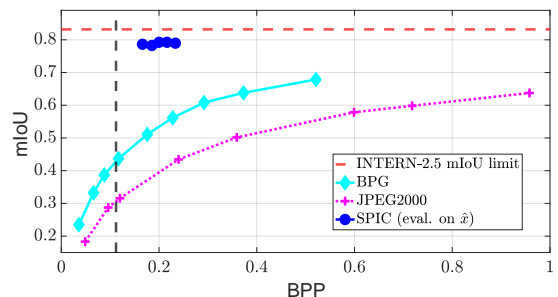
Fig. 3: Detail comparison between (top) the image reconstructed with the SOTA SR model [22] and (bottom) the image reconstructed with our model

In Figure 4 (a) and (b) we report the mIoU and the FID, vs. the BPP, used to encode the transmitted data, averaged over the whole

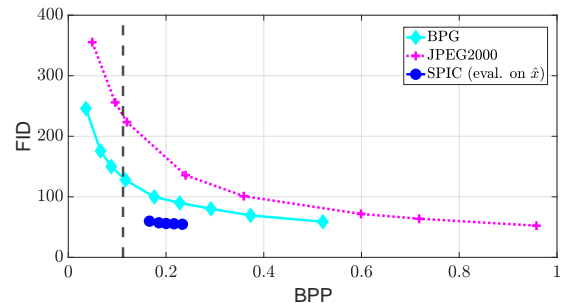
validation dataset. All the SSM are generated using the INTERN-2.5 model. In Figure 4a, the black dotted vertical line, positioned at 0.112 BPP, represents the BPP required for the lossless compression of the SSM. The blue point represents the mIoU evaluated on the reconstructed images \hat{x} obtained applying SCSRDM at a BPP given by the sum of the BPP necessary for the lossless encoding of the SSM and the lossy encoding of the coarse image. The green and magenta curves represent the results achieved with BPG and JPEG2000 compression methods. We can clearly see that both BPG and JPEG2000 exhibit worse performance than our method in terms of mIoU. To let BPG be able to achieve mIoU results akin to our model, the rate should be in the order of 1 BPP.

Looking now at Figure 4b, we can see that while being able to retain most of the semantic segmentation information, our method can reconstruct images that have a low FID score, outperforming both JPEG2000 and BPG.

In summary, our numerical results show that the proposed method, compared to alternative approaches, achieves a better balance between fidelity reconstruction and ability to extract semantic features from the reconstructed image.



(a) mIoU vs BPP



(b) FID vs BPP

Fig. 4: Comparison in terms of mIoU (a) and FID (b) vs. BPP.

4. CONCLUSIONS

In this work, we propose a novel image coding scheme, building on a doubly conditioned super-resolution diffusion model, able to better preserve the semantic content of the image than SOTA compression algorithms and SR methods while at the same time, having a better rate/quality trade-off when compared to the best compression methods. The proposed model harnesses the power of dual conditioning on a SSM and a compressed version of the original image. The double conditioning is obtained with a modular framework that allows SPIC to be easily adapted to different tasks. Future investigations include the extension to semantic video coding and the incorporation of errors due to transmission over a noisy channel.

5. REFERENCES

- [1] Xuewen Luo, Hsiao-Hwa Chen, and Qing Guo, “Semantic communications: Overview, open issues, and future research directions,” *IEEE Wireless Communications*, vol. 29, no. 1, pp. 210–219, 2022.
- [2] Sergio Barbarossa, Danilo Comminiello, Eleonora Grassucci, Francesco Pezone, Stefania Sardellitti, and Paolo Di Lorenzo, “Semantic communications based on adaptive generative models and information bottleneck,” *IEEE Communications Magazine*, vol. 61, no. 11, pp. 36–41, Nov 2023.
- [3] Jia lin Xu, Tze-Yang Tung, Bo Ai, W. Chen, Yuxuan Sun, and Deniz Gunduz, “Deep joint source-channel coding for semantic communications,” *ArXiv*, vol. abs/2211.08747, 2022.
- [4] Mohammad Akbari, Jie Liang, and Jingning Han, “Dsslic: Deep semantic segmentation-based layered image compression,” in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 2042–2046.
- [5] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros, “Image-to-image translation with conditional adversarial networks,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 5967–5976.
- [6] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro, “High-resolution image synthesis and semantic manipulation with conditional gans,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [7] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, “Generative adversarial nets,” in *Advances in Neural Information Processing Systems*, Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, Eds. 2014, vol. 27, Curran Associates, Inc.
- [8] Jonathan Ho, Ajay Jain, and Pieter Abbeel, “Denoising diffusion probabilistic models,” in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, Eds. 2020, vol. 33, pp. 6840–6851, Curran Associates, Inc.
- [9] Prafulla Dhariwal and Alexander Quinn Nichol, “Diffusion models beat GANs on image synthesis,” in *Advances in Neural Information Processing Systems*, A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, Eds., 2021.
- [10] Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song, “Denoising diffusion restoration models,” in *ICLR Workshop on Deep Generative Models for Highly Structured Data*, 2022.
- [11] Alexander Quinn 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,” in *International Conference on Machine Learning, ICML 2022*, 2022, vol. 162 of *Proceedings of Machine Learning Research*, pp. 16784–16804, PMLR.
- [12] Weilun Wang, Jianmin Bao, Wengang Zhou, Dongdong Chen, Dong Chen, Lu Yuan, and Houqiang Li, “Semantic image synthesis via diffusion models,” arXiv:2207.00050, 2022.
- [13] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu, “Semantic image synthesis with spatially-adaptive normalization,” in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 2332–2341.
- [14] Wenhai Wang, Jifeng Dai, Zhe Chen, Zhenhang Huang, Zhiqi Li, Xizhou Zhu, Xiaowei Hu, Tong Lu, Lewei Lu, Hongsheng Li, et al., “Internimage: Exploring large-scale vision foundation models with deformable convolutions,” *arXiv preprint arXiv:2211.05778*, 2022.
- [15] Fabrice Bellard, “Better portable graphics image format,” (<http://bellard.org/bpg/>), 2017.
- [16] Jon Sneyers and Pieter Wuille, “Flif: Free lossless image format based on maniac compression,” in *2016 IEEE International Conference on Image Processing (ICIP)*, 2016, pp. 66–70.
- [17] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, *U-Net: Convolutional Networks for Biomedical Image Segmentation*, pp. 234–241, Springer International Publishing, Cham, 2015.
- [18] Konpat Preechakul, Nattanat Chatthee, Suttisak Wizadwongsa, and Supasorn Suwajanakorn, “Diffusion autoencoders: Toward a meaningful and decodable representation,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [19] Jonathan Ho and Tim Salimans, “Classifier-free diffusion guidance,” in *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.
- [20] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J. Fleet, and Mohammad Norouzi, “Image super-resolution via iterative refinement,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 4, pp. 4713–4726, 2023.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [22] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer, “High-resolution image synthesis with latent diffusion models,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 10684–10695.
- [23] 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,” in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, Red Hook, NY, USA, 2017, NIPS’17, p. 6629–6640, Curran Associates Inc.
- [24] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna, “Rethinking the inception architecture for computer vision,” *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2818–2826, 2015.