

Image Outpainting and Harmonization using Generative Adversarial Networks

Basile Van Hoorick
Columbia University

basile.vanhoorick@columbia.edu

Abstract Although the inherently ambiguous task of predicting what resides beyond all four edges of an image has rarely been explored before, we demonstrate that GANs hold powerful potential in producing reasonable extrapolations. Two outpainting methods are proposed that aim to instigate this line of research: the first approach uses a context encoder inspired by common inpainting architectures and paradigms, while the second approach adds an extra post-processing step using a single-image generative model. This way, the hallucinated details are integrated with the style of the original image, in an attempt to further boost the quality of the result and possibly allow for arbitrary output resolutions to be supported.

I. INTRODUCTION

When presented with an incomplete image, humans are excellent at filling in the blanks and producing a realistic explanation for what could be missing. Image inpainting is a well-studied problem that replicates this behaviour, often tasking deep neural networks with trying to understand the semantic content of natural images in order to recover the missing regions of a photo. However, the spatially inverted variant of this problem is even more challenging and, with a small play on words, can be denoted *outpainting*. The problem statement is shown in Figure 1; essentially, the task is to extrapolate the image content rather than to interpolate within an image. More formally, we must design a generator G that converts an image x with dimensions $n * n$ into a larger image $G(x)$ with dimensions $m * m$, such that the center part of $G(x)$ looks the same as x , while the complete outpainted image $G(x)$ should be a plausible hypothesis of what could encompass the original image. In particular, the generated parts should look realistic, despite the fact that it is often impossible to know precisely what the scene contains outside of the pictured boundaries. Figure 2 gives a practical illustration of what is accomplished in this work. Possible applications of outpainting include the opportunity to experience and explore multimedia more immersively, enabling technologies similar to Ambilight to produce much richer and/or realistic surroundings, or merely as a form of computer-generated art that can be appreciated aesthetically.

A. Related Work

Deep neural networks have recently shown great performance in image completion. This section discusses previous work related to the proposed methods for outpainting.

Self-supervised learning The automatic generation of a supervision signal, by systematically omitting certain parts of the input data, grants the ability to capitalize on the vast amounts of unlabeled images and videos that are available on the Web. For example, both col- orization [1] and predicting the relative spatial position

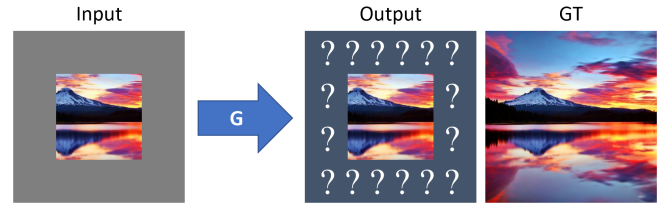


FIG. 1: Image outpainting idea.



FIG. 2: Demonstration of the outpainting task, as achieved using an adversarially trained generator.

of two patches [2, 3] have been proven to be successful pretext tasks to representation learning. These methods do not require any labellings of the dataset, since the inputs (gray-scale images resp. patches) as well as their corresponding ground truths (color images resp. locations) can be extracted procedurally.

Generative Adversarial Networks Image or video generation using convolutional neural networks with an adversarial loss [4] has attracted significant research interest. The idea is to let a *generator* network G create samples according to some distribution, and subsequently have an auxiliary network called the *discriminator* D try to distinguish whether a given sample is valid or was actually generated by G . These two components are



FIG. 3: Demonstrating Snapseed’s [11] unrealistic results upon running the ‘Expand’ feature on a natural photograph.

trained in a quick alternation, where G trying to fool D forces D to become better at telling real from fake, which in turn motivates G to synthesize increasingly convincing outputs [5]. From the perspective of game theory, GANs can be viewed as a zero-sum two-player game; as such, they are comparatively hard to train. In fact, it turns out that the conditions for their convergence is still an open research problem [6].

Inpainting and Context Encoders Recovering holes within images has important applications for the restoration of damaged media, and can be achieved by training GANs in a self-supervised way. Specifically, a generator can learn to ‘fill in the blanks’ by using an encoder-decoder network that encodes the surrounding context in order to understand the image content, and subsequently produces a plausible hypothesis for the missing square [7]. More recent improvements include attempts to complete images of arbitrary resolutions by filling in regions of any shape, and employing two types of discriminators (*global* and *local*) in order to enforce both overall consistency as well as realistic details [5].

Outpainting The problem opposite to inpainting comprises predicting which pixels reside beyond the borders of a fully intact photo, and has to our knowledge been explored only few times before by the academic research community [8–10]. One approach geometrically extrapolates the field of view of an image using another panoramic reference image of the same scene category using old-school computer vision techniques [8], while a more recent relevant paper uses a GAN to perform horizontal outpainting [9]. Results look promising, although there is room for visual improvement in terms of making the hallucinations omnidirectional and increasing the output resolution. [10] achieves impressive outcomes but tends to focus on limited domain datasets (such as faces only), and notes that generative models experience difficulties trying to fit datasets as diverse as Places2. Lastly, Google’s Snapseed application has a proprietary ‘Expand’ functionality that seems to select patches from the image and copy them to the edges [11], but a limited number of experiments suggest that this tool fails to capture the local structure of most scenes from region to region; see Figure 3.

Single-image generative models The recently introduced SinGAN framework can train an unconditional generative model on a single natural photo, captur-

Type	Kernel size	(H, W, C)
Conv + Leaky ReLU	$4 * 4$	(96, 96, 64)
Conv + BatchNorm + Leaky ReLU	$4 * 4$	(48, 48, 64)
Conv + BatchNorm + Leaky ReLU	$4 * 4$	(24, 24, 128)
Conv + BatchNorm + Leaky ReLU	$4 * 4$	(12, 12, 256)
Conv + BatchNorm + Leaky ReLU	$4 * 4$	(6, 6, 512)
Conv	$4 * 4$	(3, 3, 4000)
Up-conv + BatchNorm + ReLU	$4 * 4$	(6, 6, 512)
Up-conv + BatchNorm + ReLU	$4 * 4$	(12, 12, 256)
Up-conv + BatchNorm + ReLU	$4 * 4$	(24, 24, 128)
Up-conv + BatchNorm + ReLU	$4 * 4$	(48, 48, 64)
Up-conv + BatchNorm + ReLU	$4 * 4$	(96, 96, 64)
Up-conv + BatchNorm + Tanh	$4 * 4$	(192, 192, 3)

TABLE I: Layers of the generator network G .

ing and reproducing image statistics across various scales of the image [12]. It allows for the creation of random samples with new object configurations by starting from a low-resolution sample at the coarsest scale, and then progressively upsampling and refining the result through the pyramid of generators. By carefully modifying the initial sample, SinGAN can be used to perform several image manipulation tasks such as editing photos, splicing foreign objects into the scene and subsequently harmonizing their style with the environment, or even super-resolution.

II. METHODS

A. Datasets

During training, we decide to crop images first before feeding them into the generator. Consequently, learning can be done in a self-supervised way thanks to the fact that we are now able to enforce the output to approximate the original, uncropped image. Any sufficiently large dataset of unlabelled, natural photos will therefore suffice. We continue with the MIT CSAIL Places365-Standard dataset, which contains millions of images with landscapes, buildings and other everyday scenarios intended for scene recognition [13]. A second set of experiments was performed with a dataset consisting of images scraped from WikiArt [14]. This website holds a database of visual artwork belonging to many different categories.

The photos will first be resized to 192x192 as a preprocessing step, and the generator will then be tasked with expanding a crop of 128x128 back into a 192x192 image. Note that adding 32-pixels at every boundary might not seem significant, but we consider this configuration to be quite ambitious already: the total output size is 2.25 times that of the input, meaning that over half of all pixels will be hallucinated. In practice, the generator G maps a partially masked 192x192 color image to an out-painted variant of the same dimensions, with the masked part replaced by the model’s predictions.

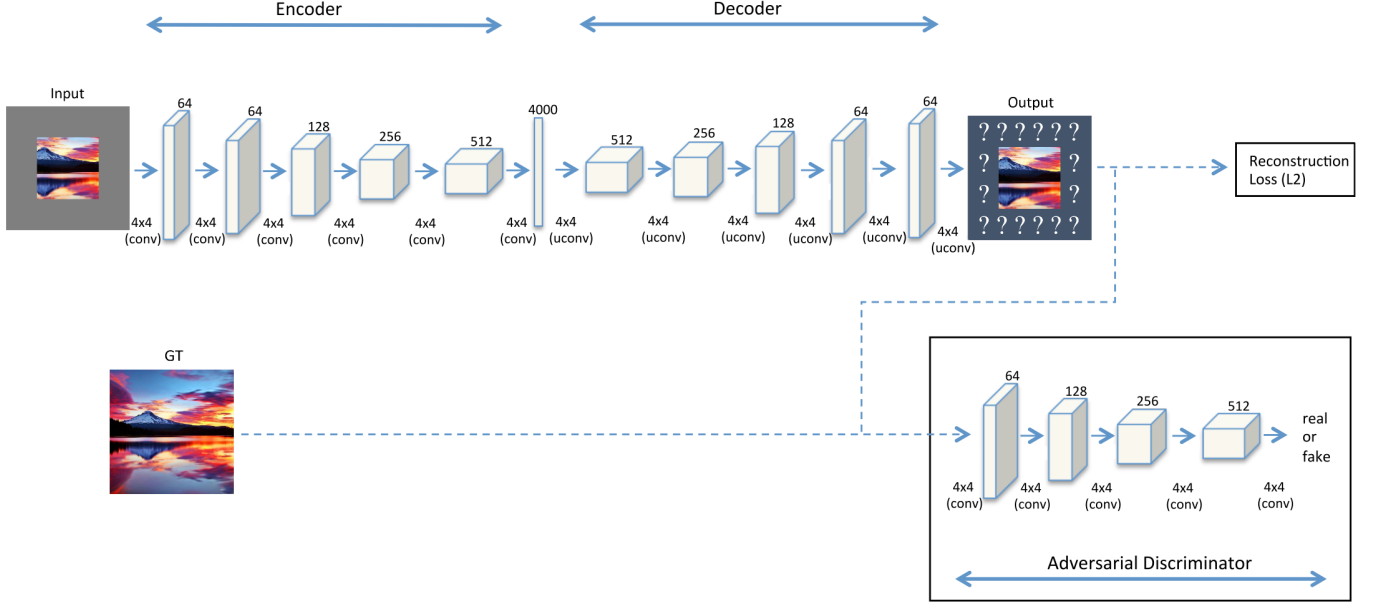


FIG. 4: Context encoder trained with joint reconstruction and adversarial loss for semantic outpainting, based on [7]. Note that as opposed to inpainting, the decoder is almost an exact mirroring of the encoder, since both input and output have the same dimensions.

Type	Kernel size	(H, W, C)
Conv + Leaky ReLU	3×3	(96, 96, 64)
Conv + InstanceNorm + Leaky ReLU	3×3	(48, 48, 128)
Conv + InstanceNorm + Leaky ReLU	3×3	(24, 24, 256)
Conv + InstanceNorm + Leaky ReLU	3×3	(24, 24, 512)
Conv	3×3	(24, 24, 1)

TABLE II: Layers of the discriminator network D .

B. Architecture

Many architectural aspects and ideas can be naturally adopted from inpainting. The context encoder part of the generator network G repeatedly downsamples the masked input through six convolutional layers, in order to efficiently capture the image content and object semantics within some embedding space. Next, the decoder consists of a special kind of layers called *up-convolutional* or *deconvolutional*, which can be understood as having a fractional stride in order to ‘undo’ the downsampling performed by the encoder [7].

The discriminator D is another deep neural network that estimates the probability of the ground truth or the hallucinated image being real. In inpainting, D typically sees the generated part only [7], although in this project we decide to operate on the full outpainted image in order to discourage G from introducing obvious perceptual discontinuities or other kinds of structural inconsistencies. The architecture we use for D results in a 24×24 grid of probabilities whose errors are averaged during the training process.

See Figure 4 as well as Tables I and II for a detailed

overview of the system architecture.

C. Training

Training is done for 200 epochs, with a fixed learning rate of $\alpha = 0.0003$ and two Adam optimizers with $\beta_1 = 0.5, \beta_2 = 0.999$. The loss functions are as follows:

$$L_{rec} = \|x - G(x)\|_1 \quad (2.1)$$

$$L_{adv} = \|D(G(x)) - 1\|_2^2 \quad (2.2)$$

$$L_G = \lambda_{rec} L_{rec} + \lambda_{adv} L_{adv} \quad (2.3)$$

$$L_D = \|D(x) - 1\|_2^2 + \|D(G(x)) - 0\|_2^2 \quad (2.4)$$

Using an L_1 reconstruction loss instead of L_2 helps produce less blurry images [15]. The weight of the adversarial loss λ_{adv} relative to the reconstruction loss $\lambda_{rec} = 1 - \lambda_{adv}$ turned out to be particularly tricky to adjust; this factor was initially set to $\lambda_{adv} = 0.001$ as in various existing works [7, 9, 16], although the GAN kept collapsing into a failure mode where the adversarial loss did not move away from 1. This means that G was unable to fool D , so D was ahead of G and could always tell real from fake successfully. A working remedy involved varying $\lambda_{adv}(n)$ throughout time as a function of the epoch n as follows:

$$\lambda_{adv}(n) = \begin{cases} 0.001, & \text{if } n \leq 10 \\ 0.005, & \text{if } 10 < n \leq 30 \\ 0.015, & \text{if } 30 < n \leq 60 \\ 0.040, & \text{otherwise} \end{cases} \quad (2.5)$$

This will punish the generator more heavily for producing unrealistic outputs as time progresses, rather than just

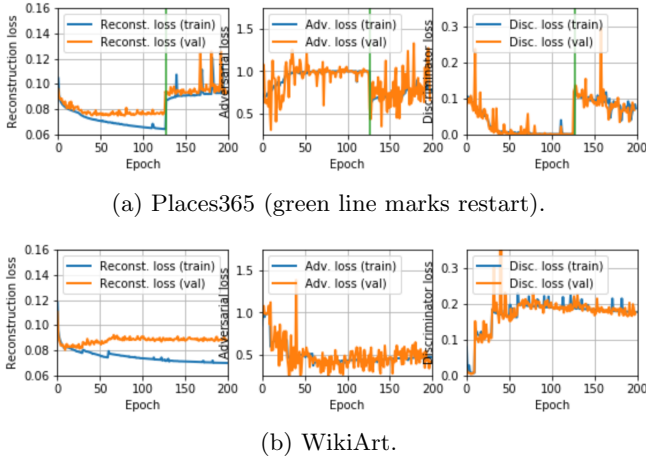


FIG. 5: Learning curves for the outpainting GAN, where the difficulty of training GANs is observed.

enforcing an accurate pixel-wise reconstruction.

D. Harmonization

Due to the visual fuzziness of our initial experiments, we came up with the idea of leveraging SinGAN’s harmonization capabilities in an attempt to improve the fidelity of the hallucinated outputs, serving as a second possible approach to the outpainting problem. To this end, we first train a SinGAN model on the original high-quality image, then propagate it forward through the outpainting generator (which produces a low-resolution output), and then try to super-resolve this result by injecting it into one of the coarser scales of SinGAN. The hope is that the model will harmonize the outpainted parts with the style of the original image that it was trained on, while simultaneously synthesizing a finer-scale, higher-resolution variant by pushing it up through the hierarchy of multi-scale generators.¹

III. RESULTS

A. Experiments

Three models were trained and evaluated, differing in the following aspects:

- 200k images of Places365, with $L_G = L_{rec}$, so only the L_1 reconstruction loss is considered.
- 200k images of Places365, with fully functional joint reconstruction and adversarial loss functions.
- 50k images of WikiArt, with fully functional joint loss functions.

¹ Note that our definition of the term *harmonization* is slightly different from what is actually meant in the referred paper, but the underlying process is exactly what we intend to perform anyway.

Model	Mean $MSE(G(x), y)$	Mean $D(G(x))$
Places365 rec	0.0181	0.0791
Places365 rec + adv	0.0230	0.1705
WikiArt rec + adv	0.0227	0.2371

TABLE III: Mean square error and realism (i.e. probability of legitimacy) statistics.

Since training takes a long time, the second model’s training process was momentarily interrupted and resumed in order to upgrade the adversarial loss weight to the better performing Equation 2.5. Selected learning curves are plotted in Figure 5. Note that the discriminator D seems to have a harder time deciding for the WikiArt dataset than for Places365, both in terms of adversarial loss L_{adv} and regular discriminator loss L_D . This suggests a priori that the artwork domain is inherently more diverse than natural scenery in a visual sense, making it harder to reliably classify the authenticity of samples.

B. Evaluations

In the final step of the outpainting pipeline, the input is pasted on top of and blended with the output using a simple 16-pixel wide gradient alpha mask. Compression artefacts arising from the autoencoder-like operation at the center of the image are thus resolved by simply hiding them. However, two quantitative performance metrics, the mean squared error and the discriminator’s output values, are calculated from the original, unmodified result. See Table III for an overview of these metrics for the three trained models as calculated on the respective test sets. Including the adversarial loss increases the reconstruction loss, and outpainted artwork fools the discriminator more often than when we stay within the domain of natural images.

Non-adversarial versus adversarial models See Figure 6 for samples of generators trained with $\lambda_{adv} = 0$ and $\lambda_{adv} = 0.040$ respectively. If only the L_1 reconstruction loss is used, then the hallucinated part looks fuzzy and unrealistic. Enabling the discriminator clearly has a sharpening effect on the output, although the results start to look excessively decorated at times. It seems that the trade-off between the reconstruction and adversarial loss functions should be more finely investigated.

Variation in reconstruction quality See Figure 7 for samples where the MSE between the output and ground truth is either particularly low or high, when including the adversarial loss. It seems that ‘smooth’ images such as landscapes perform best, while highly detailed images including indoor scenes perform worst. This is understandable: if the dimensions of objects within the image are smaller than what is being omitted, they become impossible for the outpainter to predict.

Artwork outpainting See Figure 8 for a few examples of generative outpainting performed on the WikiArt validation- and test sets. In this application, we answer the question ‘What would the artist have drawn if he/she

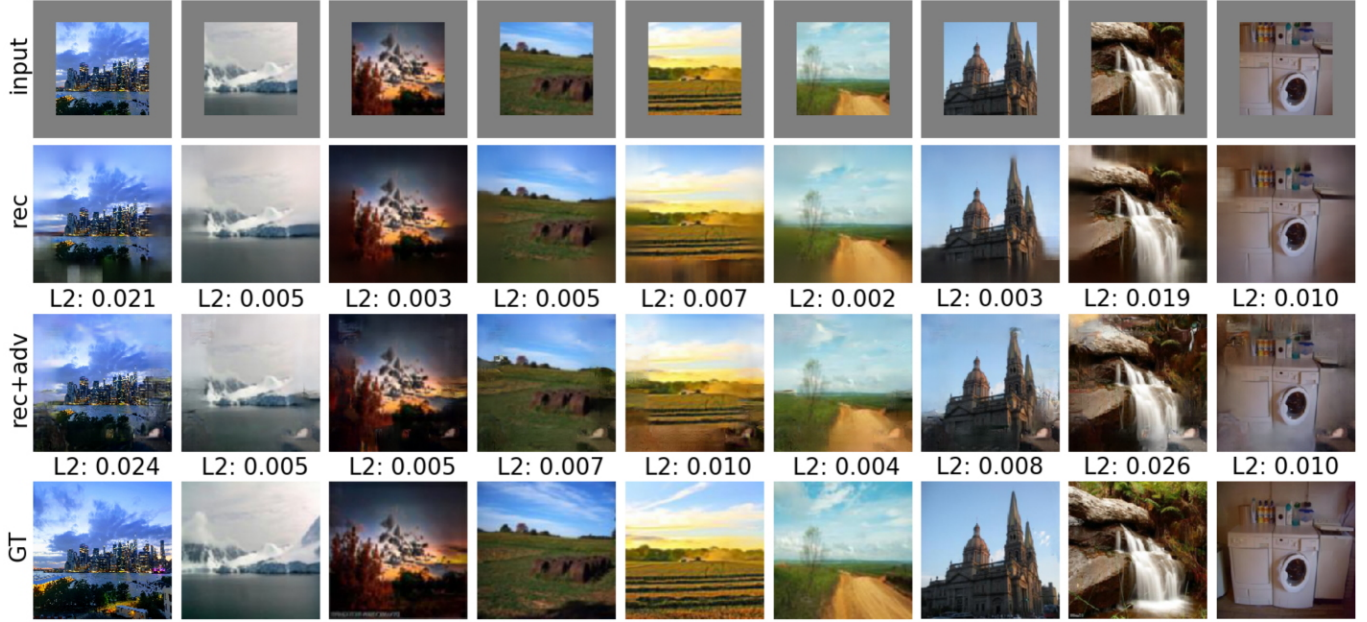


FIG. 6: Comparison between outpainting on Places365 without or with adversarial loss, demonstrating blurry results in the former case.

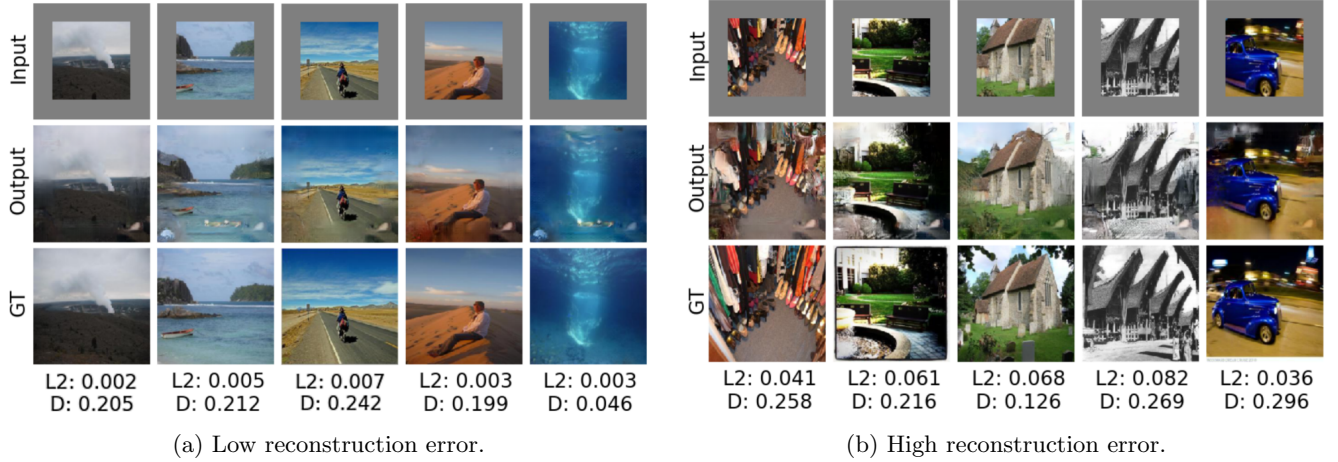


FIG. 7: Examples of outpainting on Places365 selected for varying MSEs (lower is better), revealing semantic differences within the images.

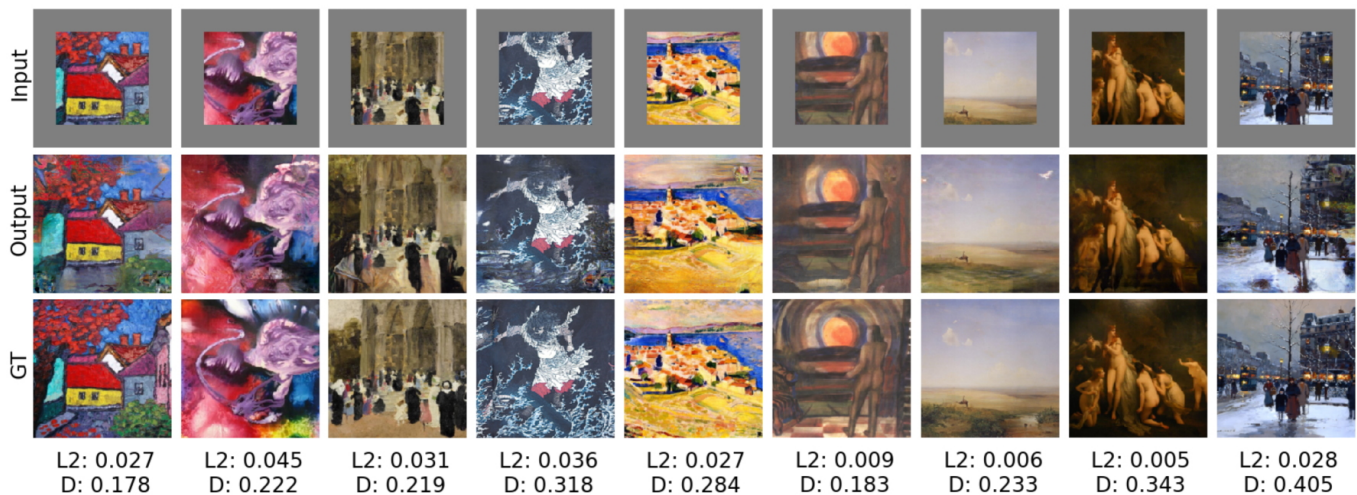


FIG. 8: Qualitative illustration of artwork outpainting on WikiArt.

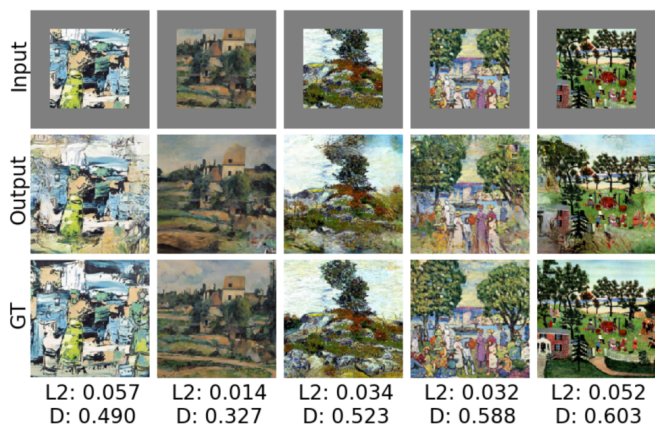


FIG. 9: Examples of outpainting on WikiArt selected for high discriminator output values, indicating a bias towards highly decorated paintings.

had a bigger canvas?'. The results look comparatively accurate and visually pleasing, for which we think the reason is that 'imperfections' in art (to the extent that term is even well-defined) are much less bothersome than in natural imagery. It is also interesting to observe that when the output's realism $D(G(x))$ is maximized, a bias towards specific types of artwork content and/or genres can be perceived as in Figure 9.

C. Recursive outpainting

In contrast to inpainting, outpainting does not have to be performed just once; in theory, there is no limit as to how many times we can extrapolate a single image. Figure 10 reveals an interesting aspect of the outpainting network: it seems that given enough iterations, the images eventually start to converge to some kind of 'eigen-mode' of the generator. We moreover observe that the painting of a park with flowers seems to slowly morph

into a landscape with clouds as we keep hallucinating outwards, suggesting that this scenery must appear (disproportionately) often in the training set.

D. Harmonization

Our second method for outpainting introduces an additional step that involves employing SinGAN as described earlier. The biggest drawback of this approach is that it takes around an hour to train SinGAN on a single input image, in contrast to the outpainter itself which is a model that can simply be applied anytime once it is trained on a dataset. Nevertheless, we tested out the joint harmonization and super-resolution process on several images that demonstrate its potential, as seen in Figure 11. When applied on the output of the reconstruction loss minimizer, it has a sharpening effect. Furthermore, despite the fact that the end result still lacks photorealism, we observe fewer glitches than the first approach which uses an adversarially trained GAN directly. Arbitrarily large resolutions could in principle be supported by SinGAN, although in practice a huge amount of GPU memory ($> 12\text{GB}$) is needed to scale beyond a maximum image dimension of around 500 pixels.

IV. CONCLUSION

Image outpainting is a novel but exciting idea that holds promise, especially when cascaded with SinGAN to further increase the output fidelity. The models trained in this project still contain some glitches, but we believe these could be alleviated by closer investigations into the subject. Non-photorealistic images such as artwork seem to produce convincing results, which we largely attribute to human judgement becoming more permissive rather than to a better-performing model.



FIG. 10: Recursive outpainting on both natural photos and artwork, revealing some glitches but also intriguing 'eigenmodes' of the network.

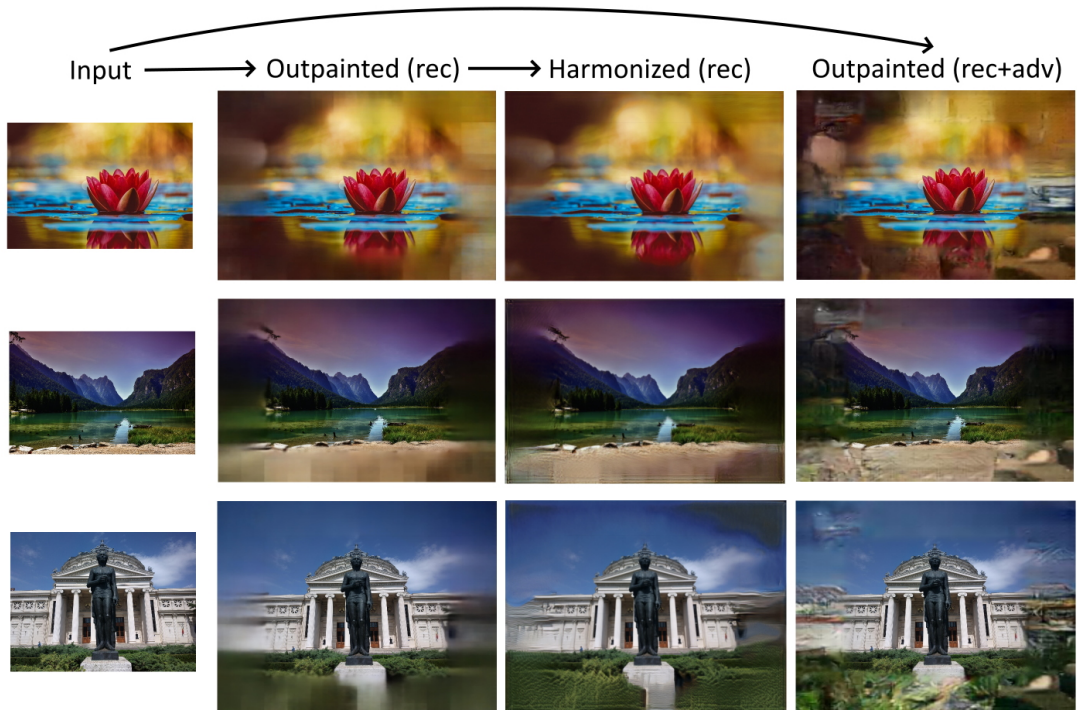


FIG. 11: Qualitative comparison of outpainting without adversarial loss followed by SinGAN harmonization (which has an adversarial aspect by nature), versus direct adversarial outpainting with the proposed GAN.

V. FUTURE WORK

Other than architectural improvements, possible directions for future research include:

- While enabling the discriminator in this project significantly sharpened the outpainting results, the training stability could be better and the model was not very successful in enforcing natural object semantics across the image. Instead of a single discriminator, utilizing both a global and local variant might increase the overall consistency of generated images as well as improve the realism of smaller-scale structures [5].
- At the present time, it is unknown how the model

reacts to slight variations in the input space. Hence, video outpainting presents a set of new interesting challenges, mostly related to temporal consistency of the hallucinations.

ACKNOWLEDGEMENTS

I would like to thank Prof. Peter Belhumeur for providing the exciting opportunity to conduct this research project, and for offering a great learning opportunity in his excellent course *Deep Learning for Computer Vision*. I am also grateful to John Daciuk for giving me access to his WikiArt dataset, thus saving me time otherwise spent scraping while also allowing for aesthetically appealing experiments to be conducted more easily.

This work was performed while being the recipient of a Belgian American Educational Foundation Fellowship.

-
- [1] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. Colorization as a proxy task for visual understanding. *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, 2017-Janua:840–849, 2017.
 - [2] Carl Doersch, Abhinav Gupta, and Alexei A. Efros. Unsupervised visual representation learning by context prediction. *Proceedings of the IEEE International Conference on Computer Vision*, 2015 Inter:1422–1430, 2015.
 - [3] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 9910 LNCS:69–84, 2016.
 - [4] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative Adversarial Networks. pages 1–9, 2014.
 - [5] Satoshi Iizuka, Edgar Simo-Serra, and Hiroshi Ishikawa. Globally and locally consistent image completion. *ACM Transactions on Graphics*, 36(4), 2017.
 - [6] Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for GANs do actually converge? *35th International Conference on Machine Learning, ICML 2018*, 8:5589–5626, 2018.
 - [7] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A. Efros. Context Encoders: Feature Learning by Inpainting. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2016-Decem:2536–2544, 2016.
 - [8] Yinda Zhang, Jianxiong Xiao, James Hays, and Ping Tan. Framebreak: Dramatic image extrapolation by guided shift-maps. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 1171–1178, 2013.
 - [9] Mark Sabini and Gili Rusak. Painting Outside the Box: Image Outpainting with GANs. 2018.
 - [10] Yi Wang, Xin Tao, Xiaoyong Shen, and Jiaya Jia. Wide-context semantic image extrapolation. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2019-June:1399–1408, 2019.
 - [11] Snapseed - apps on google play, Jun 2018.
 - [12] Tamar Rott Shaham, Tali Dekel, and Tomer Michaeli. SinGAN: Learning a Generative Model from a Single Natural Image. may 2019.
 - [13] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 Million Image Database for Scene Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(6):1452–1464, 2018.
 - [14] Wikiart: Visual art encyclopedia.
 - [15] He Huang, Philip S. Yu, and Changhu Wang. An Introduction to Image Synthesis with Generative Adversarial Nets. pages 1–17, 2018.
 - [16] Haofeng Li, Guanbin Li, Liang Lin, Hongchuan Yu, and Yizhou Yu. Context-Aware Semantic Inpainting. *IEEE Transactions on Cybernetics*, 14(8):1–12, 2018.