

Boundless: Generative Adversarial Networks for Image Extension

Piotr Teterwak Aaron Sarna Dilip Krishnan Aaron Maschinot
David Belanger Ce Liu William T. Freeman

Google Research

{pteterwak, sarna, dilipkay, amaschinot, dbelanger, celiu, wfreeman}@google.com

Abstract

Image extension models have broad applications in image editing, computational photography and computer graphics. While image inpainting has been extensively studied in the literature, it is challenging to directly apply the state-of-the-art inpainting methods to image extension as they tend to generate blurry or repetitive pixels with inconsistent semantics. We introduce semantic conditioning to the discriminator of a generative adversarial network (GAN), and achieve strong results on image extension with coherent semantics and visually pleasing colors and textures. We also show promising results in extreme extensions, such as panorama generation.

1. Introduction

Across many disparate disciplines there exists a strong need for high quality image extensions. In virtual reality, for example, it is often necessary to simulate different camera extrinsics than were actually used to capture an image, which generally requires filling in content outside of the original image bounds [19]. Panorama stitching generally requires cropping the jagged edges of stitched projections to achieve a rectangular panorama, but high quality image extension could enable filling in the gaps instead [23]. Similarly, extending videos has been shown to create more immersive experiences for viewers [3]. As televisions transition to the 16:9 HDTV aspect ratio, it is appealing to display videos filmed at a different aspect ratio than the screen. [22, 33].

We desire a seamless blending between the original and extended image regions. Moreover, the extended region should match the original at the textural, structural and semantic levels, while appearing a plausible extension. Boundary conditions are only available on one side of the extended region. This is in contrast to the image inpaint-

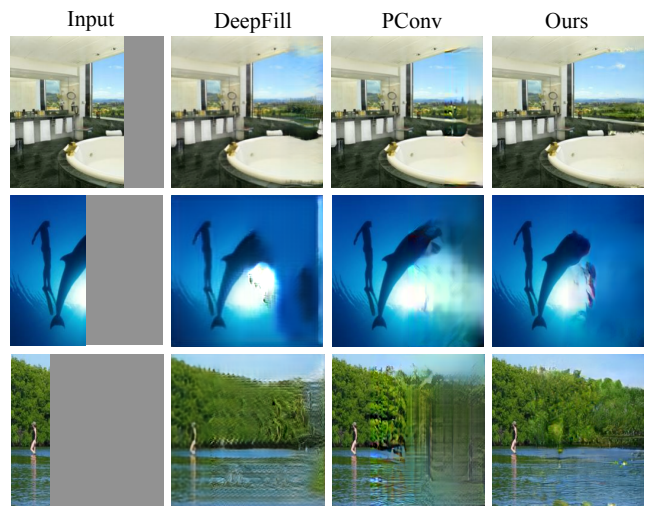


Figure 1. Some examples of image extension: Our method (right column) generates better object shapes (top/middle rows) and produce good textures (middle/bottom rows), compared with two state of the art inpainting methods: DeepFill [48] and PConv [26]. The input image is extended onto the masked area (shown in gray).

ing problem [26, 48], where the region to be filled in is surrounded in all directions by original image data, significantly constraining the problem. Therefore, inpainting algorithms tend to have more predictable and higher quality results than image extension algorithms. In fact, we demonstrate in this paper that using inpainting algorithms with no modifications leads to poor results for image extension.

In the literature, image extension has been studied using both parametric and non-parametric methods [41, 50, 31, 7, 5]. While these methods generally do a good job of blending the extended and original regions, they have significant drawbacks. They either require the use of a carefully chosen guide image from which patches are borrowed, or they mostly extend texture, without taking into account larger scale structure or the semantics of an image. These models are only applicable in a narrow range of use cases and can-

not learn from a diverse data set. In practice, we would like image extension models that work on diverse data and can extend structure.

Fast progress in deep neural networks has brought the advent of powerful new classes of image generation models, the most prominent of which are generative adversarial networks (GANs) [14] and variational autoencoders [21]. GANs in particular have demonstrated the ability to generate high quality samples. In this paper, we use GANs, modified as described below, to learn plausible image extensions from large datasets of natural images using self-supervision, similar in spirit to the use of GANs in applications such as inpainting [16] and image superresolution [24].

For the image extension problem, while state-of-the-art inpainting models [48, 47] provide us a good starting point, we find that the results quickly degrade as we extend further from the image border. We start by pruning the components that do not apply to our setting and then adopt some techniques from the broader study of GANs. Finally, we introduce a novel method, derived from [29], of providing the model with semantic conditioning, that substantially improves the results. In summary, our contributions are:

1. We are one of the first to use GAN’s effectively to learn image extensions, and do so reliably for large extrapolations (up to 3 times the width of the original).
2. We introduce a stabilization scheme for our training, based on using semantic information from a pre-trained deep network to modulate the behavior of the discriminator in a GAN. This stabilization scheme is useful for any adversarial model which has a ground truth sample for each generator input.
3. We show empirically that several architectural components are important for good image extension. We present ablation studies that show the effect of each of these components.

2. Related Work

Prior work in image inpainting can be fairly neatly divided into two subcategories: classical methods, which use non-parametric computer vision and texture synthesis approaches to address the problem, and learning-based methods, which attack the problem using parametric machine learning, generally in the form of deep convolutional neural networks. Classical methods, such as [6, 8, 13, 12] typically rely on patch similarity and diffusion to borrow information from the known regions of the image to fill in the hole. These methods work best when inpainting small holes in stationary textures and generally lack semantic understanding of the image. Perhaps the most successful of these methods are the Bidirectional Similarity [35] and PatchMatch algorithms [7, 23]. Other non-parametric approaches that specifically target image extension rely on

image patches from images other than the one to be extrapolated. [34, 41, 36] rely on having large databases of photos available during the extrapolation process, while others, such as [50], depend on a carefully selected guide image.

In recent years, deep learning based approaches have made great strides in overcoming the weaknesses of the classical methods. The first significant learning-based approach to inpainting was the Context Encoder [30], which trained an encoder-decoder model to fill in a central square hole in an image, using a combination of ℓ_2 regression on pixel values, and an adversarial loss [14]. [45] minimizes the difference of nearest neighbor activation patches in deep layers of a pretrained ImageNet classification network, for improved synthesis of highly textured content. [17] improve on the results of [30] by adding a local discriminator loss to the original global discriminator loss; the local discriminator focuses on the realism of the synthetic content, while the global discriminator encourages global semantic coherence. [48] improves on [17] further by introducing a coarse-to-fine approach. Their model has two chained encoder-decoder sections, the second of which contains a contextual attention layer, which learns the optimal locations in the unmasked regions from which the model should borrow texture patches. Other similar approaches include [26, 43, 47], while [46] train an unconditioned GAN to generate complete images from the target distribution and perform an inference-time optimization to search for the latent code that would produce the closest match to the known pixels of the masked image. The only previous fully-parametric approach to image extension that we are aware of is [40], which showed impressive results using an autoregressive model to extend 32x32 pixel images, including the ability to output multiple plausible completions. These are, however, too small for practical applications. Concurrent to our work is [44], which is similar to ours but does not condition the discriminator with pre-trained features.

3. Model

Our model uses a Wasserstein GAN framework [28] comprising a generator network that is trained with the assistance of a concurrently-trained discriminator network.

Our generator network, G has an input consisting of the image z with pixel values in the range $[-1, 1]$, which is to be extended, and a binary mask M . These are the same dimensions spatially and are concatenated channel-wise. Both z and M consist of a region of known pixels and a region of unknown pixels. In contrast to inpainting frameworks, the unknown region shares a boundary with the known region on only *one* side. z is set to 0 in the unknown region, while M is set to 1 in the unknown region and 0 in the known region. At training time,

$$z = x \odot (1 - M) \tag{1}$$

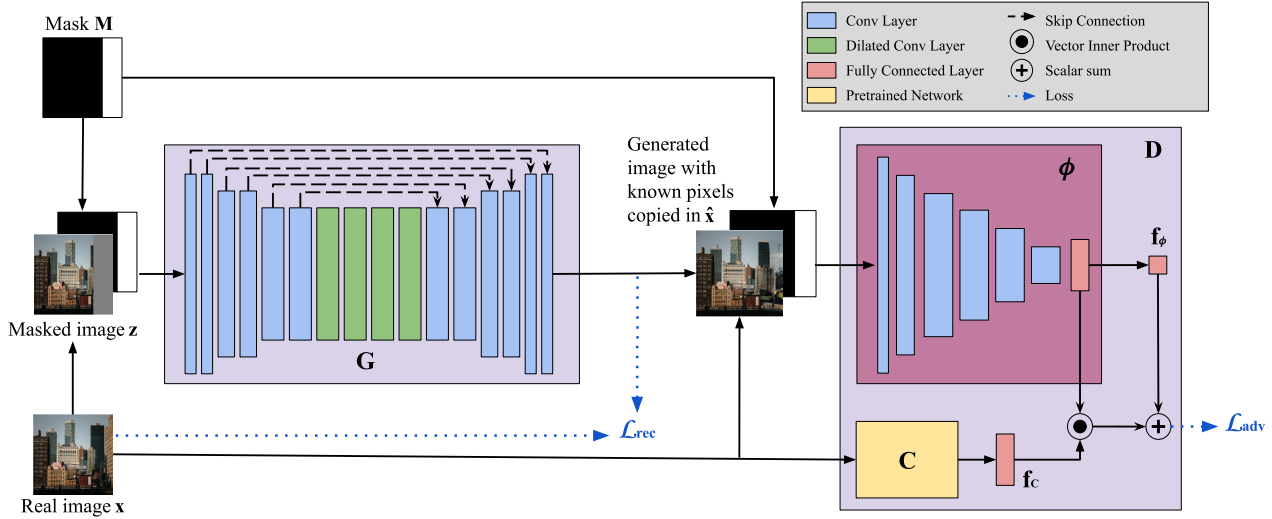


Figure 2. Model Architecture: this architecture is used for all our models. See text for further details.

where x is sampled from a natural image distribution \mathcal{X} and \odot is the element-wise multiplication operator.

The output $G(z, M)$ of G has the same dimensions as z and a pixel loss during training uses this full output. However, the last stage before feeding into the discriminator D is to replace what G synthesized in the unmasked regions with the known input pixels:

$$\hat{x} = G(z, M) \odot M + z \quad (2)$$

D is also a deep network, which transforms a real sample from \mathcal{X} or a generated sample \hat{x} to a single scalar value.

3.1. Generator

G generally follows the same fully convolutional encoder-decoder architecture as used by [47] (see Figure 2). Each layer in the generator except the last one uses gated convolutions [47] to enable the model to learn to select the contributing features for each spatial location and channel. Following the inpainting guidance in [48], each layer except the last uses an ELU activation function [10], and the final layer clips its outputs to the range $[-1, 1]$. As in [16, 47, 48], the innermost layers utilize dilated convolutions to increase their receptive field size.

To address the image extension problem, we deviated from the generator architecture proposed by [47] in a few crucial ways. We eliminated the refinement network, including the contextual attention layer, since this layer is biased towards copying patches from the unmasked portion of the input. While borrowing patches is a useful property for inpainting of images [7], in the case of image extension, it is less likely that repeated patterns will result in convincing extension. Figure 3 shows the effect of the contextual attention layer of [48, 47]. We also compare to

Adobe Photoshop’s PatchMatch-based [7] Content Aware Fill tool, which generates similar artifacts due to copying patches. These copying artifacts occur on a large fraction of the output images.

We also introduced skip connections [32] between the non-dilated layers, since we found that they improved the network’s ability to synthesize high frequency information. In Figure 6, we show the typical benefit of using skip connections. We additionally added instance normalization [39] after every generator layer besides the output layer, finding that it significantly reduces the number of artifacts in the generated images.



Figure 3. The contextual attention layer from DeepFill [48, 47] tends to repeat patches and structures. The original image (top left) is extended to the right (top-middle and top-right). DeepFill creates a copy of the door handle, whereas our extension extends the structure in a semantically and geometrically more plausible manner. Similarly, Photoshop’s Content Aware Fill (bottom row) often creates artifacts since it is based on PatchMatch [7].

3.2. Discriminator

The objective of the discriminator network (see Figure 2) is determining whether an image is generator-produced or real. In our problem setup, the concern is not just whether the output of G appears real, but also that it is a plausible extension of G 's inputs. To this end, we design our discriminator to be conditioned on the specific generator inputs when evaluating whether what is fed into the discriminator is real or fake. We condition the discriminator in two ways.

First, when a generated image is input, we copy the known pixels from z to overwrite the corresponding generated pixels, as described in eq. 2, and we additionally input the mask M itself. This on its own provides a major advantage to the discriminator in the adversarial game, since it can focus in on the area right around the seam at the edge of the real content and easily determine that an image is fake if there is any abrupt change in image statistics along that seam. We see this play out during training, as the generated image content close to the seam is the first to improve and the quality improvement gradually spreads towards the opposite edge of the image as training progresses. On its own, this form of conditioning produces seamless results, but the quality of generated content still deteriorates as it moves further from the real content.

To address this, we add another form of conditioning, which is a modified version of the conditional projection discriminator (cGAN) [29]. In the original cGAN paper, a one-hot class label y is passed into the discriminator in addition to the image x^* to be classified as real or fake. The discriminator output is then

$$D(x^*, y) = f_\phi(\phi(x^*)) + \langle \phi(x^*), f_y(y) \rangle \quad (3)$$

where ϕ is a learned function mapping an image to a vector, f_ϕ is a learned fully-connected layer that maps that vector to a scalar, f_y is a learned fully-connected layer mapping y to a vector of the same size as the output of ϕ , and $\langle \cdot, \cdot \rangle$ denotes an inner product. The cGAN paper shows that this parameterization of the GAN objective enables the model to simultaneously learn the distributions $p(x)$ and $p(y|x)$.

In our setting we don't necessarily have class labels available, and we also want our conditioning vectors to contain more information than class labels would provide. To this end, we were inspired by previous work on perceptual metrics [18, 49] to replace y with the activations of a pretrained image classification network, C , when applied to x (the ground truth image). We chose to instantiate C as an InceptionV3 [37] network trained on ImageNet [11] with the final softmax removed. We found that it helps to normalize these activations by subtracting the mean activation over the dataset and then dividing the result by its ℓ_2 norm. Note that since the discriminator is only used during training, we can condition on the full unmasked image (x), which also means that these activations can be precomputed

before training. This conditioning encourages the generated content to semantically match the target image, which especially helps avoid semantic drift in larger extensions. Formally, we replace eq. 3 with

$$D(x^*, M, x) = f_\phi(\phi(x^*, M)) + \langle \phi(x^*, M), f_C(C(x)) \rangle \quad (4)$$

The architecture of ϕ is based on [47] and consists of six strided convolutional layers, followed by a fully connected layer. Each convolutional layer uses a leaky ReLU activation function [27] and all layers apply spectral normalization [28] to satisfy the Lipschitz constraints of Wasserstein GANs [4]. The output dimensions of ϕ and f_C are both 256.

3.3. Training

The model is trained via a combination of a reconstruction loss and an adversarial loss. The reconstruction loss optimizes for coarse image agreement and is implemented as an ℓ_1 loss imposed on the full output of G . The full equation is below:

$$\mathcal{L}_{rec} = \|x - G(z, M)\|_1 \quad (5)$$

For the adversarial loss, which refines the coarse prediction, we use a Wasserstein GAN hinge loss [25, 38]:

$$\mathcal{L}_{adv,D} = \mathbb{E}_{x \sim P_{\mathcal{X}}(x)} [\text{ReLU}(1 - D(x, M, x)) + \text{ReLU}(1 + D(\hat{x}, M, x))] \quad (6)$$

$$\mathcal{L}_{adv,G} = \mathbb{E}_{x \sim P_{\mathcal{X}}(x)} [-D(\hat{x}, M, x)]$$

where ReLU is the rectified linear unit function. The total loss on the generator is

$$\mathcal{L}_{total} = \mathcal{L}_{rec} + \lambda \mathcal{L}_{adv,G} \quad (7)$$

In all our experiments we set $\lambda = 10^{-2}$.

Our model is implemented in TensorFlow [1]. The generator and discriminator are trained jointly using the Adam optimizer [20] with parameters $\alpha = 10^{-4}$, $\beta_1 = 0.5$, $\beta_2 = 0.9$; (the discriminator has a slightly larger $\alpha = 10^{-3}$, but other parameters are the same). Unlike many previous papers, we did not see improvement from training the discriminator for multiple steps per each generator step. Based on the findings of [9] on the benefits of training GANs with large batch sizes, we trained on 32 cores of a Google Cloud TPUv3 Pod with batch size 256.

4. Experimental Results

For all experiments we train our model on a dataset composed of the top 50 classes (measured by number of samples in the training set) of the Places365-Challenge dataset [51], producing a training set of just under 2 million images, which we scaled to the spatial dimensions of 257x257

	FID (25%)	PSNR (25%)	FID (50%)	PSNR (50%)	FID (75%)	PSNR (75%)	FID (Inp)	PSNR (Inp)
DF	1.87	7.11	11.65	6.69	31.21	9.74	4.96	14.31
PC	1.40	11.10	11.20	6.63	31.83	8.94	3.70	13.78
NCnd	0.85	8.96	5.01	7.55	19.17	9.08	2.73	14.24
Ours	0.79	10.17	3.46	8.63	8.79	8.07	2.53	14.17

Table 1. Quantitative metrics on 500 test images. The mask types are: 25% extension (3:1 ratio of context to mask), 50% extension (1:1 ratio), 75% (1:3 ratio) and inpainting a central square mask comprising 25% of image pixels (Inp). We compare DeepFill (DF), PartialConv (PC), ours without conditioning (NCnd), and our model.

	FID (25%)	PSNR (25%)	FID (50%)	PSNR (50%)	FID (75%)	PSNR (75%)
Prcptl	0.40	9.95	2.32	8.31	14.15	9.65
FM	0.75	9.42	3.14	8.97	14.74	8.87
Ours	0.79	10.17	3.46	8.63	8.79	8.07

Table 2. Comparisons with other methods for stabilizing GAN training. We provide PSNR for reference, but found that FID correlates with perceptual quality best. Based on FID on 25% and 50% extensions, feature matching and perceptual losses outperform our conditioning, but the difference is fairly small. On 75% extensions, our conditioning provides the best results and the difference is large.

pixels. The use of 50 classes allows us to test how well our model can generalize to multiple categories. We used a held-out set of 500 images from the same set of classes in the Places365 dataset, approximately 10 images per class, to compute quantitative scores and to visualize the image extension results.

4.1. Image Extension

We compare our model (which we call **Ours**) with various baselines both qualitatively and quantitatively on the task of image extension. Specifically, we evaluate each algorithm’s ability to fill in masked out image content for three different image extension tasks (where the rightmost 25%, 50% and 75% of pixels in the image are respectively masked) and one inpainting task (a central square mask comprising 25% of image pixels). For each experiment, our model is retrained with masks of the appropriate position, shape and size. The baselines we compare against are:

No-Cond: A model that is identical to “Ours,” but without discriminator conditioning.

DeepFill: Our re-implementation of DeepFillv2 [47], which is a state of the art inpainting model. We confirmed that our reimplementation achieves inpainting results nearly identical to that of the original papers (see Figure 6). For each experiment, we retrain the model with the same masks and data as ours. We follow the authors’ guidance of training for 5 days on an NVIDIA P100 GPU. We note that this results in the model being trained for many fewer

steps than “Ours” because it trains much more slowly (0.8 steps/sec vs 4.7 steps/sec for “Ours”). Comparison against this model shows the benefits of our approach compared to simply repurposing an architecture suited for inpainting tasks.

PConv: The authors of another state of the art inpainting work [26] generated results for us based on provided masks, but the models were not retrained specifically for these tasks. The model was trained on the full Places2 dataset, which is a superset of our training set. They use a database of free-form masks, some of which are very large (up to 50% of the image size), but are often non-contiguous and non-convex, which means that at training time the model may not have needed to generate pixels that were very far from known context. While our comparisons to this paper are not exactly apples-to-apples, we believe that this still provides a strong baseline against which to compare our performance.

CAF: Results from Adobe Photoshop’s content aware fill, which is based on the PatchMatch algorithm [7]. Content aware fill is a very powerful tool used for image extension, and represents a strong classical baseline. However, due to the use of only patch level information, it does not provide semantically meaningful extensions.

We provide quantitative performance metrics for each mask-type and each algorithm in Table 1. We agree with the authors of [48, 26] that there are really no good metrics that capture the goals of these experiments, but we nonetheless report the best approximations. Specifically, we report Fréchet Inception Distance (FID) [15] on the full output image and PSNR of the masked regions only. For FID we used a diagonal covariance matrix, due to having few samples. Based on our own qualitative evaluations, we feel that FID of the entire output image best correlates with what we perceive as quality image extension. We additionally performed a qualitative analysis. We show results on a few images from our test set in Figures 4 and 6. We show many more results, including on free-form masks, in the Supplementary Material. Overall we see that all methods, other than “CAF,” perform admirably for inpainting. On the extension tasks, as we move towards larger extensions with smaller context, “CAF” and “DeepFill” degrade into just repeating textures, while “PConv” gets blurrier and more artifact-filled the further away from the context it gets. The “NoCond” version of our model maintains higher quality for the larger extensions but does show some blurring and semantic drift. Meanwhile, our full model remains semantically consistent, with mostly photorealistic and seamless synthesis.

Furthermore, we experiment with replacing our conditioning with perceptual [18] and feature matching [42] losses, see Table 2 and Figure 5. In the perceptual loss,

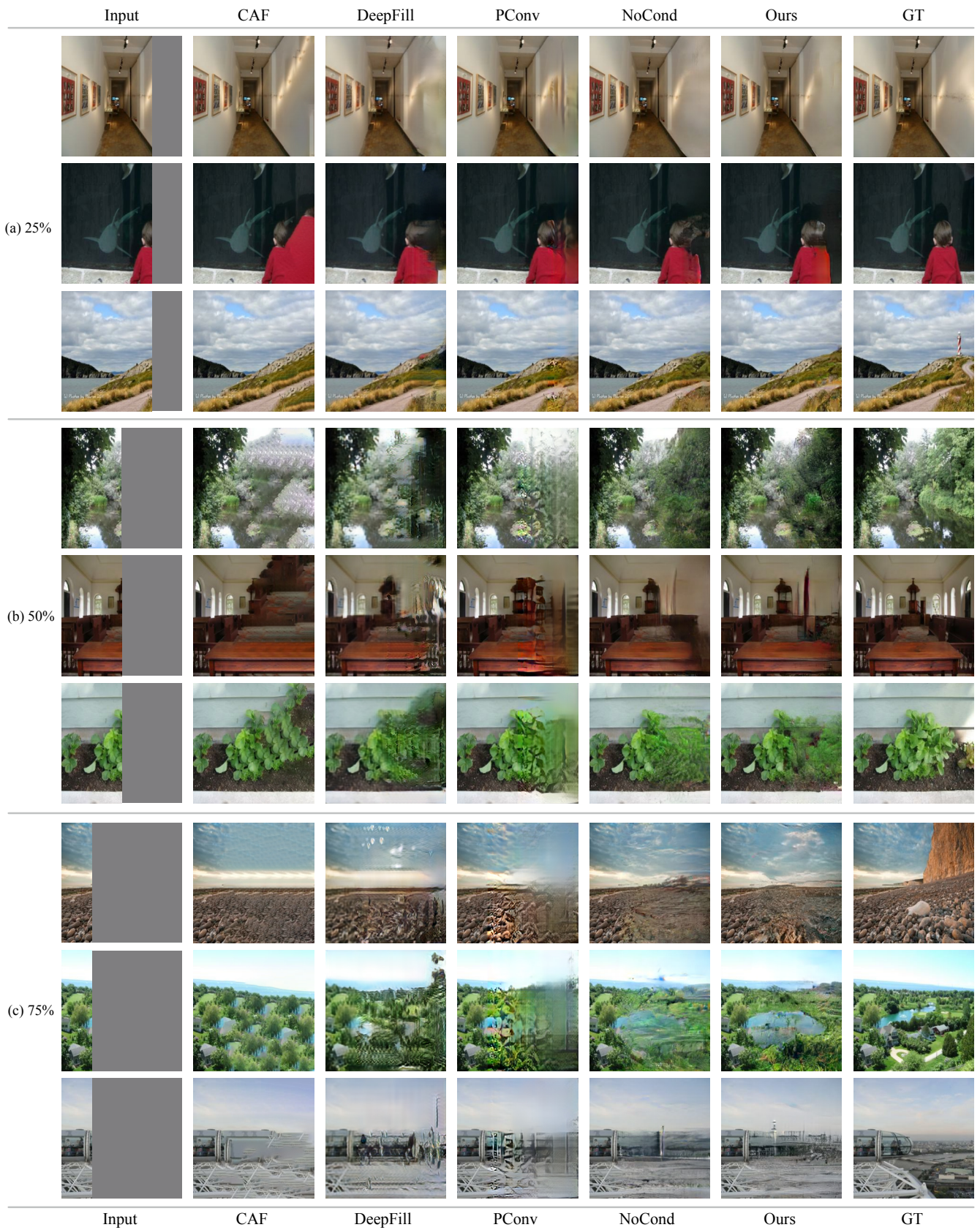


Figure 4. Extending images from masks of (a) 25%, (b) 50% and (c) 75% of the image width using multiple algorithms. From left to right: DeepFill [48], PConv [26], Photoshop Content Aware Fill, our model with no conditioning, our full model and ground truth.

the generator is optimized to produce images that are close to the ground truth images in the activation space of a pre-trained classification network. Similar to our conditioning, perceptual loss uses a pre-trained network to guide the training towards plausible extensions. Unlike our conditioning, the perceptual loss modifies only the generator objective to bias it towards semantic coherence, whereas our conditioning figures into both the generator and discriminator objectives and adding semantic information to the whole adversarial optimization game. Feature matching is similar in principle to a perceptual loss, but it minimizes the distance between activations in a hidden layer of the discriminator, rather than a pre-trained classification network, and similarly only figures into the generator objective.

In our experiments on 75% extensions, we found that our conditioning performs significantly better than feature matching and perceptual loss, while on smaller extensions they perform similarly. Our perceptual loss implementation tries to match pre-softmax logits, while our feature matching follows [42] and tries to match convolutional feature maps in the discriminator at multiple scales. Preliminary experiments indicate that combinations of all the three result in even better results, chiefly with fewer GAN-style artifacts, but we leave a more thorough analysis to future work.

We also qualitatively compare against PixelCNN in Figure 5. It is clear PixelCNN performs much worse than our method. Furthermore, it takes about 12 minutes to do inference for a single 64x64x3 image, on a Tesla P100. On our higher resolution 256x256 images, this would translate to over 2.5 hours. In contrast, our method takes 0.1 seconds/image.

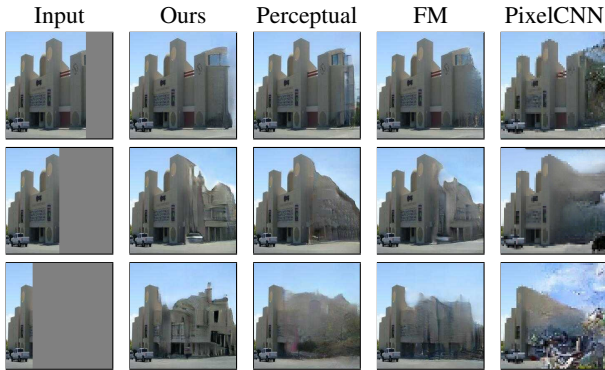


Figure 5. Although FM and Perceptual give slightly better results on short range extensions, our model far outperforms the others in the 75% case.

4.2. Ablations

In order to validate the contributions of each aspect of our model, we trained models each ablating a single feature of our network. We experimented with removing skip

connections in the generator, removing feature conditioning from the discriminator, removing instance norm from the generator and reducing batch size to 64.

We show an example of the results in Figure 6. We see that without skip connections, the model has a hard time synthesizing high frequency textures, which often results in a perceptible seam between the real and generated content. Without discriminator conditioning the quality degrades more rapidly as we move further away from the edge with known context. Removing instance norm causes an increase in white and over-saturated artifacts. The smaller batch model generally produces blurrier results.

4.3. Panorama Generation

We evaluate our model on its ability to generate panoramic images from a much narrower seed image. We use the same dataset and models as in Section 4.1, but whereas in those experiments we ran each model’s forward pass only once, this time we apply the model recursively, using the synthesized content from the previous step as the known content for the current step in a sliding window approach. More specifically, we take a 257x192 image, pad it with 65 columns of zeros, and extend it using our our model trained for 25% extension. We then extend the image again by selecting the right most 192 pixels of the image, padding with zeros, and feeding it to the extension model. This is repeated 5 times, ultimately extending the image 2.7 times out to a width of 582 pixels. We show some results in Figure 7 and more results in the Supplementary Material. While the results are for the most part plausible and aesthetically pleasing, we do see some degradation and semantic drift as we move away from the original image.

4.4. Exploring the Space of Plausible Extensions

To visually explore the space of results generated by our extension model, we took sample videos from the YouTube8m dataset [2] and applied our model to extend both the right and left sides horizontally. Videos are a natural domain to test our model since consecutive frames are closely related to each other, and the small variations can generate interesting plausible outcomes. This allows us to verify that our model has not memorized a fixed completion for closely related images or collapses with small natural variations in the input. We scaled the videos to have a height of 257 pixels and respectively selected the right and leftmost 192 columns. We then padded each of them with 65 columns of zeros on the side to be extended. We ran the resulting images, each frame independently, through our model trained with 25% masks. We then took the model output for the extended region and concatenated it back on the original video frame. Please refer to the supplementary video https://drive.google.com/open?id=1x6FCYPmoqSuCdeLJTD0UpQ_MQhBPv7_e.

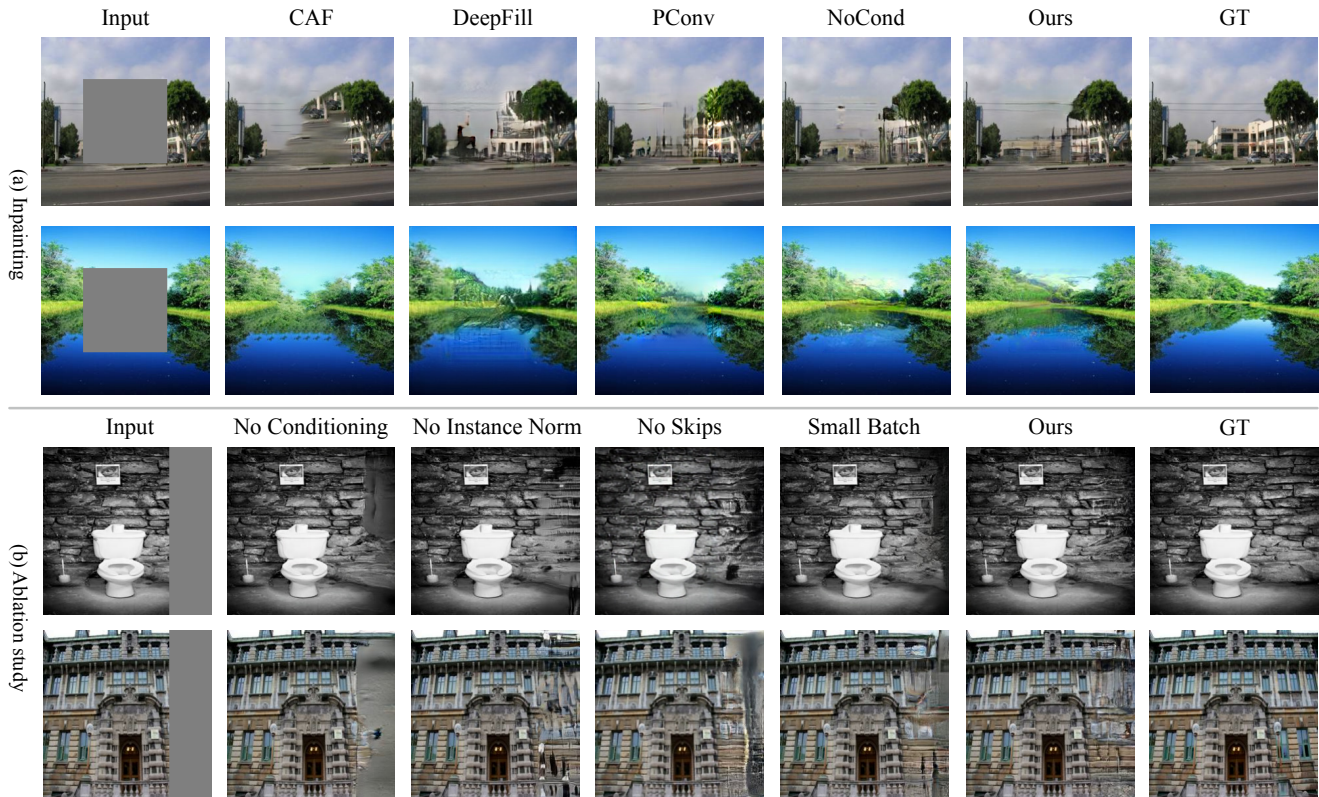


Figure 6. Further analysis: (a) Comparing the different models on inpainting tasks; our conditioned model performs on par with the state of the art models such as PConv and DeepFill for the inpainting problem. (b) Ablation tests: we remove (from second column to fifth column) only one of the following: discriminator conditioning, instance norm, skip connections, and reduce batch size.



Figure 7. Our models can also be used generate image panoramas. This can be viewed as a stress test for image extension tasks. We recursively apply the 25% model to create a very large output image of about 3 times the original width.

5. Acknowledgements

The authors would like to acknowledge Amol Kapoor and Dr. Huiwen Chang for helpful discussion and sharing code, and Dr. Guilin Liu for helping with running Partial Convolution comparison experiments.

References

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*, pages 265–283, 2016. 4
- [2] Sami Abu-El-Haija, Nisarg Kothari, Joonseok Lee, Paul Natsev, George Toderici, Balakrishnan Varadarajan, and Sudheendra Vijayanarasimhan. Youtube-8m: A large-scale video classification benchmark. *arXiv preprint arXiv:1609.08675*, 2016. 7, 12
- [3] Aseem Agarwala, Ke Colin Zheng, Chris Pal, Maneesh Agrawala, Michael Cohen, Brian Curless, David Salesin, and Richard Szeliski. Panoramic video textures. In *ACM Transactions on Graphics (TOG)*, volume 24, pages 821–

827. ACM, 2005. 1
- [4] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International Conference on Machine Learning*, pages 214–223, 2017. 4
- [5] Shai Avidan and Ariel Shamir. Seam carving for content-aware image resizing. In *ACM Transactions on graphics (TOG)*, volume 26, page 10. ACM, 2007. 1
- [6] Coloma Ballester, M Bertalmio, V Caselles, Guillermo Sapiro, and Joan Verdera. Filling-in by joint interpolation of vector fields and gray levels. *IEEE TRANSACTIONS ON IMAGE PROCESSING*, 10(8), 2001. 2
- [7] Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. Patchmatch: A randomized correspondence algorithm for structural image editing. In *ACM Transactions on Graphics (ToG)*, volume 28, page 24. ACM, 2009. 1, 2, 3, 5
- [8] Marcelo Bertalmio, Guillermo Sapiro, Vincent Caselles, and Coloma Ballester. Image inpainting. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pages 417–424. ACM Press/Addison-Wesley Publishing Co., 2000. 2
- [9] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. In *International Conference on Learning Representations*, 2019. 4
- [10] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). In *International Conference on Learning Representations*, 2016. 3, 11
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 4, 11
- [12] Alexei A Efros and William T Freeman. Image quilting for texture synthesis and transfer. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 341–346. ACM, 2001. 2
- [13] Alexei A Efros and Thomas K Leung. Texture synthesis by non-parametric sampling. In *Proceedings of the seventh IEEE international conference on computer vision*, volume 2, pages 1033–1038. IEEE, 1999. 2
- [14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc., 2014. 2
- [15] 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 *Advances in Neural Information Processing Systems*, pages 6626–6637, 2017. 5
- [16] Satoshi Iizuka, Edgar Simo-Serra, and Hiroshi Ishikawa. Globally and locally consistent image completion. *ACM Transactions on Graphics (ToG)*, 36(4):107, 2017. 2, 3
- [17] Satoshi Iizuka, Edgar Simo-Serra, and Hiroshi Ishikawa. Globally and Locally Consistent Image Completion. *ACM Transactions on Graphics (Proc. of SIGGRAPH 2017)*, 36(4):107:1–107:14, 2017. 2
- [18] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016. 4, 5
- [19] Biliana Kaneva, Josef Sivic, Antonio Torralba, Shai Avidan, and William T Freeman. Infinite images: Creating and exploring a large photorealistic virtual space. *Proceedings of the IEEE*, 98(8):1391–1407, 2010. 1
- [20] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015. 4
- [21] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations*, 2014. 2
- [22] Mike Knee and Roberta Piroddi. Aspect ratio problems in television today and some new solutions. *SMPTE Motion Imaging Journal*, 119(1):35–41, 2010. 1
- [23] Johannes Kopf, Wolf Kienzle, Steven Drucker, and Sing Bing Kang. Quality prediction for image completion. *ACM Transactions on Graphics (TOG)*, 31(6):131, 2012. 1, 2
- [24] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4681–4690, 2017. 2
- [25] Jae Hyun Lim and Jong Chul Ye. Geometric gan. *arXiv preprint arXiv:1705.02894*, 2017. 4
- [26] Guilin Liu, Fitsum A. Reda, Kevin J. Shih, Ting-Chun Wang, Andrew Tao, and Bryan Catanzaro. Image inpainting for irregular holes using partial convolutions. In *The European Conference on Computer Vision (ECCV)*, 2018. 1, 2, 5, 6
- [27] Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml*, volume 30, page 3, 2013. 4, 11
- [28] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. In *International Conference on Learning Representations*, 2018. 2, 4, 11
- [29] Takeru Miyato and Masanori Koyama. cGANs with projection discriminator. In *International Conference on Learning Representations*, 2018. 2, 4, 11
- [30] Deepak Pathak, Philipp Krähenbühl, Jeff Donahue, Trevor Darrell, and Alexei A. Efros. Context encoders: Feature learning by inpainting. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2536–2544, 2016. 2
- [31] Patrick Pérez, Michel Gangnet, and Andrew Blake. Poisson image editing. *ACM Transactions on graphics (TOG)*, 22(3):313–318, 2003. 1

- [32] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 3
- [33] Michael Rubinstein, Ariel Shamir, and Shai Avidan. Improved seam carving for video retargeting. In *ACM transactions on graphics (TOG)*, volume 27, page 16. ACM, 2008. 1
- [34] Qi Shan, Brian Curless, Yasutaka Furukawa, Carlos Hernandez, and Steven M Seitz. Photo uncrop. In *European Conference on Computer Vision*, pages 16–31. Springer, 2014. 2
- [35] Denis Simakov, Yaron Caspi, Eli Shechtman, and Michal Irani. Summarizing visual data using bidirectional similarity. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2008. 2
- [36] Josef Sivic, Biliana Kaneva, Antonio Torralba, Shai Avidan, and William T Freeman. Creating and exploring a large photorealistic virtual space. In *2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, pages 1–8. IEEE, 2008. 2
- [37] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016. 4, 11
- [38] Dustin Tran, Rajesh Ranganath, and David Blei. Hierarchical implicit models and likelihood-free variational inference. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5523–5533. Curran Associates, Inc., 2017. 4
- [39] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6924–6932, 2017. 3
- [40] Aäron Van Den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning-Volume 48*, pages 1747–1756. JMLR. org, 2016. 2
- [41] Miao Wang, Yukun Lai, Yuan Liang, Ralph Robert Martin, and Shi-Min Hu. Biggerpicture: data-driven image extrapolation using graph matching. *ACM Transactions on Graphics*, 33(6), 2014. 1, 2
- [42] 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. 5, 7
- [43] Yi Wang, Xin Tao, Xiaojuan Qi, Xiaoyong Shen, and Jiaya Jia. Image inpainting via generative multi-column convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 329–338, 2018. 2
- [44] Yi Wang, Xin Tao, Xiaoyong Shen, and Jiaya Jia. Wide-context semantic image extrapolation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [45] Chao Yang, Xin Lu, Zhe Lin, Eli Shechtman, Oliver Wang, and Hao Li. High-resolution image inpainting using multi-scale neural patch synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6721–6729, 2017. 2
- [46] Raymond A. Yeh*, Chen Chen*, Teck Yian Lim, Schwing Alexander G., Mark Hasegawa-Johnson, and Minh N. Do. Semantic image inpainting with deep generative models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. * equal contribution. 2
- [47] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Free-form image inpainting with gated convolution. *arXiv preprint arXiv:1806.03589*, 2018. 2, 3, 4, 5
- [48] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Generative image inpainting with contextual attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5505–5514, 2018. 1, 2, 3, 5, 6, 12
- [49] 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*, pages 586–595, 2018. 4
- [50] Yinda Zhang, Jianxiong Xiao, James Hays, and Ping Tan. Framebreak: Dramatic image extrapolation by guided shift-maps. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1171–1178, 2013. 1, 2
- [51] 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*, 2017. 4, 11

Boundless: Generative Adversarial Networks for Image Extension - Supplementary Material

6. Network Training and Architecture Details

6.1. Generator Network G

Layer ID	Type	Act.	K	S	D	Out	Skip
1	Gated Conv	ELU[10]	5	1	1	32	None
2	Gated Conv	ELU	3	2	1	64	None
3	Gated Conv	ELU	3	1	1	64	None
4	Gated Conv	ELU	3	2	1	128	None
5	Gated Conv	ELU	3	1	1	128	None
6	Gated Conv	ELU	3	1	1	128	None
7	Gated Conv	ELU	3	1	2	128	None
8	Gated Conv	ELU	3	1	4	128	None
9	Gated Conv	ELU	3	1	8	128	None
10	Gated Conv	ELU	3	1	16	128	None
11	Gated Conv	ELU	3	1	1	128	5
12	Gated Conv	ELU	3	1	1	128	4
13	Resize (2x)	n/a	n/a	n/a	n/a	n/a	n/a
14	Gated Conv	ELU	3	1	1	64	3
15	Gated Conv	ELU	3	1	1	64	2
16	Resize (2x)	n/a	n/a	1	n/a	n/a	n/a
17	Gated Conv	ELU	3	1	1	32	1
18	Gated Conv	ELU	3	1	1	16	None
19	Conv	None	3	1	1	3	None
20	Clip	n/a	n/a	n/a	n/a	n/a	n/a

Table 3. The generator architecture. Act. stands for activation type, K stands for kernel size, S for stride, D for dilation, Out for number of channels in convolutional layers and number of units in fully connected units, and Skip represents the layer-id which is concatenated into the output of the given layer. All resize operations use bilinear interpolation. In the Generator, all convolutional layers use 'Same' padding.

6.2. Discriminator Network D

The discriminator applies spectral normalization [28] at all layers, and consists of the the common tower (D_N , Table 4), which feeds into the non-conditional branch (f_N , Table 5) and projection discriminator branch (f_C , Table 6). These two branches produce scalars, which are then summed to produce a single network output. We invite the reader to see Section 3 of the main paper for more in depth discussion of the model.

The scalar outputs of the main and projection discriminator are summed and passed to the adversarial loss.

Common Tower D_N

Layer ID	Type	Act.	K	S	Padding	Out Size
1	Conv	LeakyReLU[27]	5	2	Same	64
2	Conv	LeakyReLU	5	2	Same	128
3	Conv	LeakyReLU	5	2	Same	256
4	Conv	LeakyReLU	5	2	Same	256
5	Conv	LeakyReLU	5	2	Same	256
6	Conv	LeakyReLU	5	2	Same	256
7	Conv	LeakyReLU	5	1	Valid	256
8	Flatten	n/a	n/a	n/a	n/a	n/a

Table 4. The base of the discriminator. It takes generated and ground truth images as input. Act. stands for activation type, K stands for kernel size, S for stride, Out for number of channels in convolutional layers and number of units in fully connected units.

Non-Conditional Branch f_N

Layer ID	Type	Act.	Out Size
1	Fully Connected No Bias	None	1

Table 5. The non-conditional branch of the discriminator, taking the common tower from Table 2 as input and outputting a single scalar value. Act. stands for activation type.

Projection Discriminator Branch f_C

Layer ID	Type	Act.	Out Size
1	Normalize	None	1000
2	Fully Connected No Bias	None	256
3	Inner Product w/Common Tower	None	1

Table 6. The projection discriminator [29] branch of the network. The input is logits of a pretrained classification network, for which we used an InceptionV3 [37] network trained on ImageNet [11]. The output is a single scalar, which is summed with the output of the non-conditional branch and passed to the hinge loss.

6.3. Training details:

We take the training set of Places365-Challenge dataset [51], select the top 50 classes by number of samples, and create a holdout validation set from this. This creates about 39,000 training and 930 test samples per class, for a total training set size of 1,953,624 and testing set size of 46376. The classes selected are:

- amusement park
- aquarium
- athletic field
- baseball field
- bathroom
- beach
- bridge
- building facade
- car interior
- church - indoor

- church - outdoor
- cliff
- coast
- corridor
- dining room
- embassy
- forest
- forest path
- golf course
- harbor
- highway
- industrial area
- lagoon
- lake
- lighthouse
- living room
- lobby
- mansion
- mountain
- ocean
- office building
- palace
- parking lot
- pier
- pond
- porch
- railroad track
- rainforest
- river
- skyscraper
- stadium
- staircase
- swamp
- swimming hole
- swimming pool
- train station
- underwater
- valley
- vegetable garden
- water park

Before passing the training image into the generator we resize the image to 257×257 , and also concatenate the mask channel. The mask size is randomly sampled from a uniform distribution, which is the target size plus/minus 4 pixels, so the model doesn't overfit to a specific mask size.

Following the code of DeepFill [48], we concatenate a channel of 1's to the input of the generator. This enables the

generator to see the edge of the image after 0 padding the inputs, although we do not verify this in this work.

We take generator and discriminator steps in a 1:1 ratio, with the steps executed jointly.

Please see Section 3 of the main paper for more discussion of the loss and optimizer.

7. Qualitative Results

We show additional samples from on the 25%, 50%, and 75% mask image extension experiments, and refer the reader to Figures 9, 10, and 11. We also show additional results from in-painting experiment in Figure 12 and more panorama results in Figure 13. We also demonstrate the suitability of our method on freeform masks in Figure 8.

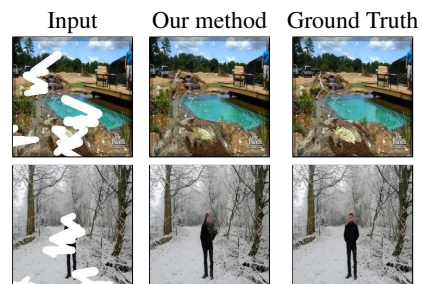


Figure 8. Results on freeform masks.

8. Exploring the Space of Plausible Extensions

We invite the reader to view the accompanying video derived from a sample from the YouTube8m dataset [2] at https://drive.google.com/file/d/1x6FCYPmoqSuCdeLJTD0UpQ_MQhBPv7_e/view?usp=sharing. Please refer to the main paper for details on how it was created. We encourage the reader to pause the video at arbitrary frames to see how the model produces different plausible completions as the result of tiny perturbations of the original frame.

9. Failure Cases

In Figure 14 we examine some of the failure modes of our image extension model. We note that our model is much better at textures than objects; for example vehicles, people, and furniture are challenging for the model. Addressing this is left to future work.

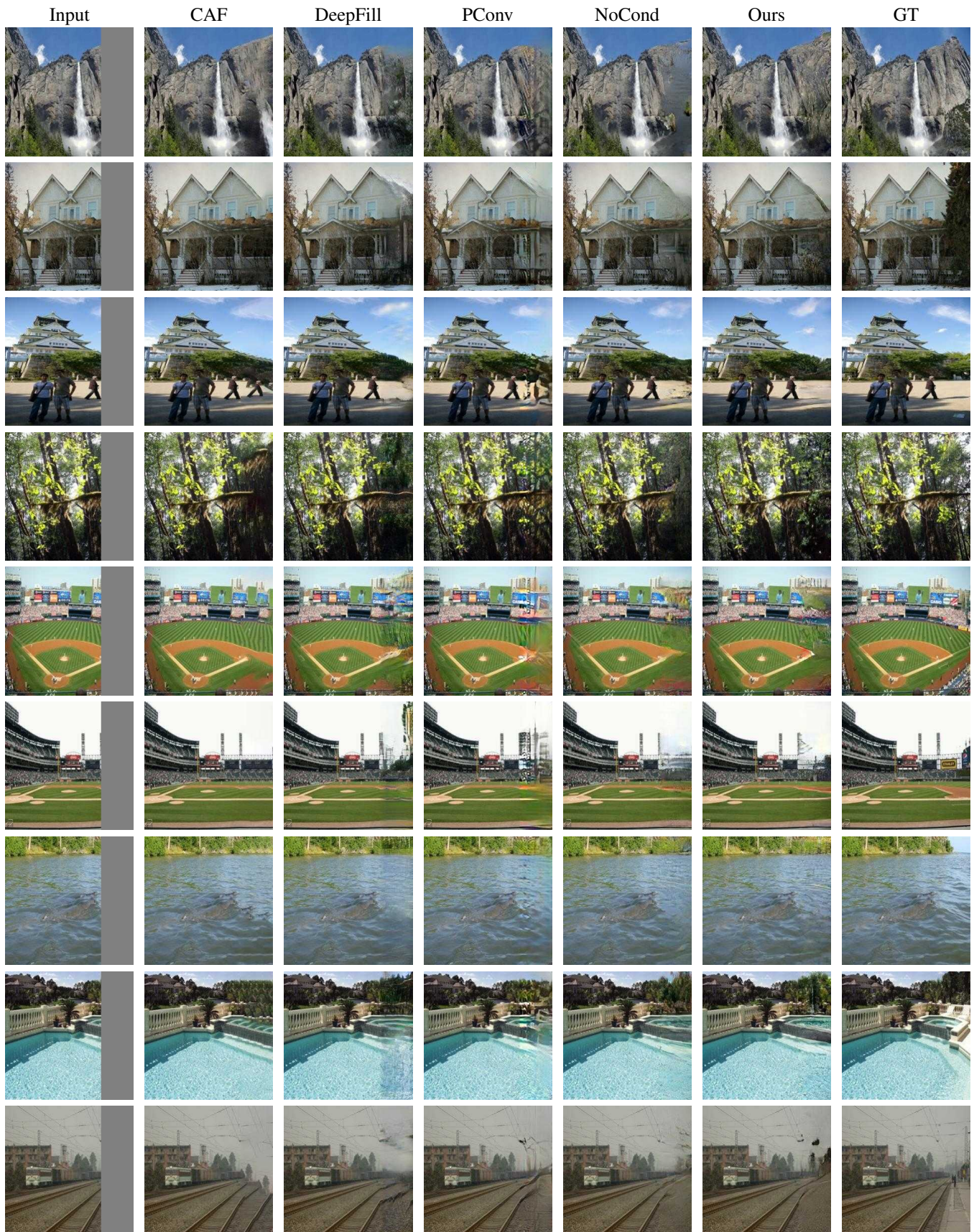


Figure 9. Extending images from masks which are 25% of the image width. We note that edges and structure are better defined in our method. For instance, edge of the roof in the second row.

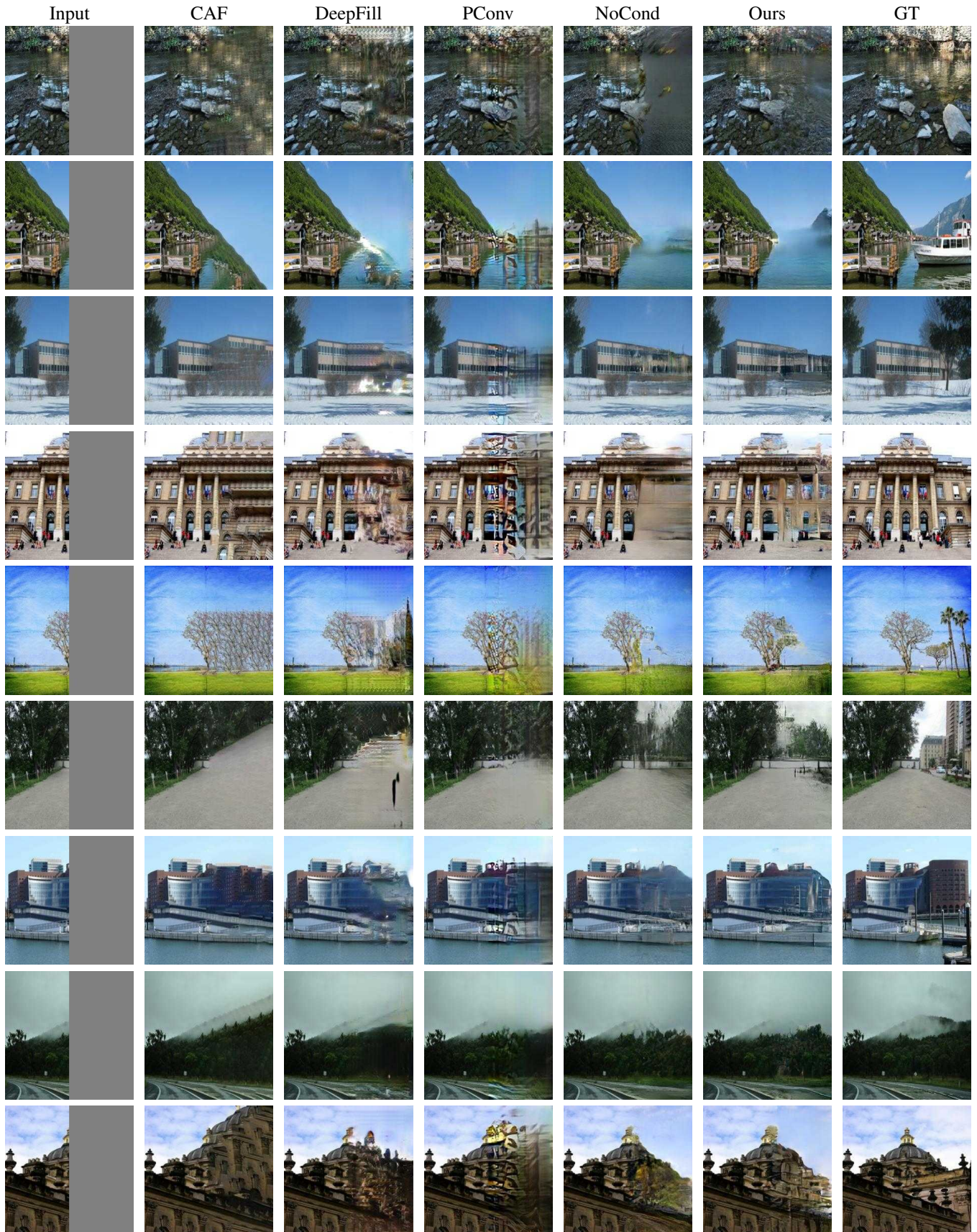


Figure 10. Extending images from masks which are 50% of the image width.

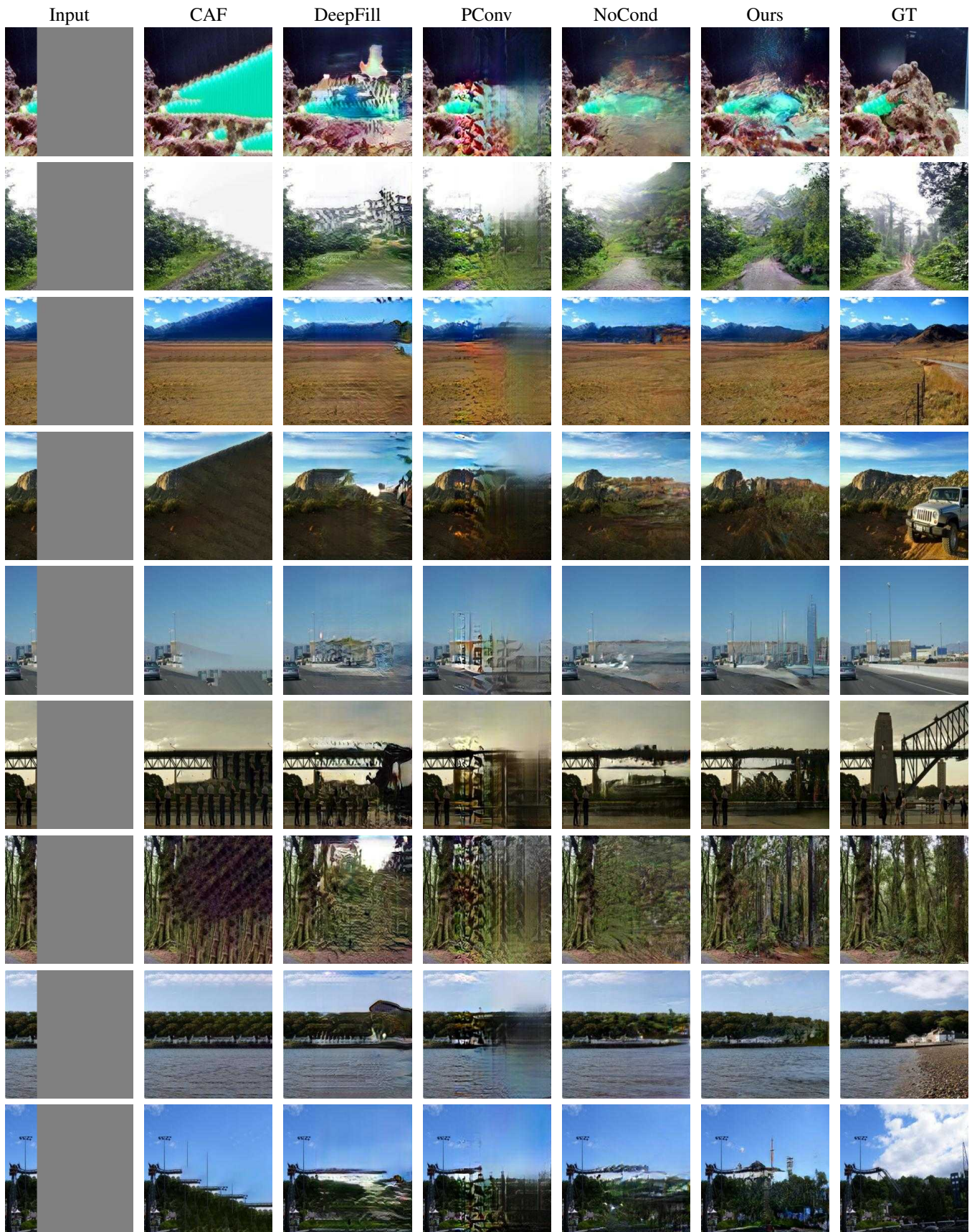


Figure 11. Extending images from masks which are 75% of the image width.

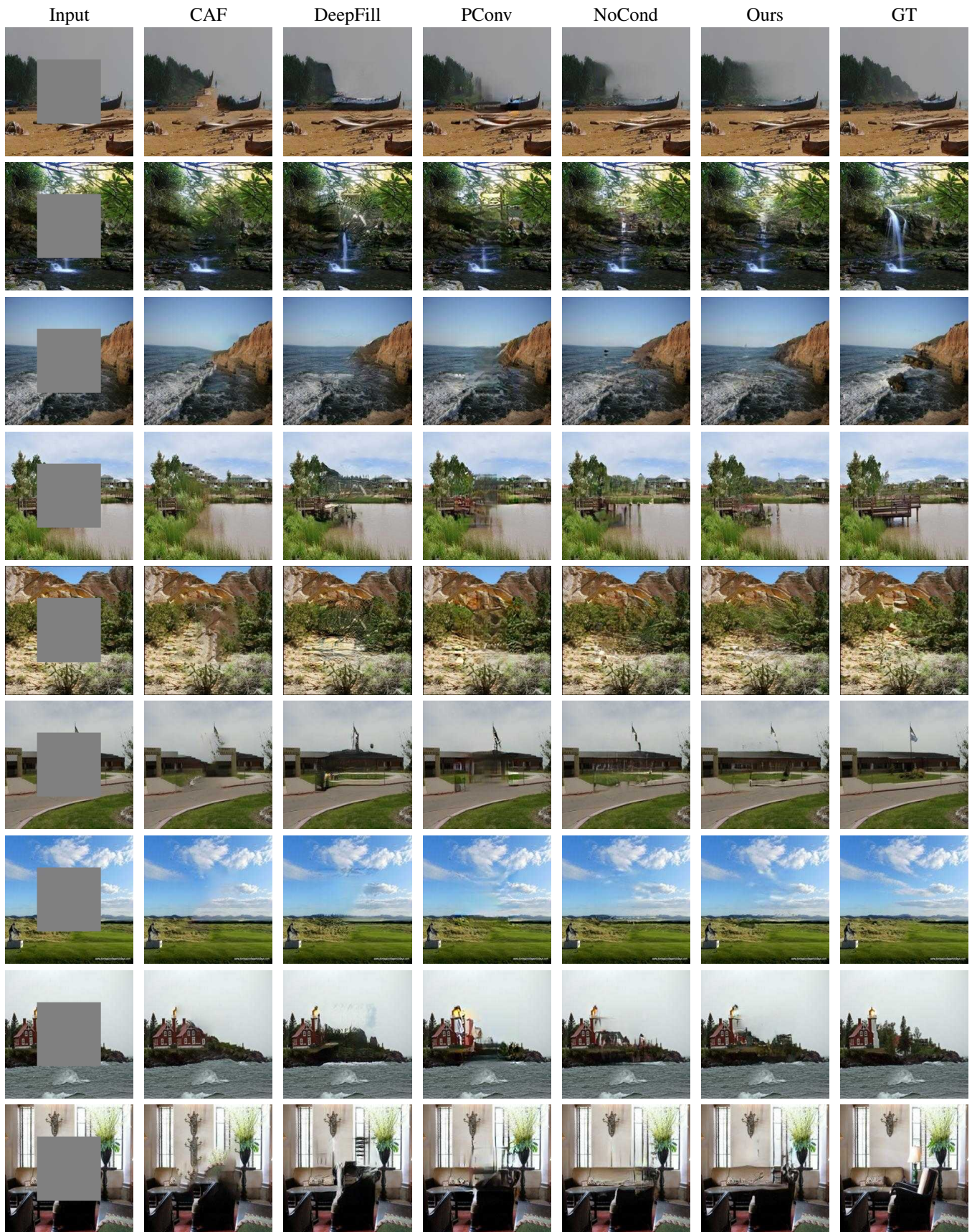


Figure 12. Center Inpainting.

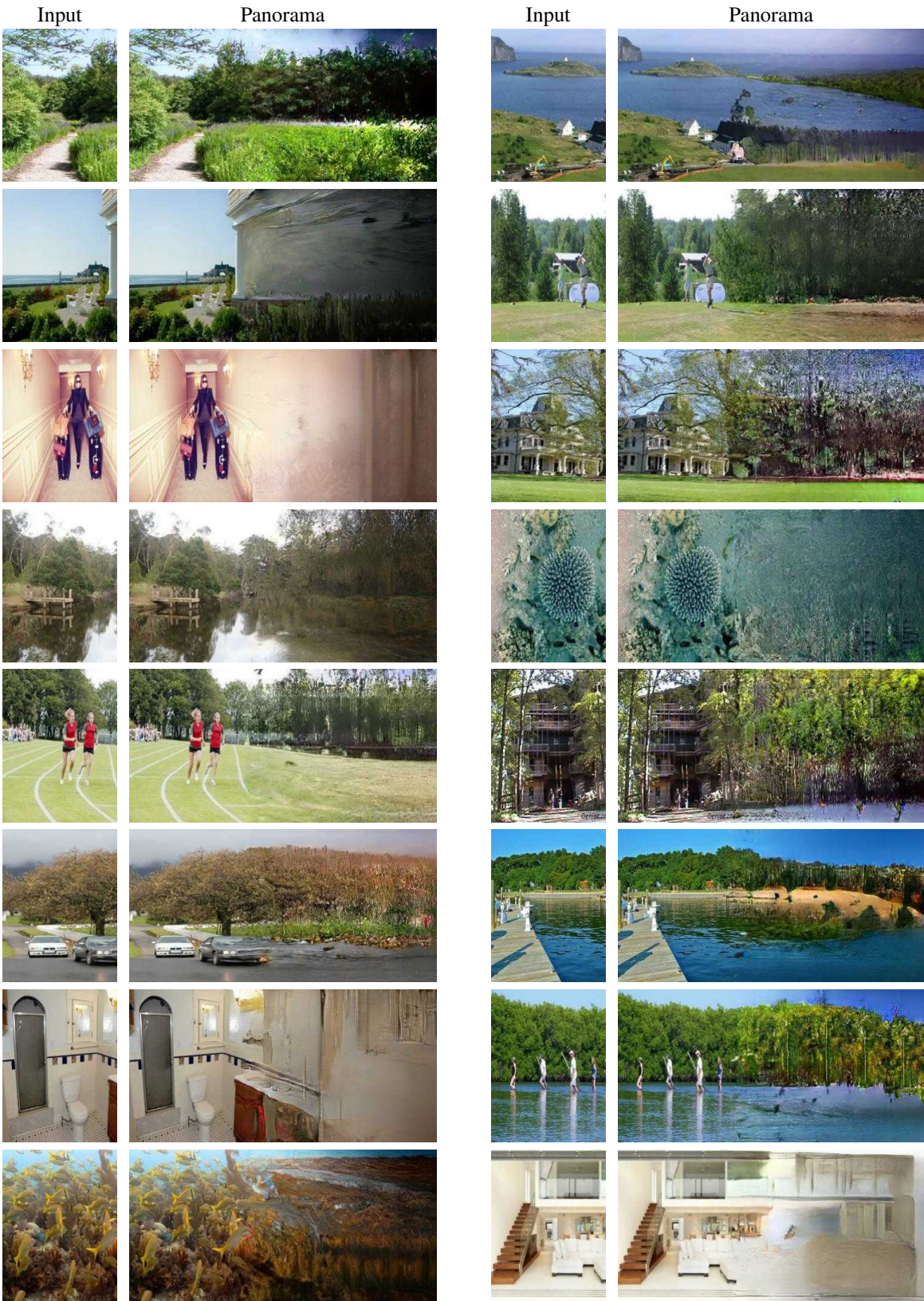


Figure 13. Additional panorama results

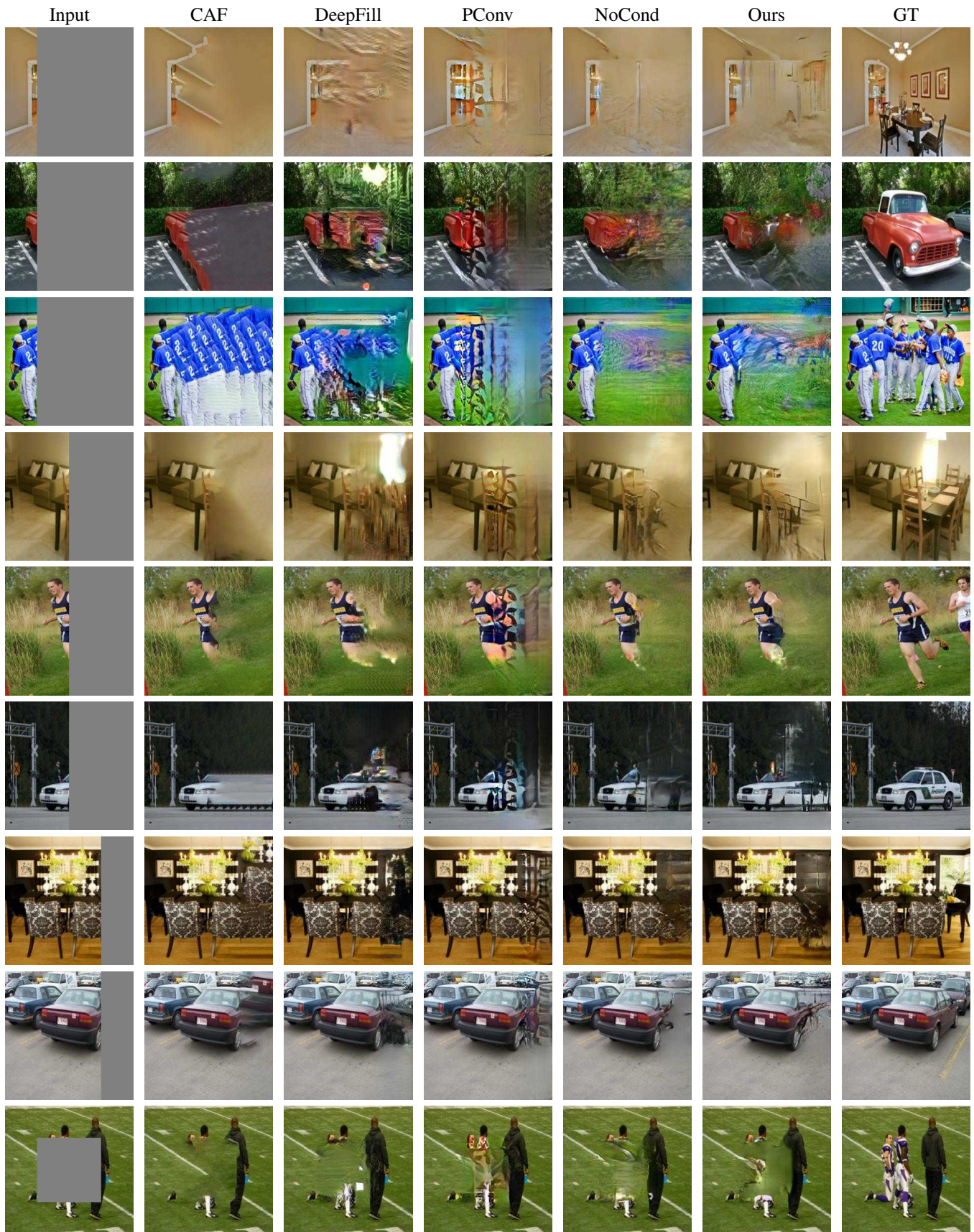


Figure 14. Failure cases. The network struggles with objects; especially cars, humans, and furniture.