

Image Inpainting

Jin Woo, Baik

Stanford University

jwbaik96@stanford.edu

Abstract

There have been some notable progress in image inpainting field by using learning-based methods. However, these methods often produce blurry or awkward results. In order to solve this problem, I have proposed three ideas. First, Coarse-to-refine structure [4] was adopted and Recurrent Convolution was used instead of vanilla convolution. Next, I have implemented Sobel Edge Discriminator, which is a discriminator using Sobel Edge detection algorithm [3]. Finally, frequency separation loss were added to the total loss. Since the high frequency components are less important in the coarse network, a low pass filter was applied to get low frequency components. In contrast, a high pass filter was applied in refine stage so the network can focus on details. From these proposed methods, my model was able to achieve better results than other methods.

1. Introduction

Image inpainting is a task which fills missing pixels with semantically and perceptually plausible contents. In order to produce appropriate reconstruction, it should be semantically plausible and matched with the existing backgrounds. This process could be easily done with deep learning techniques.

There are many deep learning based inpainting techniques. I have adopted coarse-to-refine structure [4] for my model. This structure fills the masked region through coarse network and refine network. The coarse network takes an original image and the mask as an input to generate coarse output. Then, the refine network takes coarse output as an input to generate refined output. This structure allows to expand the receptive field and stabilize the training stage.

Based on this coarse-to-refine structure, I have proposed 3 ideas; **Recurrent Convolution, Sobel Edge Discriminator, and frequency separation loss**. Recurrent Convolution (or RNN) is often used when the network deals with sequential data such as music, movie, and natural languages. My model can be treated as a sequential model because it is based on coarse-to-refine structure.

Sobel Edge Discriminator originates from Sobel Edge Detection [3]. The overall quality can be greatly improved by using the edge information. In my model, Sobel edge discriminator takes the edge map as input and determines whether the edge map is real or fake.

Frequency separation loss is designed to catch the structure of an image more precisely. Since the basic information of an image is concentrated in low frequency regions, coarse network focuses on low frequency components by eliminating the high frequency details. This method makes the network simple, while maintaining core information of an image. In contrast, refine network focuses on high frequency information(details).

The performance was greatly improved comparing to other models [5, 7, 8] by these techniques. Especially, my model shows the best performance at background images.

2. Related Work

There are some popular methods for inpainting task. **Globally and Locally Consistent Image completion** [2] has two discriminators: global discriminator and local discriminator. The global discriminator makes an output consistent with the global context and the local discriminator makes an output consistent with local regions which corresponds to the masked region. This model is composed of three networks: completion networks, local discriminator, and global discriminator. The completion network fills the masked region using CNN structure. Then, the global and local discriminator evaluate whether the output is consistent with the original image globally and locally.

Generative Image Inpainting with Contextual Attention [6] introduces the concept of ‘contextual attention layer’. First, the input feature is divided into two regions. One is called ‘Foreground’ which corresponds to the masked region and the other is called ‘Background’. The objective of this layer is to find a patch in background which matches with foreground. To achieve this, background should be transformed into convolutional filters with cosine similarity. Then, the convolution operation is applied to the foreground and convolutional filters. By applying the softmax operation to the output, attention scores for each

pixel can be achieved. Finally, the foreground region can be reconstructed by applying deconvolution to the attention score.

Free-Form Image Inpainting with Gated Convolution [7] solves the problem of vanilla convolution which treats all pixels as valid ones. While most of image inpainting networks used hard-gating mask, this model used soft mask. Gating corresponds to a soft mask which is updated automatically with convolution. By using this mechanism, the network can selectively learn masked region.

Image Fine-Grained Inpainting [1] introduced ‘dense multi-scale fusion network (DMFN)’. Unlike common dilated convolutions, DMFN uses hierarchical features which come from various kinds of convolution having different dilation rates. In addition, the author proposes a new concept called ‘self-guided regression loss’ which corrects low-level features by using normalized discrepancy map.

Recurrent Neural Network(RNN) is designed to deal with sequential data such as music, movies, and natural language. Since sequential data has different length, traditional neural network is not adequate. The main difference between RNN and traditional neural network is that traditional neural network only uses current input, while RNN uses past information by using hidden states.

Edgeconnect: Generative Image Inpainting with Adversarial Edge Learning [5] uses edge information when generating an image. It first draws an edge map of masked region. Then the network uses edge map to fill masked region. The network is composed of two stages. First stage is called edge generation stage. In this stage, the network generates edge map of the masked region. The next stage is called image completion stage. It fills the masked region by using the edge map generated at the first stage.

3. Methods

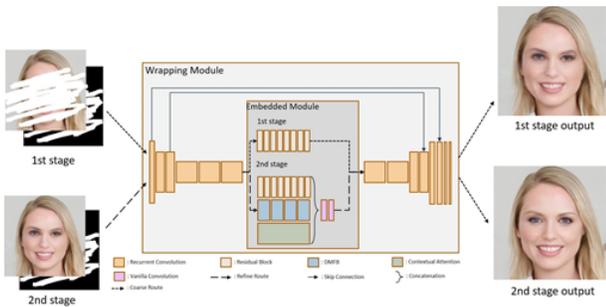


Figure 1. Overall flow of the network

Figure 1 shows the overall flow of my model. First, feed masked image and the mask into the network to get coarse output. Then, feed coarse output and the mask into the network to get final output.

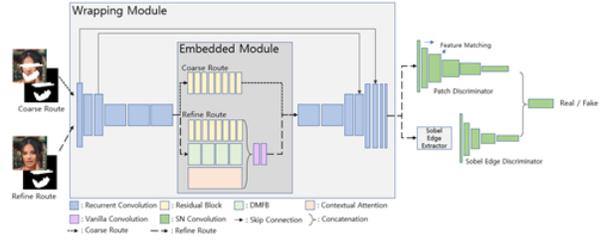


Figure 2. Overall structure of network

This is an overall structure of my model. There are two main modules: wrapping module and embedded module. After wrapping module, the result goes through patch discriminator and Sobel edge discriminator which determines whether the generated image is real or fake. The results from these discriminators are concatenated and make the final decision whether the generated image is real or fake.

Wrapping module is composed of encoder-decoder network. Inputs(original image and mask) are encoded using recurrent convolution and passes through embedded module. Then, the output of embedded module is decoded using recurrent convolution.

Embedded module is composed of coarse route and refine route. Coarse route is composed of residual block and refine route is composed of residual block, DMFB [1] and contextual attention [6] which the concatenation of these blocks goes through vanilla convolution. Further details are elaborated in **Appendix B**.

3.1. Recurrent Convolution

The concept of ‘recurrent convolution’ comes from Recurrent Neural Network(RNN). Coarse-refine structure is adopted in my model. One of the problem with this structure is that network becomes too long since coarse stage and refine stage are connected in series. From this reason, generator loss might not reach to the end of a stage. To handle this problem, I adopted RNN structure. Features that come from the coarse stage are reused at refine stage.

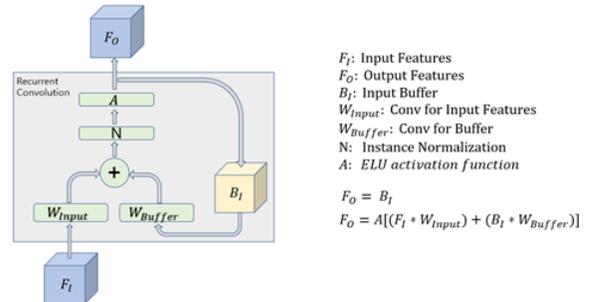


Figure 3. Recurrent convolution structure

Figure 3 shows the design of a recurrent convolution. Initially, buffer is set to zero. Convolution operation is applied to the inputs (input features and buffer) and their weights. Then, sum them up and apply instance normalization and activation function (ELU function) as followed. Finally, the output is achieved, and outcome is stored in buffer. By using this method, coarse and refine networks could be used simultaneously, while solving the gradient vanishing problem.

3.2. Sobel Edge Discriminator

Sobel Edge Discriminator is a discriminator which uses Sobel edge detection [3]. It uses Sobel edge extractor which role is to extract edge information of an image. Horizontal and vertical edge information could be achieved by using horizontal line detector and vertical line detector. From this information, the whole edge of an image could be achieved.

It looks similar to Edgeconnect [5] since both models use edge information. However, there is a difference between my model and Edgeconnect. Edgeconnect uses edge information in the second stage. It draws the edge first and puts it into subsequent network to get a complete image. So, if the edge map is awkward, the completed image will be awkward as well. However, in my model, the network learns edge information and output simultaneously.

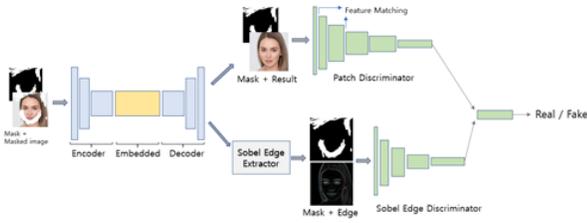


Figure 4. Sobel Edge Discriminator

Figure 4 shows how Sobel edge discriminator is applied. After the network produces an image, the result is fed into the patch discriminator and Sobel edge extractor. Then, the outputs of two discriminators are concatenated and determine whether the generated image is real or fake.

3.3. Frequency Separation Loss

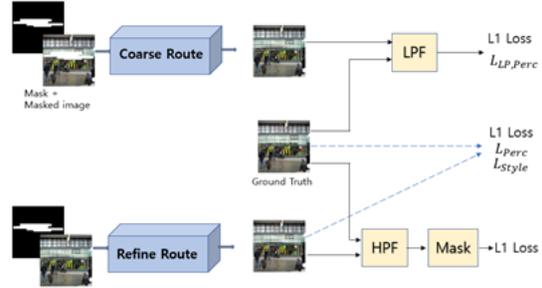


Figure 5. Sobel Edge Discriminator

I have also added frequency separation loss. High frequency details are less important in the coarse stage, so I used low pass filter to focus on low frequency information to catch basic information. In contrast, I used high pass filter at the refine stage to achieve high frequency details. In addition, the L1 loss between the ground truth and the reconstructed image was applied at the refine stage.

For the low pass filter, I did not select Gaussian filter, which is also called as mean filter, since it blocks high frequency components like edge. Instead, I chose the low pass filter above so that the coarse network can consider both high frequency and low frequency components.

4. Experiments

4.1. Dataset

The network was trained by CelebA HQ dataset for human face inpainting and Places365 dataset was used for place inpainting. For the evaluation, Google search face images were used for face inpainting and Places365 testset were used for place inpainting. I have compared the result with ‘Free-form image inpainting with gated convolution’ [7], ‘Edgeconnect: Generative image inpainting with adversarial edge learning’ [5], and ‘Pluralistic image completion’ [8].

4.2. Settings

Hinge loss is used for discriminator with learning rate of 0.004, which is commonly used in GAN.

$$L_D = \mathbb{E}_{x \sim P_{data}(x)} [ReLU(1 - D(x))] + \mathbb{E}_{z \sim P_z(z)} [ReLU(1 + D(G(z)))]$$

The total generator loss can be represented as below. The learning rate of generator is 0.001.

$$L_{Total} = 0.01L_G + 5L_1 + L_{lpf} + L_{hpf} + L_{LP, Perc} + L_{Perc} + 250L_{Style} + L_{FM}$$

The generator loss is composed of hinge loss, L1 loss, perceptual loss and style loss. These losses make the net-

work robust to variations. The details of these losses are shown as below.

$$\begin{aligned}
 L_G &= -\mathbb{E}_{z \sim P_Z(z)} [D(G(z))] \\
 L_{D1} &= \mathbb{E} \left[\sum_{i=1}^N \frac{1}{N} \sum_{t=1}^N |D_i(I_{gt}) - D_i(I_{pred.refine})| \right] \\
 L_{L1} &= \frac{1}{N} \sum_{i=1}^N |F_i(I_{gt}) - F_i(I_{pred.refine})| \\
 L_{hpf} &= \frac{1}{N} \sum_{i=1}^N M * |F_i(I_{gt}) - F_i(I_{pred.refine})| \\
 L_{l1} &= \frac{1}{N} \sum_{i=1}^N |I_{gt} - I_{pred.refine}| \\
 L_{perc} &= \mathbb{E} \left[\sum_t \|\Phi_i(I_{gt}) - \Phi_i(I_{pred.refine})\|_1 \right] \quad t = \text{relu1_1, relu2_1, relu3_1, relu4_1, relu5_1} \\
 G_j^{\Phi}(x)_{c,w} &= \frac{1}{C_H W_j} \sum_{h=1}^H \sum_{w=1}^W \Phi_j(x)_{h,w,c} \Phi_j(x)_{h,w,c} \quad j = \text{relu2_2, relu3_4, relu4_4, relu5_2} \\
 L_{style} &= \sum_j \|G_j^{\Phi}(I_{gt}) - G_j^{\Phi}(I_{pred.refine})\|_1
 \end{aligned}$$

For the optimizer, I used Adam optimizer with a learning rate scheduler of beta1 and beta2 as 0.5 and 0.9. For the learning, I have iterated 200,000 times and set the learning rate decay as 0.75 after 100,000 iterations.

5. Results

For the evaluation metrics, L1 error, L2 error, perceptual loss, PSNR, and SSIM were used. In addition, the masks were divided into 3 parts by mask rates from 10~20%, 20~30%, and 30~40%.

5.1. Face inpainting

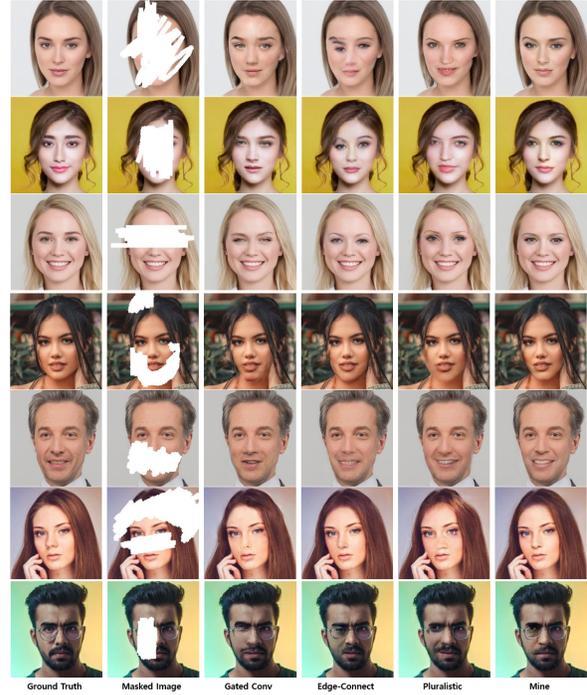


Figure 6. Qualitative comparison

	Mask	Gated	Edge	Plural	Mine
11 (%)	10-20%	4.95	5.79	4.11	5.76
	20-30%	6.54	6.70	6.32	6.54
	30-40%	9.86	9.19	10.12	8.76
12 (%)	10-20%	1.47	1.32	1.47	1.28
	20-30%	2.03	1.54	2.72	1.43
	30-40%	4.38	2.79	4.81	2.58
PL	10-20%	0.431	0.554	0.358	0.542
	20-30%	0.592	0.632	0.512	0.604
	30-40%	0.813	0.794	0.739	0.762
PSNR	10-20%	29.575	29.803	29.693	29.920
	20-30%	27.853	28.748	27.847	29.152
	30-40%	24.026	25.888	24.468	26.012
SSIM	10-20%	0.933	0.917	0.941	0.920
	20-30%	0.904	0.898	0.912	0.905
	30-40%	0.849	0.854	0.854	0.863

Table 1. Quantitative Comparison

My model performs best at 30~40% mask rate. My model shows the best performance for all categories except for the perceptual loss. And for the 10~20% and 20~30% rates, my model shows best performance at L2 error and PSNR.

5.2. Background inpainting

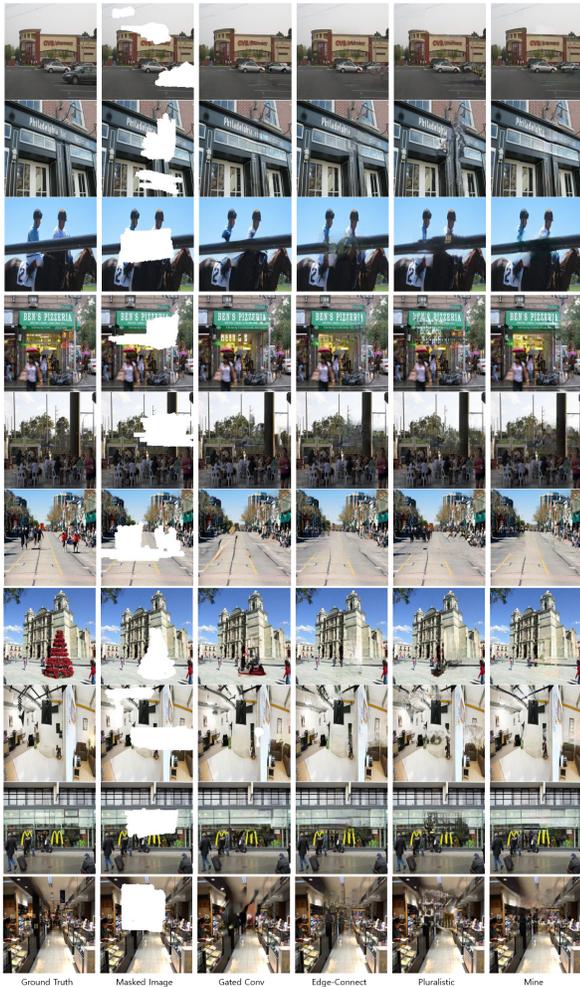


Figure 7. Qualitative comparison

For Places365 dataset, my model outperforms other models in all categories.

	Mask	Gated	Edge	Plural	Mine
I1 (%)	10-20%	14.25	12.70	14.60	11.82
	20-30%	16.47	15.52	18.89	14.02
	30-40%	28.24	27.55	32.087	24.840
I2 (%)	10-20%	12.25	10.03	10.53	9.99
	20-30%	12.80	11.17	13.89	10.53
	30-40%	25.62	22.06	26.217	19.999
PL	10-20%	0.682	0.592	0.718	0.546
	20-30%	0.927	0.847	1.001	0.775
	30-40%	1.311	1.205	1.413	1.118
PSNR	10-20%	23.311	24.112	22.835	24.237
	20-30%	21.617	22.270	21.313	22.572
	30-40%	19.114	20.104	18.890	20.308
SSIM	10-20%	0.899	0.898	0.880	0.907
	20-30%	0.858	0.856	0.832	0.871
	30-40%	0.769	0.768	0.729	0.789

Table 2. Quantitative Comparison

6. Conclusion

As shown in the **5. Results 5**, my model shows better performance than other methods. This performance can be attributed to three main ideas. First, recurrent convolution allows the network size to be reduced, so it can prevent gradient vanishing problem. Next, the network can simultaneously learn edge information and image generation from Sobel edge discriminator. Finally, frequency separation loss makes the coarse network to focus on low frequency components and the refine network to focus on high frequency components.

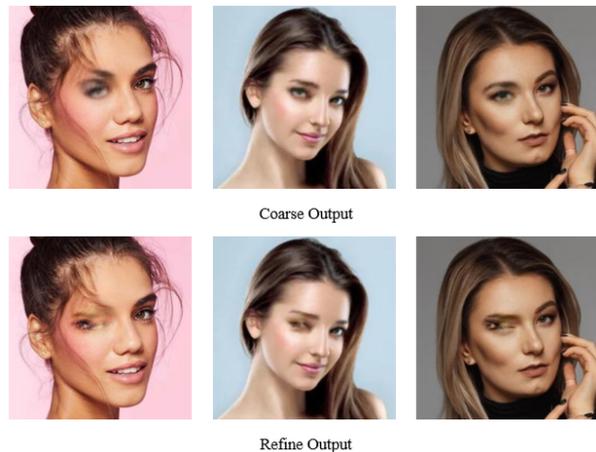


Figure 8. Examples of cases which refined output is worse than coarse output

Even though the model shows good performance, there is a limitation. There are some cases (Figure 8) which re-

finer output is worse than coarse output. For the future work, the refine network could be modified to prevent this issue.

Contributions/Acknowledgements

For comparing the qualitative result, Github codes from the baseline papers were used. The below are the links to the code.

Gated-convolution [7] https://github.com/JiahuiYu/generative_inpainting

Pluralistic [8] <https://github.com/lyndonzheng/Pluralistic-Inpainting>

Edgeconnect [5] <https://github.com/knazeri/edge-connect>

A. Code

You can find the entire code in <https://github.com/jwbaik96/Image-inpainting.git>

B. Network composition

B.1. Horizontal and vertical filters

```
filter_h = [[[[[1, 2, 1], [0, 0, 0], [-1, -2, -1]]]]]
filter_v = [[[[[-1, 0, 1], [-2, 0, 2], [-1, 0, 1]]]]]
```

B.2. Sobel edge extractor

```
R_edge_h = conv2d(R_imgs, filter_h, stride=1, padding=1)
G_edge_h = conv2d(G_imgs, filter_h, stride=1, padding=1)
B_edge_h = conv2d(B_imgs, filter_h, stride=1, padding=1)
R_edge_v = conv2d(R_imgs, filter_v, stride=1, padding=1)
G_edge_v = conv2d(G_imgs, filter_v, stride=1, padding=1)
B_edge_v = conv2d(B_imgs, filter_v, stride=1, padding=1)

R_edge = R_edge_h ** 2 + R_edge_v ** 2
G_edge = G_edge_h ** 2 + G_edge_v ** 2
B_edge = B_edge_h ** 2 + B_edge_v ** 2
```

B.3. Patch discriminator

Type	Kemel	Stride	Outputs	Padding	Dilation rate	Nonlinearity
Disconv	5×5	2	48	SAME	1	LeakyReLU(0.2)
Disconv	5×5	2	96	SAME	1	LeakyReLU(0.2)
Disconv	5×5	2	192	SAME	1	LeakyReLU(0.2)
Disconv	5×5	2	192	SAME	1	LeakyReLU(0.2)
Disconv	5×5	2	192	SAME	1	LeakyReLU(0.2)
Disconv	5×5	2	384	SAME	1	LeakyReLU(0.2)

B.4. Sobel edge discriminator

Type	Kemel	Stride	Outputs	Padding	Dilation rate	Nonlinearity
Disconv	5×5	4	24	SAME	1	LeakyReLU(0.2)
Disconv	5×5	2	48	SAME	1	LeakyReLU(0.2)
Disconv	5×5	2	96	SAME	1	LeakyReLU(0.2)
Disconv	5×5	2	192	SAME	1	LeakyReLU(0.2)
Disconv	5×5	2	192	SAME	1	LeakyReLU(0.2)

B.5. Low pass filter / High pass filter

```

low_pass_filter = ones(size=(1, 1, 5, 5)) / 25.0
high_pass_filter = [[[-1, -1, -1], [-1, 9, -1], [-1, -1, -1]]]

R_low_imgs = conv2d(R_imgs, low_pass_filter, stride=1)
G_low_imgs = conv2d(G_imgs, low_pass_filter, stride=1)
B_low_imgs = conv2d(B_imgs, low_pass_filter, stride=1)

R_high_imgs = conv2d(R_imgs, high_pass_filter, stride=1)
G_high_imgs = conv2d(G_imgs, high_pass_filter, stride=1)
B_high_imgs = conv2d(B_imgs, high_pass_filter, stride=1)

```

B.6. Generator

Layer	Type	Ker- nel	Str- ide	Out- puts	Padding	Nonli- nearity
Encoder_layer1 (skip connection)	RecurrentConv	5	1	16	SYMMETRIC	ELU
Encoder_layer2	RecurrentConv	3	2	32	SYMMETRIC	ELU
Encoder_layer3 (skip connection)	RecurrentConv	3	1	32	SYMMETRIC	ELU
Encoder_layer4	RecurrentConv	3	1	64	SYMMETRIC	ELU
Encoder_layer5	RecurrentConv	3	1	64	SYMMETRIC	ELU
Encoder_layer6	RecurrentConv	3	1	64	SYMMETRIC	ELU
Decoder_layer1	RecurrentConv	3	1	64	SYMMETRIC	ELU
Decoder_layer2	RecurrentConv	3	1	64	SYMMETRIC	ELU
Decoder_layer3	RecurrentDeConv	3	1	32	SYMMETRIC	ELU
Decoder_layer4 (skip connection)	RecurrentConv	3	1	32	SYMMETRIC	ELU
Decoder_layer5	RecurrentDeConv	3	1	16	SYMMETRIC	ELU
Decoder_layer6 (skip connection)	RecurrentConv	3	1	16	SYMMETRIC	ELU
Decoder_layer7	RecurrentConv	3	1	8	SYMMETRIC	ELU
Decoder_layer8	RecurrentConv	3	1	3	SYMMETRIC	NONE

References

- [1] Zheng Hui, Jie Li, Xiumei Wang, and Xinbo Gao. Image fine-grained inpainting. *arXiv preprint arXiv:2002.02609*, 2020. **2**
- [2] Satoshi Iizuka, Edgar Simo-Serra, and Hiroshi Ishikawa. Globally and locally consistent image completion. *ACM Transactions on Graphics (ToG)*, 36(4):1–14, 2017. **1**
- [3] FG Irwin et al. An isotropic 3x3 image gradient operator. *Presentation at Stanford AI Project*, 2014(02), 1968. **1, 3**
- [4] Yuqing Ma, Xianglong Liu, Shihao Bai, Lei Wang, Dailan He, and Aishan Liu. Coarse-to-fine image inpainting via region-wise convolutions and non-local correlation. In *IJ-CAI*, pages 3123–3129, 2019. **1**
- [5] Kamyar Nazeri, Eric Ng, Tony Joseph, Faisal Z Qureshi, and Mehran Ebrahimi. Edgeconnect: Generative image inpainting with adversarial edge learning. *arXiv preprint arXiv:1901.00212*, 2019. **1, 2, 3, 6**
- [6] 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**
- [7] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Free-form image inpainting with gated convolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4471–4480, 2019. **1, 2, 3, 6**
- [8] Chuanxia Zheng, Tat-Jen Cham, and Jianfei Cai. Pluralistic image completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1438–1447, 2019. **1, 3, 6**