

Image Inpainting Through Coherent Semantic Attention

Jialiang Yan

Abstract—The latest deep learning-based approaches have shown promising results for the challenging task of inpainting missing regions of an image. However, the existing methods often generate contents with blurry textures and distorted structures due to the discontinuity of the local pixels. From a semantic-level perspective, the local pixel discontinuity is mainly because these methods ignore the semantic relevance and feature continuity of hole regions. To handle this problem, we investigate the human behavior in repairing pictures and propose a refined deep generative model-based approach with a novel coherent semantic attention (CSA) layer, which can not only preserve contextual structure but also make more effective predictions of missing parts by modeling the semantic relevance between the holes features. The task is divided into rough, refinement as two steps and model each step with a neural network under the U-Net architecture, where the CSA layer is embedded into the encoder of refinement step. To stabilize the network training process and promote the CSA layer to learn more effective parameters, we propose a consistency loss to enforce the both the CSA layer and the corresponding layer of the CSA in decoder to be close to the VGG feature layer of a ground truth image simultaneously. The experiments on CelebA, Places2, and Paris StreetView datasets have validated the effectiveness of our proposed methods in image inpainting tasks and can obtain images with a higher quality as compared with the existing state-of-the-art approaches.

Index Terms—convolutional neural network, deep learning, generative adversarial network, image inpainting.

I. INTRODUCTION

Image inpainting is the task to synthesize the missing or damaged parts of a plausible hypothesis, and can be utilized in many applications such as removing unwanted objects, completing occluded regions, restoring damaged or corrupted parts. The core challenge of image inpainting is to maintain global semantic structure and generate realistic texture details for the missing regions.

Traditional works [2, 3, 11, 12, 34] mostly develop texture synthesis techniques to address the problem of hole filling. In [2], Barnes et al. propose the Patch-Match algorithm which iteratively searches for the best fitting patches from hole boundaries to synthesize the contents of the missing parts. Wilczkowiak et al. [34] take further steps and detect desirable search regions to find better match patches. However, these methods fall short of understanding high-level

semantics and struggle at reconstructing patterns that are locally unique. In contrast, early deep convolution neural networks based approaches [17, 24, 30, 39] learn data distribution to capture the semantic information of the image, and can achieve plausible inpainting results. However, these methods fail to effectively utilize contextual information to generate the contents of holes, often leading to the results containing noise patterns.

Some recent studies effectively utilize the contextual information and obtain better inpainting results. These methods can be divided into two types. The first type [32, 36, 42] utilizes spatial attention which takes surrounding image features as references to restore missing regions. These methods can ensure the semantic consistency of generated content with contextual information. However, they just focus on rectangular shaped holes, and the results always tend to show pixel discontinuous and have semantic chasm (See in Fig 1(b, c)). The second type [26, 41] is to make the prediction of the missing pixels condition on the valid pixels in the original image. These methods can handle irregular holes properly, but the generated contents still meet problems of semantic fault and boundary artifacts (See in Fig 1(g, h)). The reason that the above mentioned methods do not work well is because they ignore the semantic relevance and feature continuity of generated contents, which is crucial for the local pixel continuity.

II. PROCEDURE FOR PAPER SUBMISSION

In order to achieve better image restoration effect, we investigate the human behavior in inpainting pictures and find that such process involves two steps as conception and painting to guarantee both global structure consistency and local pixel continuity of a picture. To put it more concrete, a man first observes the overall structure of the image and conceives the contents of missing parts during conception process, so that the global structure consistency of the image can be maintained. Then the idea of the contents will be stuffed into the actual image during painting process. In the painting process, one always continues to draw new lines and coloring from the end nodes of the lines drawn previously, which actually ensures the local pixel continuity of the final result.

Inspired by this process, we propose a coherent semantic attention layer (CSA), which fills in the unknown regions of the image feature maps with the similar process. Initially, each unknown feature patch in the unknown region is initialized with the most similar feature patch in the known regions. Thereafter, they are iteratively optimized by considering the spatial consistency with adjacent patches. Consequently, the global semantic consistency is guaranteed by

Manuscript received October 23, 2021

Jialiang Yan, School of Computer Science and Technology, Tiangong University, Tianjin, China

the first step, and the local feature coherency is maintained by the optimizing step.

Similar to [42], we divide the image inpainting into two steps. The first step can be constructed by training a rough network to rough out the missing contents. A refinement network with the CSA layer in encoder guides the second step to refine the rough predictions. In order to make network training process more stable and motivate the CSA layer to learn more effective features, we propose a consistency loss to measure not only the distance between the VGG feature layer and the CSA layer but also the distance between the VGG feature layer and the corresponding layer of the CSA in decoder. Meanwhile, in addition to a patch discriminator [18], we improve the details by introducing a feature patch which is simpler in formula, faster and more stable for training than conventional one [29]. Except for the consistency loss, reconstruction loss, and relativistic average LS adversarial loss [28] are incorporated as constraints to instruct our model to learn meaningful parameters. We conduct experiments on standard datasets CelebA [27], Places2 [44], and Paris StreetView [8]. Both the qualitative and quantitative tests demonstrate that our method can generate higher-quality inpainting results than existing ones. (See in Fig 1(d, i)). Our contributions are summarized as follows: We propose a novel coherent semantic attention layer to construct the correlation between the deep features of hole regions. No matter whether the unknown region is irregular or centering, our algorithm can achieve state-of-the-art inpainting results.

To enhance the performance of the CSA layer and training stability, we introduce the consistency loss to guide the CSA layer and the corresponding decoder layer to learn the VGG features of ground truth. Meanwhile, a feature patch discriminator is designed and joined to achieve better predictions.

Our approach achieves higher-quality results in comparison with [26,36,41,42] and generates more coherent textures. Even the inpainting task is completed in two stages, our full network can be trained in an end to end manner.

III. APPROACH

Our model consists of two steps: rough inpainting and refinement inpainting. This architecture helps to stabilize training and enlarge the receptive fields as mentioned in [42]. The overall framework of our inpainting system is shown in Fig 2. Let I_{gt} be the ground truth images, I_{in} be the input to the rough network, the M and M' denote the missing area and the known area in feature maps respectively. We first get the rough prediction I_p during the rough inpainting process. Then, the refinement network with CSA layer takes the I_p and I_{in} as input pairs to output final result I_r . Finally, the patch and feature patch discriminators work together to obtain higher resolution of I_r .

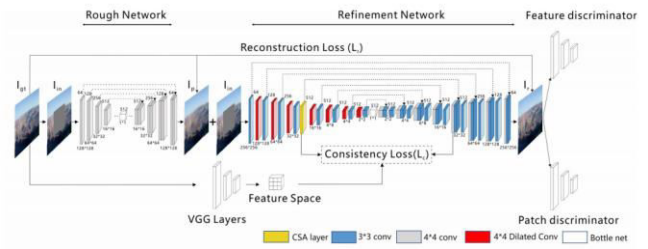


Figure 1: our network

IV. EXPERIMENTS

Our proposed model is evaluated on the datasets CelebA. Results are compared against the current state-of-the-art methods both qualitatively and quantitatively. Shown in Figure 2 and Figure 3. The image inpainting indicators are compared in Table 1



Figure 2 Our results show in celebA

Method	PSNR	SSIM
PC	21.34	0.814
GC	21.42	0.819
Our	21.75	0.823

V. CONCLUSION

In this paper, we proposed a refined deep generative model based approach which designed a novel Coherent Semantic Attention layer to learn the relationship between features of missing region in image inpainting task. The consistency loss is introduced to enhance the CSA layer learning ability for ground truth feature distribution and training stability. Moreover, a feature patch discriminator is joined into our model to achieve better predictions. Experiments have verified the effectiveness of our proposed methods. In future, we plan to extend the method to other tasks, such as style transfer and single image super-resolution.

TRANSFER AND SINGLE IMAGE
SUPER-RESOLUTION.REFERENCES

- [1] Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. 1, 3
- [2] M. Eitz, J. Hays, and M. Alexa. How do humans sketch objects? *ACM Transactions on graphics (TOG)*, 31(4):44 – 1, 2012. 2
- [3] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015. 2, 5
- [4] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017. 2, 5
- [5] J. Canny. A computational approach to edge detection. *IEEE Transactions on pattern analysis and machine intelligence*, (6):679 – 698, 1986. 3
- [6] S. Nowozin, C. H. Lampert, et al. Structured learning and prediction in computer vision. *Foundations and Trends R ? in Computer Graphics and Vision*, 6(3 – 4):185 – 365, 2011. 3
- [7] S. Xie and Z. Tu. Holistically-nested edge detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1395 – 1403, 2015. 3, 7, 12
- [8] P. Sangkloy, J. Lu, C. Fang, F. Yu, and J. Hays. Scribbler: Controlling deep image synthesis with sketch and color. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 2, 2017. 3
- [9] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 3, 11
- [10] M. Mirza and S. Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 3
- [11] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *The IEEE International Conference on Computer Vision (ICCV)*, 2017. 3, 11
- [12] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672 – 2680, 2014. 1, 3
- [13] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision (ECCV)*, pages 694 – 711. Springer, 2016. 3, 4, 11
- [14] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770 – 778, 2016. 3
- [15] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 3, 11
- [16] D. Ulyanov, A. Vedaldi, and V. 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 (CVPR)*, 2017. 3