# A ONE-SHOT TEXTURE-PERCEIVING GENERATIVE ADVERSARIAL NETWORK FOR UNSUPERVISED SURFACE INSPECTION

*Lingyun Gu*[†]       *Lin Zhang*[⋆]       *Zhaokui Wang*[†]

Tsinghua University, Beijing, China[†]
University of Cincinnati, Cincinnati, Ohio, USA[⋆]

## ABSTRACT

Visual surface inspection is a challenging task owing to the highly diverse appearance of target surfaces and defective regions. Previous attempts heavily rely on vast quantities of training examples with manual annotation. However, in some practical cases, it is difficult to obtain a large number of samples for inspection. To combat it, we propose a hierarchical texture-perceiving generative adversarial network (HTP-GAN) that is learned from the one-shot normal image in an unsupervised scheme. Specifically, the HTP-GAN contains a pyramid of convolutional GANs that can capture the global structure and fine-grained representation of an image simultaneously. This innovation helps distinguishing defective surface regions from normal ones. In addition, in the discriminator, a texture-perceiving module is devised to capture the spatially invariant representation of normal image via directional convolutions, making it more sensitive to defective areas. Experiments on a variety of datasets consistently demonstrate the effectiveness of our method.

***Index Terms***— One-shot learning, texture-perceiving module, visual surface inspection, generative adversarial network

## 1. INTRODUCTION

Due to the rapid development of deep neural networks [1, 2, 3, 4], visual surface inspection [5, 6] has attracted increasing attention as an important technology in many intelligent industrial applications. Visual surface inspection aims to detect the abnormal regions on the surface of material using visual images. It is a challenging task owing to various image noises, texture variations of the target surface, and highly diversified appearance of abnormal regions.

Typical visual inspection approaches can be categorized into two main groups: traditional methods [6, 7] and learning-based methods [5, 8, 9, 10]. Traditional methods adopt hand-crafted features to perform surface inspection, which cannot be well generalized to new scenarios. Learning-based approaches achieve significant performance when equipped with large amounts of manually annotated training data. However, in some practical cases, such as surface inspection under planes [11], only a small number of normal samples, or even a single sample, are available. Therefore, how to design the one-shot surface inspection method is very important for practical scenarios.

We present the problem of one-shot unsupervised surface inspection: given an example of a normal image as training data, all defective regions should be detected and segmented for arbitrary images with the same texture category as the normal image. The challenge lies in 1) how to design an adaptable perceiving model that is prone to handle the texture variations; 2) how to perform a more generalized surface inspection model in the one-shot way.

To deal with the challenge, we propose a novel hierarchical texture-perceiving generative adversarial network (HTP-GAN) that is learned from a one-shot normal image in an unsupervised scheme. Specifically, HTP-GAN model contains a pyramid of convolutional GANs which utilize a single image to simultaneously extract the global and fine-grained representation of the image. This guides the model to indeed learn the normal texture representation of the category of the corresponding image and can distinguish the defective surface regions. In addition, in the discriminator of HTP-GAN, texture-perceiving module is devised to capture the spatially invariant representation of normal image via directional convolutions, making it more sensitive to defective areas. Experiments on a variety of datasets consistently demonstrate the effectiveness of our method.

We summarize our main contributions as follows:

- We introduce a novel one-shot surface inspection network with a pyramid of convolutional GANs, achieving unsupervised surface inspection with a single example of normal image.
- A texture-perceiving module is devised to capture the spatially invariant representation of normal image via directional convolutions, making it more sensitive to defective areas.
- Experimental results achieve the state-of-the-art performance on two public datasets, which demonstrate the effectiveness of the proposed approaches.
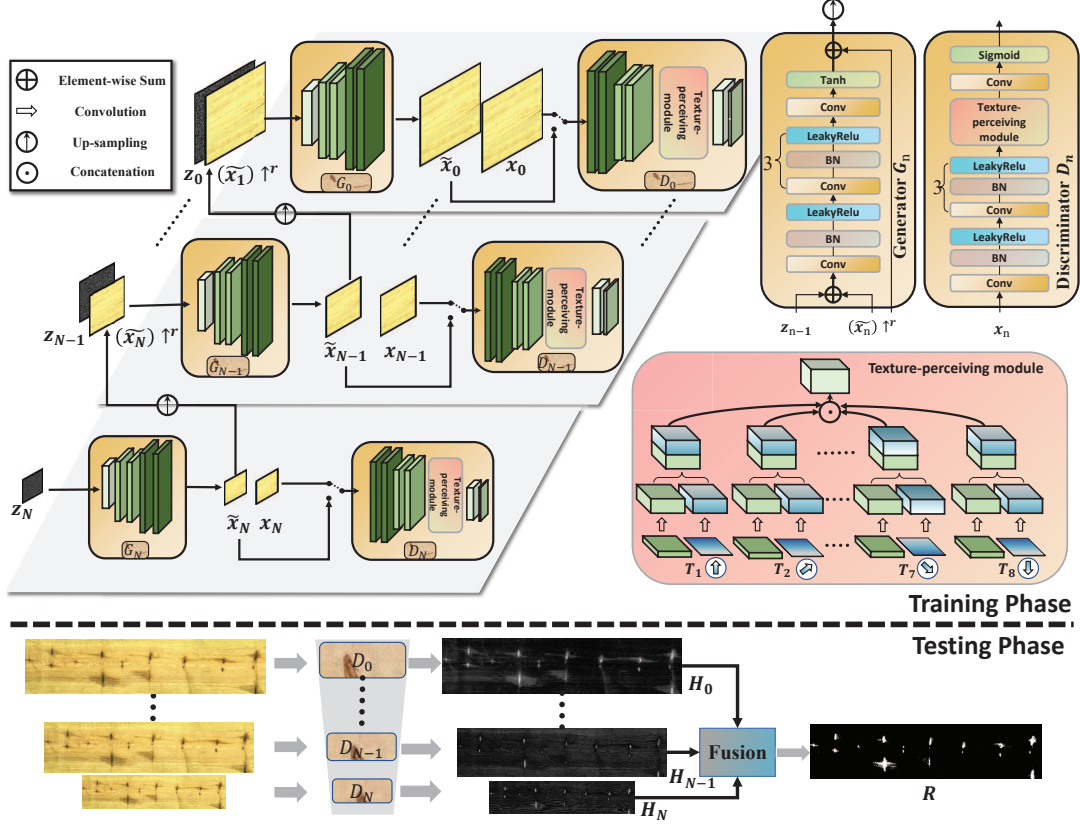
**Fig. 1**. Overview of the proposed hierarchical texture-perceiving generative adversarial network (HTP-GAN). In the training phase, the HTP-GAN is devised to capture the spatially invariant representation of a single normal image via directional convolutions at multiple scales. In the testing phase, the well-trained discriminators can indeed memorize the latent distribution properties of normal texture representation and distinguish the defective surface regions.

## 2. HIERARCHICAL TEXTURE-PERCEIVING GENERATIVE ADVERSARIAL NETWORK

Our goal is to detect the abnormal regions on the surface of material from a single image. When large amounts of training data are available, a well-trained generative adversarial network (GAN) can easily learn a representation of the distribution of the target samples via a generator $G$ and a discriminator $D$. But in the case of the limited number of samples, conventional GAN is hard to learn such a representation due to insufficient training data [12]. In the context of learning from a single image, inspired by the SinGAN [1], we adopt the downsampling strategy to generate different scales of images from a single image as training samples. Then the network is able to learn internal distributions at different scales. At the same time, a pyramid of convolutional GANs is utilized to simultaneously capture the statistics of complex image structures and perceive the fine-grained representation of the normal image. To our best knowledge, this is the first time that the SinGAN being applied in the one-shot surface inspection task.

### 2.1. Hierarchical Fully Convolution Architecture

In this section, we will detail the overall architecture of the proposed HTP-GAN. As shown in Fig. 1, the network is based on SinGAN [1] to achieve the one-shot learning task. We change the traditional GAN from two perspectives: 1) we resize the input image to generate the multi-scale images. Specifically, given an image of normal surface $x_0$, we adopt the down-sampled strategy to generate an image pyramid of $x : \{x_0, \ldots, x_N\}$. Multi-scale inputs contain more fine details and texture information. 2) a pyramid of convolutional GANs are utilized to simultaneously capture the characteristics of complex image structures and perceive fine-grained representation of the training image. So we need to design a pyramid of generators $\{G_0, \ldots, G_N\}$ and discriminators $\{D_0, \ldots, D_N\}$ to handle the multi-scale inputs.

At the training phase, the proposed HTP-GAN is a multi-stage training process from coarse-grain to fine-grain. Firstly, a noise map is injected into the generator $G_N$ at the coarsest scale to generate the $\tilde{x}_N$. Then a combination between the generated image and the noise map sequentially passes through all other generators up to the finest scale:

$$\tilde{x}_n = G_n\left(z_n, (\tilde{x}_{n+1}) \uparrow^r\right) = (\tilde{x}_{n+1}) \uparrow^r + \psi_n\left(z_n + (\tilde{x}_{n+1}) \uparrow^r\right), \quad (1)$$

where $\psi_n$ is a fully convolutional net with 5 conv-blocks and $\uparrow^r$ denote to the up-sampling operation. As shown in Fig. 1, the first 4 conv-blocks consist of a $3 \times 3$ convolutional layer following with a BatchNorm and a LeakyReLU. The last conv-block consist of a $3 \times 3$ convolutional layer and a tanh activation layer. In addition, the $\psi_n$ is able to learn the residual feature of a generated image at a finer scale.

## 2.2. Textual-perceiving Discriminator

A textual-perceiving module in the discriminator is devised to capture the spatially invariant representation of a normal image via directional convolutions, making it more sensitive to defective areas. Specifically, the textual-perceiving discriminator consists of the several conv-blocks $f_n$, the textual-perceiving module [13] ($c_n$, $g_n$) and a sigmoid activation layer. As shown in Fig. 1, the several conv-blocks $f_n$ are used to extract preliminary features $F_n$ from the generated image. Then, the $c_n$ and $g_n$ are designed to capture the spatially invariant representation of normal image via eight directional texture features under the guidance of corresponding directional feature map $T_i$. Specifically,

$$P_n = c_n\left(\text{cat}\left(f_n(\tilde{x}_n), T_i\right)\right) (i = 1, 2 \cdots 8), \quad (2)$$

where $c_n$ denotes several convolutions which contains a $3 \times 3$ convolution layer with a BatchNorm and a ReLU activation layer, $cat$ denotes the concatenation operation. In our design, total eight directions are adopted, including top, bottom, left, right, top left, bottom left, top right, and bottom right. Each directional map $T_i$ is a generated trend square matrix that decreases from 1 to 0 in a certain direction. Finally, the output features of different branches are concatenated to form the whole spatial invariant feature $M_n$, that is,

$$M_n = Sigmoid(g_n\left(\text{cat}\left(P_1, P_2 \cdots P_8\right)\right)), \quad (3)$$

where $g_n$ represents a set of standard convolution block, $M_n \in \mathbb{R}^{1 \times H \times W}$ is the distinguish map that is used to calculate the loss of discriminator. Since the proposed directional convolution unit is sensitive to local variations of the image along each direction, it can make the network well adaptable to spatial distortions and scale variations. Intuitively, the discriminator of a well-trained GAN memories the patch distribution of a normal image. So it should be insensitive to the normal regions but varies drastically in the abnormal regions.

## 2.3. Hierarchical Fusion for Surface Inspection

In this section, we introduce a hierarchical fusion strategy to fuse the multi-scale distinguish maps for producing the final segmented result.

Inspired by [5], the information entropy is a suitable metric to represent the output of discriminator for abnormal region segmentation. Thus, the information entropy of multi-scale distinguish maps are expressed as following:

$$H_n = M_n * \log M_n, \quad (4)$$

$H = \{H_n\}|_{n=0}^N$ are a set of information entropy of multi-scale distinguish maps and we refer to $H$ as the coarse inspection map set, which reveals the coarse abnormal regions of the input image.

The GANs have small receptive fields and limited capacity, preventing them from representing the single image. To capture global textual structure and fine-grained texture information, we fuse the multi-scale coarse inspection maps together to better distinguish the abnormal regions. Thus, the fusion process expresses as following:

$$R = \sum_{n=0}^N \alpha_n * H_n, \quad (5)$$

where $R$ denotes the final fusion and $\alpha_n = \frac{1}{N+1}$ are weighting factors.

## 3. EXPERIMENTS

To evaluate the effectiveness of HTP-GAN in one-shot surface inspection task, we design extensive experiments based on WOOD Defect Database (WOOD) [14] and Road Crack Database (CRACK) [15]. Furthermore, we compare our method with four representative methods as following: (1) Unsupervised visual surface inspection method proposed in [5] (ICASSP 2018); (2) Surface defect detection method based on positive samples and artificial defects [16] (Prical 2018); (3) Automated surface inspection method proposed in [17] (JOIM 2019); (4) Semantic image segmentation method proposed in [18] (ECCV 2018).

### 3.1. Datasets and Evaluation Metrics

**Dataset and Setting**: WOOD Defect Database (WOOD) and Road Crack Database (CRACK) are two surface inspection datasets which are annotated with segmentation labels of abnormal regions [5]. For our HTP-GAN, a single image with normal surface is adopted as the training data, which is the embodiment of the one-shot surface inspection. To compare our method with unsupervised methods: [5] and [16], we adopt the same image with normal surface to train the models. For supervised methods: [17] and [18], several abnormal samples and normal samples are used to train the networks. All images are resized to $256 \times 256$ size.

**Evaluation Metrics**: We adopt two evaluation metrics to compare and analyze experimental results: Intersection-over-union (IOU) and Pixel Accuracy (pixel acc).
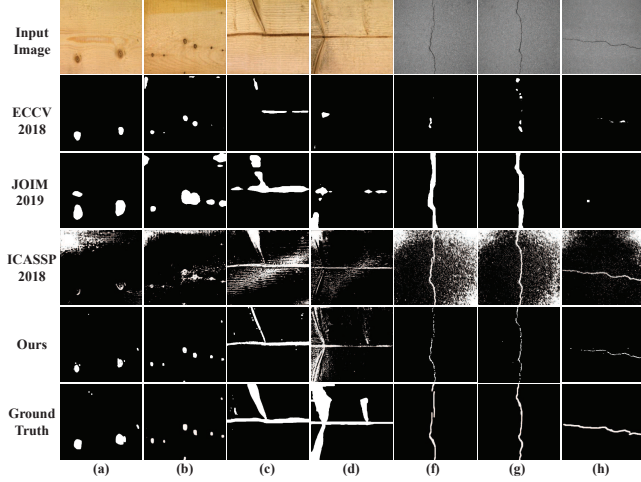
**Fig. 2**. Visualization comparison with the state-of-the-art methods for the surface inspection task.

**Table 1**. Quantitative comparisons. (IoU(%) / pixel acc(%))

|  | [5] | [18] | [17] | Ours |
|---|---|---|---|---|
| WOOD | 44.83/81.96 | 47.23/94.45 | 47.37/93.82 | **59.83/96.54** |
| CRACK | 31.06/58.21 | 48.68/90.46 | 33.02/61.56 | **56.97/96.32** |

### 3.2. Comparison with the State-of-the-art Approaches

Tab. 1 shows the performance comparison of our HTP-GAN against the other four methods in terms of IoU and pixel acc on the WOOD and CRACK dataset. As we can see in Tab. 1, our method achieved higher IoU and pixel acc than the four methods.

In the one-shot setting, when the single normal image or an image pair (a normal image and an abnormal image) is given as training data, all four methods can not predict any abnormal regions or produce a noise map. This result demonstrates that these methods rely on a tremendous amount of training data and these models are easy to cause overfitting on one-shot surface inspection task. Thus, we try to add the training data on these methods in the supervised learning methods or adopt the fully convolution network to improve the unsupervised learning methods [1].

Fig. 2 shows the visualization of different methods for surface inspection. [17] has a good performance until we adopt 10 normal-abnormal image pairs for training while [18] performs well on 4 abnormal images. Although improved [5] detects the more accurate abnormal surface regions than other methods, our proposed model has a better performance. More importantly, HTP-GAN only needs a normal image as training data and achieves the best performance.

### 3.3. Ablation Studies

**Effectiveness of Scale Number.** The number of scales in HTP-GAN architecture has a strong influence on the results. A small number of scales only can capture the local tex-
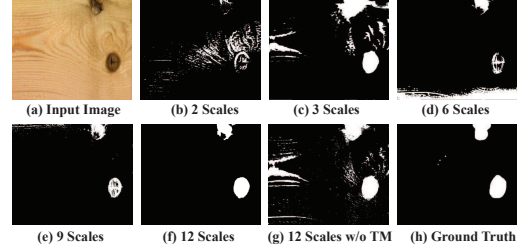


**Fig. 3**. Visualization of influence of different scales and textual-perceiving module in our proposed HTP-GAN. "TM" refers to the textual-perceiving module.

tures, leading to poor segmentation results. As the number of scales increases, Fig. 3 (a)-(f) demonstrates that our proposed method manages to capture larger structures as well as the fine-grained representation of the image. It indicates that a strong representation can help the model to segment most of the anomalous regions very well.

**Effectiveness of Textual-perceiving Module.** As shown in Fig. 3 (g) and (f), by adding the textual-perceiving module, our method segments more abnormal regions and alleviates to split normal regions into defects. It demonstrates that the proposed textual-perceiving module can learn spatially invariant representation of normal images and be more sensitive to defective areas.

**Influence of Different Image Variations.** Our proposed method can handle most variations (*i.e.*Translation, Mirror, Scaling). However, for rotation variation, when the training image contains mostly horizontal textures (such as the wood in Fig. 3, a 45-degree change in rotation only brings a slight performance drop 3% in term of IOU score, and no performance drop in term of pixel accuracy.

### 4. CONCLUSION

In this paper, we propose a hierarchical texture-perceiving generative adversarial network (HTP-GAN) to achieve the one-shot surface inspection task. By applying a pyramid of convolutional GANs, HTP-GAN can indeed learn the global and fine-grained representation of normal surface from a single image. Thus, it enables discriminators in well-trained HTP-GAN to exhibit more active responses for the defective regions. Furthermore, the texture-perceiving module is devised to capture the spatially invariant representation of normal, making the discriminator more sensitive to defective areas. More importantly, the training of the HTP-GAN only relies on a single image, which explores a more general and practical solution to the surface inspection task.

## Acknowledgments

# 5. REFERENCES

[1] Tamar Rott Shaham, Tali Dekel, and Tomer Michaeli, "Singan: Learning a generative model from a single natural image," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 4570–4580.

[2] Kecheng Zheng, Wu Liu, Jiawei Liu, Zheng-Jun Zha, and Tao Mei, "Hierarchical gumbel attention network for text-based person search," in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 3441–3449.

[3] Kecheng Zheng, Wu Liu, Lingxiao He, Tao Mei, Jiebo Luo, and Zheng-Jun Zha, "Group-aware label transfer for domain adaptive person re-identification," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021.

[4] Kecheng Zheng, Cuiling Lan, Wenjun Zeng, Zhizheng Zhan, and Zheng-Jun Zha, "Exploiting sample uncertainty for domain adaptive person re-identification," *Association for the Advancement of Artificial Intelligence*, 2021.

[5] W. Zhai, J. Zhu, Y. Cao, and Z. Wang, "A generative adversarial network based framework for unsupervised visual surface inspection," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2018, pp. 1283–1287.

[6] Choon-Woo Kim and Antti J Koivo, "Hierarchical classification of surface defects on dusty wood boards," *Pattern Recognition Letters*, vol. 15, no. 7, pp. 713–721, 1994.

[7] W Wen and Aihua Xia, "Verifying edges for visual inspection purposes," *Pattern recognition letters*, vol. 20, no. 3, pp. 315–328, 1999.

[8] R. Liu, Q. Gu, X. Wang, and M. Yao, "Region-convolutional neural network for detecting capsule surface defects," *Boletin Tecnico/Technical Bulletin*, vol. 55, no. 3, pp. 92–100, 2017.

[9] Kecheng Zheng, Zheng-Jun Zha, and Wei Wei, "Abstract reasoning with distracting features," in *Advances in Neural Information Processing Systems*, 2019, pp. 5842–5853.

[10] Zhiyang Yu, Xiaojun Wu, and Xiaodong Gu, "Fully convolutional networks for surface defect inspection in industrial environment," in *International Conference on Computer Vision Systems*. Springer, 2017, pp. 417–426.

[11] Igor Jovančević, Huy-Hieu Pham, Jean-José Orteu, Rémi Gilblas, Jacques Harvent, Xavier Maurice, and Ludovic Brèthes, "3d point cloud analysis for detection and characterization of defects on airplane exterior surface," *Journal of Nondestructive Evaluation*, vol. 36, no. 4, pp. 74, 2017.

[12] Shengyu Zhao, Zhijian Liu, Ji Lin, Jun-Yan Zhu, and Song Han, "Differentiable augmentation for data-efficient gan training," *Advances in Neural Information Processing Systems*, vol. 33, 2020.

[13] Yanzhao Zhou, Qixiang Ye, Qiang Qiu, and Jianbin Jiao, "Oriented response networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 519–528.

[14] Olli Silvén, Matti Niskanen, and Hannu Kauppinen, "Wood inspection with non-supervised clustering," *Machine Vision and Applications*, vol. 13, no. 5-6, pp. 275–285, 2003.

[15] Henrique Oliveira and Paulo Lobato Correia, "Crackit—an image processing toolbox for crack detection and characterization," in *IEEE International Conference on Image Processing*. IEEE, 2014, pp. 798–802.

[16] Zhixuan Zhao, Bo Li, Rong Dong, and Peng Zhao, "A surface defect detection method based on positive samples," in *Pacific Rim International Conference on Artificial Intelligence*. Springer, 2018, pp. 473–481.

[17] Domen Tabernik, Samo Šela, Jure Skvarč, and Danijel Skočaj, "Segmentation-based deep-learning approach for surface-defect detection," *Journal of Intelligent Manufacturing*, vol. 31, no. 3, pp. 759–776, 2020.

[18] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 801–818.