# Semantic Image Translation for Repairing the Texture Defects of Building Models

Qisen Shang, Han Hu<sup>\*</sup>, Haojia Yu, Bo Xu, Libin Wang, Qing Zhu

Faculty of Geosciences and Environmental Engineering, Southwest Jiaotong University, Chengdu, China

#### Abstract

The accurate representation of 3D building models in urban environments is significantly hindered by challenges such as texture occlusion, blurring, and missing details, which are difficult to mitigate through standard photogrammetric texture mapping pipelines. Current image completion methods often struggle to produce structured results and effectively handle the intricate nature of highly-structured facade textures with diverse architectural styles. Furthermore, existing image synthesis methods encounter difficulties in preserving high-frequency details and artificial regular structures, which are essential for achieving realistic facade texture synthesis. To address these challenges, we introduce a novel approach for synthesizing facade texture images that authentically reflect the architectural style from a structured label map, guided by a ground-truth façade image. In order to preserve fine details and regular structures, we propose a regularity-aware multi-domain method that capitalizes on frequency information and corner maps. We also incorporate SEAN blocks into our generator to enable versatile style transfer. To generate plausible structured images without undesirable regions, we employ image completion techniques to remove occlusions according to semantics prior to image inference. Our proposed method is also capable of synthesizing texture images with specific styles for facades that lack pre-existing textures, using manually annotated labels. Experimental results on publicly available facade image and 3D model datasets demonstrate that our method yields superior results and effectively addresses issues associated with flawed textures. The code and datasets will be made publicly available for further research and development.

*Keywords:* Oblique Photogrammetry, 3D Building Model, Texture Mapping, Image Translation, Generative Adversarial Network (GAN)

## 1. Introduction

Three-dimensional (3D) building models are fundamental to the development of digital cities, serving as a crucial component in high-precision mapping, autonomous driving, and urban planning (Lin et al., 2013; Tao and Qi, 2019). The realism of 3D models is conveyed through their geometry and texture (Chen et al., 2020; Buyukdemircioglu and Kocaman, 2020). At present, high-precision building models are predominantly created manually, while textures are primarily sourced from various imagery. However, in densely populated urban areas, aerial imagery is often obstructed by building occlusions, and ground-level photography is impeded by adjacent objects, such as vegetation and billboards. As a result, acquiring unobstructed imagery from any angle and distance proves challenging, rendering traditional texture mapping pipelines inadequate for

Preprint submitted to ISPRS Journal of Photogrammetry and Remote Sensing

<sup>\*</sup>Corresponding Author: han.hu@swjtu.edu.cn

processing or de-occluding building façade images (Zhou et al., 2020; Zhu et al., 2021a; Zhang et al., 2021; Li et al., 2023). As depicted in Figure 1, defective textures are prevalent in various realistic building models.

Current models frequently utilize occluded textures directly or substitute them with manually crafted repetitive textures. In certain cases, textures are selected from pre-existing material libraries. These processing techniques substantially constrain the visualization quality of models (Li et al., 2023). To improve the visual fidelity of realistic 3D building models, it is crucial to tackle the issue of defective textures. As a result, this paper concentrates on repairing defective textures in realistic 3D building models to facilitate high-precision urban 3D reconstruction and advanced applications. Although there have been advancements in texture processing and occlusion removal, a number of outstanding challenges still need to be addressed.



(a) Occlussions of SWJTU model. (b) Texture missing of London model.

Figure 1: Defective textures of 3D building models. The red rectangles denote enlarged regions.

1) Inaccessible façade textures in built-up areas. Multi-view aerial camera systems, mobile measurement systems (MMS), and handheld cameras are widely employed and effectively integrated in urban realistic 3D modeling for tasks such as feature matching and texture mapping (Remondino and Gerke, 2015; Zhu et al., 2020b, 2021b). However, there remain many areas that are inaccessible to these sensors, particularly buildings in built-up areas. Some researchers have attempted to mitigate occlusions by selecting optimal pixels from various aerial images for texture remapping, but this approach is limited to areas that are not fully occluded (Zhou et al., 2020; Yang et al., 2021). A mesh completion method has been proposed in previous work, which processes 3D textures in 2D space and employs image completion methods to repair textures in occluded road areas. Nevertheless, this method cannot handle unseen pixels or generate textures from scratch, which is a frequent issue in built-up areas (Zhu et al., 2021b). Moreover, we observe that recovering the semantics of occluded buildings is considerably more manageable than completing textures. As such, synthesizing textures from recovered or specified semantic labels presents a practical solution.

2) Highly structured and diverse styles of building façades. Windows, doors, balconies, and other components on façades contribute to the textures of buildings. An ideal building façade texture should exhibit structured man-made components and a realistic style with details. However, neither data-driven nor patch-based methods possess sufficient generalization for different buildings with various styles, and they cannot generate highly structured de-occluded or synthesized textures (Criminisi et al., 2004; Zhu et al., 2021b). Moreover, the results of existing image synthesis methods lack texture details, such as bricks and window frames, which are crucial for building façades (Cai et al., 2021). The regularization and generalization of current methods

cannot fulfill the requirements of building façade texture repair. As such, one viable solution is to first recover or provide the semantics of façades to control regularization, and then translate the semantic label map into textures by combining patch-based and data-based methods.

To address the issues outlined above, we propose a method for repairing building façade textures to enhance the realism of 3D building models. Our method involves synthesizing realistic façade textures by using unoccluded semantic labels and ground-truth façade images, a process also known as image translation (Isola et al., 2017). Our approach is capable of handling both occluded and missing façade textures. The difference in processing occlusion and missing façade textures lies primarily in the acquisition of semantic labels. For occlusion, we use an image completion algorithm to recover the semantics of occluded regions. For missing textures, we manually annotate the semantics of the corresponding façade. We then utilize the semantic labels as content and the ground-truth façade images as style to train a generative adversarial network (GAN) to synthesize pseudo façade textures (Goodfellow et al., 2020). If the wall style is too structured to be synthesized by the GAN, our method falls back on image quilting to generate more regular results (Efros and Freeman, 2001).

In summary, this paper offers the following two contributions to repairing defective textures of building façades: 1) a practical solution for texture occlusion or missing building façades through image translation from unoccluded semantic labels, and 2) a novel arbitrary label-to-image translation method with rich details and regular structures. The rest of this paper is organized as follows: Section 2 provides a brief review of related work. Section 3.1 introduces the workflow of the proposed façade image synthesis approach. Sections 3.2, 3.3, and 3.4 elaborate on the details of the proposed method. Experimental evaluations are presented in Section 4. Finally, Section 5 concludes the paper.

## 2. Related work

In the following, we only discuss the most relevant literature, including 1) texture mapping and de-occlusion, 2) semantic recovery and 3) image translation.

1) Texture mapping and de-occlusion. Significant progress has been made in bundle adjustment (Verykokou and Ioannidis, 2018) and dense image matching (Hirschmuller, 2007; Hu et al., 2016), enabling realistic urban modeling. To further advance applications such as autopilot, researchers have obtained monolithic building models using parametric (Kelly et al., 2017, 2018) or interactive modeling (Vanegas et al., 2012; Kelly and Wonka, 2011) approaches. These models achieve realism through texture mapping of images acquired from different platforms. The concept of texture mapping was first proposed by Catmull (1974), and subsequent studies by Sinha et al. (2008) and Tan et al. (2008) realized texture mapping for 3D models through interactive approaches. One approach to automated texture mapping is to use a Markov Random Field (MRF) energy function to select optimal images for each facet, as demonstrated by Lempitsky and Ivanov (2007) and Waechter et al. (2014). Gal et al. (2010) further improved upon this method by introducing clarity measurement and translation vectors to achieve a smoother textured mesh. However, these single-view-based methods have limitations in terms of texture alignment with geometric structures. In contrast, multi-view-based methods fuse multiple images to obtain more consistent textures for each facet (Callieri et al., 2008; Grammatikopoulos et al., 2007). However, these methods have strict requirements on reconstruction accuracy and image resolution, and their low efficiency limits their widespread application (Waechter et al., 2014).

Removing occlusions for 3D models with discontinuous textures is a challenging problem. While Grammatikopoulos et al. (2007) addressed this issue by automatically filtering out texture outliers using statistical tests, and Yu et al. (2019); Yang et al. (2021) used deep-learning-based target detection methods to detect occlusions and remap textures to eliminate some of them, these methods cannot solve the inherent problem embedded in the texture mapping pipeline of multi-view images, which lacks imagination for invisible areas (Zhu et al., 2021b). Our previous work (Zhu et al., 2021b) successfully solved this problem for road areas using offscreen rendering and image completion, but this method has limitations when applied to building façades. Specifically, it cannot generate highly structured textures like those found in building façades and lacks generalization ability for different architectural styles, e.g., Bauhaus and Baroque (Zhu et al., 2021b). To overcome these limitations, we propose a practical approach to de-occlusion by synthesizing texture from semantic labels and ground-truth images using GAN (Goodfellow et al., 2020). Here, semantic labels control the structure of synthetic content, while ground-truth images specify the style (Isola et al., 2017).

2) Semantic recovery. Building façades always exhibit a highly structured character, which was exploited by Stiny (1975); Ripperda and Brenner (2009) to parse building façades for reconstruction. Koutsourakis et al. (2009) proposed a façade parsing method guided by MRF, and subsequent work by Teboul et al. (2011) and Cao et al. (2017) improved upon this method to achieve better results with higher efficiency. While this approach is less susceptible to occlusion, it is also complex and difficult to apply to different architectural styles. To address this limitation, researchers have turned to supervised learning methods for façade parsing using labeled data (Martinovic and Van Gool, 2013; Gadde et al., 2017; Dehbi et al., 2017). However, these methods struggle with the complexity of façades with varying styles. To alleviate this issue, researchers have explored the use of repeating patterns in buildings (Müller et al., 2007; Friedman and Stamos, 2012; Zhang et al., 2013; Fan et al., 2014; Cohen et al., 2017).

In recent years, deep learning-based object detection has revolutionized the field of computer vision, with R-CNN (Girshick et al., 2014) being a pioneering method. Mask R-CNN (He et al., 2017) extended Faster R-CNN (Ren et al., 2015) by adding a mask prediction branch, enabling semantic segmentation. While these methods have shown promising results in building façade annotation, they are still limited by occlusions. To address this issue, Lin et al. (2019) and Hu et al. (2020) utilized multi-source data, such as infrared or panorama data, to alleviate occlusion using semantic segmentation methods. However, the cost of acquiring multi-source data is often high, and there may still be areas that are inaccessible to sensors.

We have observed that most de-occlusion requirements are for reconstructed textured models. Semantic labels can be easily and accurately captured by existing methods or through manual annotation, and the semantics of occlusions are easy to determine. Therefore, we propose generating masks based on label maps, and then recovering the masked regions using image completion to achieve the goal of de-occlusion.

Image inpainting and completion techniques are commonly used to fill in missing or undesirable parts of an image with plausible pixels. Traditional methods for this task include partial differential equations (PDEs) and sampling-based approaches (Bertalmio et al., 2000; Criminisi et al., 2004). PDE-based approaches lack attention to global information and cannot fill large holes (Zhu et al., 2021b). On the other hand, sampling-based approaches fill in the void regions by using global similar patches, which are translated and rotated, and can repair large regions. The patch matching algorithm proposed by Barnes et al. (2009) significantly accelerated the search for similar patches and made sampling-based methods state-of-the-art. Subsequently, He and Sun (2012), Huang et al. (2014), and Zhu et al. (2021b) improved the patch-match based methods by incorporating offsets statistics, affine deformation, and linear patterns, respectively. While GAN-based approaches have shown impressive results on benchmark datasets, they rely on massive labeled data such as the ffhq-dataset (Karras et al., 2019). However, unlike other easily accessible datasets, building façade images require complex processing such as geometric deformation correction, making it difficult to obtain massive training data (Zhu et al., 2020a). Moreover, the complexity of semantic labels for building façades is much lower than that of photorealistic images. Therefore, in this paper, we choose the patch-match based algorithm for the semantic recovery of occluded regions.

3) Image translation. GANs employ an adversarial strategy to train two networks: a generator that simulates the probability distribution of training data from random signals and a discriminator that discerns whether the generated samples are real or fake (Goodfellow et al., 2020). Unlike CNNs, GANs use a zero-sum game between the generator and discriminator to reach Nash equilibrium, thereby enhancing the generator's ability (Goodfellow et al., 2020). Vanilla GAN generates sharper samples from random signals than Variational Auto Encoder (VAE) (Pu et al., 2016), paving the way for a new image synthesis approach (Goodfellow et al., 2020). DCGAN introduces CNN structure for stable training (Radford et al., 2015). WGAN replaces the Kullback-Leibler (KL) and Jensen-Shannon (JS) divergence with Wasserstein distance for measurement, solving the vanishing gradient problem (Arjovsky et al., 2017). LSGAN sets the objective function to the squared difference form, resulting in a more stable training process and better results (Mao et al., 2017). PGAN adopts a progressive training strategy to generate higher-resolution images (Karras et al., 2017).

Controlling the inference process can be challenging since GANs generate results from random input signals. Conditional GANs, designed to address this issue, modify the random inputs into conditional maps (Mirza and Osindero, 2014). In this paper, semantic labels are used as the conditional maps. Pix2Pix introduces PatchGAN, which evaluates generated results with a patch-based discriminator (Isola et al., 2017). Additionally, Isola et al. (2017) designs a U-Net (Ronneberger et al., 2015) generator, improving performance on several benchmark datasets. Pix2PixHD (Wang et al., 2018) builds on Pix2Pix by designing a multi-scale network to generate higher resolution images and adding perceptual (Johnson et al., 2016) and feature matching (Wang et al., 2018) loss to control the style of generated images by learning in the latent space. The network is also trained with boundary maps to obtain clearer results.

For generating realistic and plausible façade images, consistency with real images at a high level and accurate capture of low-level details are crucial. Style transfer is useful in this context as it adapts the style of a content image to match another image's style, acting as a form of domain adaptation for a single image (Jing et al., 2019). Early style transfer methods, such as the one proposed by Gatys et al. (2016), used deep features extracted through a DCNN and represented content and style using the Mean Squared Error (MSE) of the feature map and its Gram Matrix. However, this method only supported single style transfer. To address this limitation, StyleBank (Chen et al., 2017) was developed to train multiple styles simultaneously, but it still couldn't transfer inputs to arbitrary styles outside the training dataset. Instance Normalization (IN) (Ulyanov et al., 2016) and Adaptive Instance Normalization (AdaIN) (Huang and Belongie, 2017) were introduced to enable arbitrary style transfer. StyleGAN (Karras et al., 2019) built on this idea by injecting style information as AdaIN into the network, achieving impressive results on face datasets. Another approach by Li et al. (2018) used a whitening and coloring transform for arbitrary style transfer, operating in the feature space extracted by a pre-training network (encoder) and requiring only a few reconstruction networks (decoder) for good transfer results.

In recent years, the combination of image translation and style transfer has seen remarkable results. SPADE (SPAtially-Adaptive DE-normalization) (Park et al., 2019) treated style transfer as a process of de-normalization and designed a SPADE ResBlk module to replace the ResBlk in Pix2PixHD (Wang et al., 2018), injecting style information extracted by a trainable encoder. SEAN (Zhu et al., 2020) improved SPADE by performing style transfer separately for different classes, achieving state-of-the-art performance. In this paper, we consider frequency and regular

structure information to enhance the texture details and regularity of SEAN, meeting the façade texture mapping needs of photorealistic building models. However, GANs may not be effective on all datasets due to their data-based implicit probability density estimation nature. To address this issue, the proposed algorithm falls back to image quilting in regions annotated as walls when deep learning methods are ineffective (Efros and Freeman, 2001).

# 3. Structured realistic image synthesis method for building façades

- 3.1. Overview and problem setup
- 3.1.1. Overview of the approach



Figure 2: Workflow of proposed approach. Rectangles with different colors denote different phases, i.e. training or inference, and black rectangles denote the shared parts of these phases.

To address the challenges of building façade texture occlusion and missing problems in builtup areas, we develop a label-to-image deep neural network that utilizes the ground-truth image as an additional input for style. This approach enables the generation of realistic façade texture images from complete semantic labels. In addition to employing universal losses, such as GAN loss, L1 loss, and perceptual loss, we propose regularity loss and detail loss from multiple domains to enhance the regularity and detail of the results. Furthermore, we apply image completion to recover occluded semantics and use image quilting when the generation capacity of the trained network is insufficient. The overall workflow, consisting of training and inference phases, is illustrated in Figure 2.

Training. During each step of the training phase, we use a ground-truth image and its corresponding semantic label as inputs, providing style and content information, respectively. A trainable style encoder processes these inputs, extracting style vectors that represent different label styles. We then feed the semantic label and style vectors into a trainable image generator to synthesize a stylized image. To enhance the regularity and detail of the generated images, we transform both input and output images into frequency maps, spectrum maps, and corner maps. We calculate regularity loss and detail loss separately based on these maps during the back-propagation optimization. As in vanilla GANs, we also train a discriminator to compete with the generator (Goodfellow et al., 2020).

*Inference.* During the inference phase, our approach can address both texture occlusion and missing problems, with the primary differences arising from the acquisition of semantic labels fed into the image generator. For the occlusion problem, we identify occluded regions based on semantics and apply an image completion algorithm to recover semantics. For the missing

problem, we manually provide a semantic label map to specify the content of the synthetic image. In this paper, the image used to specify style is set to an occluded façade image for the occlusion problem and a desirable style image for the missing problem. Following a similar process to the first half of the training phase, we feed de-occluded or manually annotated semantic maps and style vectors into the trained image generator to synthesize realistic façade textures with the actual style. Finally, if the generated result is not satisfactory, our approach employs image quilting to composite wall textures, improving the overall quality of the final result.

#### 3.1.2. Problem setup

As illustrated in Figure 2, our approach consists of three trainable components: the style encoder E, the generator G, and the discriminator D. The overall objective of our study is more formally presented in Equation 1.

$$\min \mathcal{V}(\mathbf{R}', \mathbf{R}) \tag{1}$$

where  $\mathcal{V}$  denotes the measurement between two samples.  $\mathbf{R}' \in \mathbb{R}^{H \times W \times 3}$  and  $\mathbf{R} \in \mathbb{R}^{H \times W \times 3}$  denote the generated image and ground-truth image respectively.  $\mathbf{R}'$  can be expressed as  $G(\mathbf{S}, \mathbf{M}')$ , where  $\mathbf{S} \in \mathbb{R}^{H \times W \times 3 \times N}$  is the style vectors calculated by equation  $\mathbf{S} = E(\mathbf{R}, \mathbf{M})$ . In  $E(\mathbf{R}, \mathbf{M})$ ,  $\mathbf{R}$  and  $\mathbf{M} \in \mathbb{B}^{H \times W \times 3}$  are the most primitive inputs of network, specifically, the ground-truth sample and its corresponding semantic label.  $\mathbf{M}' \in \mathbb{B}^{H \times W \times 3}$  is the inputted semantic label map of generator  $G(\mathbf{M}' = \mathbf{M}$  in training phase).

#### 3.2. Direction guided semantics completion

Due to the difficulty in obtaining structured results through direct image completion, the semantic label map of a building façade, particularly manually labeled semantic maps, provides a highly structured alternative that can be easily used to synthesize a desirable façade image. However, to achieve de-occlusion through image translation from the label map, the first challenge is semantics recovery. We identify the occluded regions in the label map based on predefined semantics, and then intuitively employ an image restoration method for semantics completion. As label maps are simpler than natural images and data-driven methods are excessive, we opt for a patch-based algorithm to accomplish our image restoration goal. The pipeline of the proposed semantics completion is illustrated in Figure 3.



Figure 3: Workflow of direction guided semantics completion.

Our objective for semantics completion is to recover the pixels of the void region in  $\mathcal{R}_m$ , which is generated from occluded or missing semantics. The algorithm establishes a pyramid to

progressively complete the image by searching for the best nearest neighbor field (NNF)  $\mathcal{N}$ . The algorithm adopts a scanline-based expansion strategy in the patch match of each level (Barnes et al., 2009). For each pixel p in  $\mathcal{R}_m$ , we optimize the NNF  $\mathcal{N}(p)$  using an improved strategy guided by the direction  $\Theta$  (Zhu et al., 2021b).

## 3.2.1. Similarity measure of patches

For efficiency without sacrificing accuracy, the vanilla patch-match algorithm compares the similarity between the current pixel p and its four neighbors at offset  $v = \mathcal{N}(p)$ . More specifically,

$$T(\boldsymbol{p}) = \{\mathcal{R}_c(\boldsymbol{p} + \boldsymbol{s}) | \boldsymbol{s} \in [-\frac{W}{2}, \frac{W}{2}] \times [-\frac{W}{2}, \frac{W}{2}] \}$$
  

$$S(\boldsymbol{p}, \boldsymbol{v}) = T(\boldsymbol{p} + \boldsymbol{v})$$
(2)

where T and S are pixel patches centered on p and p + v in the target and source domains, respectively. The target domain and the source domain correspond to the void and known regions in  $\mathcal{R}_m$ , as well as the occluded and unoccluded regions in  $\mathcal{R}$ . During the random expansion introduced in 3.2.2, v can be either  $\mathcal{N}(p)$  or its four neighbors  $\mathcal{N}_4(p)$ . The input label map size in this paper is 256 × 256, and in order to balance effect and efficiency, we set the patch-size W to 7. Based on our previous research, the measurement E of pixel similarity in this paper is shown in equation 3.

$$E = E_a + \lambda_1 E_p + \lambda_2 E_d \tag{3}$$

The terms  $E_a$ ,  $E_p$ , and  $E_d$  represent appearance, proximity, and direction costs, respectively, and will be detailed in the following paragraphs. Additionally,  $\lambda_1$  and  $\lambda_2$  are empirically set to  $5 \times 10^{-4}$  and 0.5, respectively.

1) Appearance cost. We measure the appearance similarity between patches using the following equation 4 which computes a Gaussian weighted sum of the absolute value difference. Here,  $w_i$  is an isotropic weight generated from a Gaussian kernel (Huang et al., 2014).

$$E_a(\boldsymbol{p}, \boldsymbol{v}) = \sum_i w_i |T_i(\boldsymbol{p}) - T_i(\boldsymbol{p} + \boldsymbol{v}))|$$
(4)

2) Proximity cost. Researchers have shown that better pixels tend to appear in closer patches. Proximity cost is used to penalize the selection of nearby pixels and is shown in equation 5.

$$E_p(\boldsymbol{p}, \boldsymbol{v}) = \frac{||\boldsymbol{v}||^2}{\sigma_d(\boldsymbol{p})^2 + \sigma_c^2}$$
(5)

where  $\sigma_d(\cdot)$  calculate the minimum distance from the current pixel to the void region boundary, and  $\sigma_c = \max(w, h)/8$  (Zhu et al., 2021b).

3) Direction cost. The building façade textures are corrected orthophotos with horizontally or vertically distributed components. Thus, we can utilize this pattern to evaluate the selected pixels for better results. This can be formulated as the following equation:

$$E_r(\boldsymbol{v}) = \min_{\boldsymbol{\theta} \in \Theta} \cos(\boldsymbol{\theta}_{\boldsymbol{v}} - \boldsymbol{\theta}) \tag{6}$$

where  $\theta_{v}$  denotes the direction of current offset v,  $\theta$  is the element of  $\Theta = \{\pi/2, \pi\}$ .

#### 3.2.2. Direction guided expansion

The optimization of the NNF  $\mathcal{N}(\boldsymbol{p})$  is an iterative process of random expansion. In a single iteration, for every pixel  $\boldsymbol{p} \in \Omega$ , we compare  $E(\boldsymbol{p}, \mathcal{N}(\boldsymbol{p}))$  to  $E(\boldsymbol{p}, \mathcal{N}(\boldsymbol{q}))$  twice. The first time,  $\boldsymbol{q}$  comes from  $N_4(\boldsymbol{p})$ , which are the four neighbors of the current pixel  $\boldsymbol{p}$ . Subsequently,  $\boldsymbol{q}$ comes from a set of random pixels  $R(\boldsymbol{p})$ . The elements of  $R(\boldsymbol{p})$  are selected from the pixels distributed in the known regions within a radius r centered on  $\boldsymbol{p}$ . The value of r is initialized with  $\max(w, h)$  and is halved per iteration until it reaches 1 or other predetermined stopping conditions. Additionally, we constrain the areas of random expansion by the rectangular buffer generated along the direction  $\Theta$ , as we did in our previous work (Zhu et al., 2021b).

#### 3.3. Regularity-aware multi-domain universal image translation

Buildings are artificial objects, and their façades have regularly distributed components. While previous data-driven deep learning image synthesis or translation methods have made great progress on many datasets, these methods usually aim at non-artificial images, such as faces and natural scenes, which do not require clear boundaries (Karras et al., 2019). They have limitations on images with highly structured, man-made objects, including building façades. Additionally, existing methods cannot preserve enough high-frequency details on synthesized results, which affects the expression of architectural style and realism (Cai et al., 2021). Therefore, to address these issues, we use frequency and corner information from the pixel and spectral domains to improve the synthesized results. Moreover, due to the different styles of buildings, this paper also utilizes a style encoder to specify the style and embeds the SEAN module (Zhu et al., 2020) in the generator to achieve the façade texture translation of arbitrary style buildings from a semantic label. The objective of our network can be formally summarized by the following equations.

$$\min_{\boldsymbol{E},\boldsymbol{G}} \max_{\boldsymbol{D}} \Upsilon(\boldsymbol{E},\boldsymbol{D},\boldsymbol{G}) \tag{7}$$

where E, D, G are trainable style encoder, discriminator and generator respectively.  $\Upsilon$  is the measurement of difference between ground-truth images and synthetic images. The meaning of  $\Upsilon$  is shown in Equation 8.

$$\Upsilon(E, D, G) = \mathcal{V}_{GAN} + \mathcal{V}_{detail} + \mathcal{V}_{regularity} \tag{8}$$

 $\Upsilon$  comprises of  $\mathcal{V}_{GAN}$ ,  $\mathcal{V}_{detail}$ , and  $\mathcal{V}_{regularity}$ , which measure the dissimilarity between the synthesized image and the ground-truth image from different perspectives.  $\mathcal{V}_{GAN}$  denotes the conventional GAN measurement, while  $\mathcal{V}_{detail}$  and  $\mathcal{V}_{regularity}$  enhance the synthetic results by enriching details and improving regularity, respectively. Their specific forms are shown in Equation 9.

$$\begin{cases} \mathcal{V}_{detail} = \mathcal{V}\left(\mathcal{F}\left(\mathbf{R}'\right), \mathcal{F}(\mathbf{R})\right) \\ \mathcal{V}_{regularity} = \mathcal{V}\left(\mathcal{C}\left(\mathbf{R}'\right), \mathcal{C}(\mathbf{R})\right) \end{cases}$$
(9)

where  $\mathcal{V}$  denotes the difference between two samples;  $\mathcal{F}(\cdot)$  and  $\mathcal{C}(\cdot)$  are frequency transformation and corner map extraction respectively.  $\mathcal{F}(\cdot)$  is calculated in multiple domains by Equation 13 (pixel domain) and 20 (spectral domain). Equations 14 to 18 are the specific calculation process of  $\mathcal{C}(\cdot)$ .

#### 3.3.1. Network architecture

Figure 4 depicts the architecture of our image translation network, which takes different inputs during training and inference phases. In the training phase, the inputs are a façade image  $\mathbf{R}$  with its corresponding semantic label map  $\mathbf{M}$ . In the inference phase, another semantic label



Figure 4: Network architecture.

map  $\mathbf{M}'$  is required to provide the synthesis content. The network comprises a trainable style encoder  $\mathbf{E}$ , generator  $\mathbf{G}$ , and discriminator  $\mathbf{D}$ .  $\mathbf{E}$  encodes the façade image into a style vector  $\mathbf{S}$ for every class based on the semantic label map.  $\mathbf{G}$  synthesizes an image from the input label map  $\mathbf{M}'$ , which is equivalent to  $\mathbf{M}$  during training.  $\mathbf{D}$  is utilized to implicitly evaluate the synthesized result. The structure of the generator and discriminator is shown in Figure 4. The employed generator, with residual blocks, and the multi-scale discriminator follow the SPADE (Park et al., 2019) and Pix2PixHD (Wang et al., 2018) architectures, respectively. SPADE manipulates style transfer at the feature space, considering the statistical characteristics of the feature map as style and the normalized feature map as content (Park et al., 2019). Specifically, this approach first encodes the style image to a vector and then uses its statistical characteristics to de-normalize the input Gaussian-distributed vector of the generator, achieving state-of-the-art synthesized results with given styles. Our paper adopts this method to achieve universal image translation.

Semantic de-normalization. As shown in Figure 4, we employ the SEAN module (Zhu et al., 2020) which is an improved version of SPADE (Park et al., 2019). The structure of the SEAN module is shown in Figure 5. It embeds style and semantics in the generator through de-normalization operation. As depicted in Figure 5, the SEAN module consists of two parts. The upper part embeds image style into the network to generate images of the same style, while the lower part uses the SPADE structure to embed semantic information into the generation network to improve image generation quality (Park et al., 2019). Specifically, the activation value at position  $(n \in N, c \in C, y \in H, x \in W)$  is calculated by Equation 10:

$$\gamma_{c,y,x}(\mathbf{S},\mathbf{M})\frac{h_{n,c,y,x}-\mu_c}{\sigma_c} + \beta_{c,y,x}(\mathbf{S},\mathbf{M})$$
(10)

where h is the activation value before normalization,  $\mu$  and  $\sigma$  are the mean and variance of the activation value on channel c,  $\gamma$  and  $\beta$  are calculated by Equation 11:

$$\begin{cases} \gamma_{c,y,x}(\mathbf{S}, \mathbf{M}) = \alpha_{\gamma} \gamma^{s}_{c,y,x}(\mathbf{S}) + (1 - \alpha_{\gamma}) \gamma^{o}_{c,y,x}(\mathbf{M}) \\ \beta_{c,y,x}(\mathbf{S}, \mathbf{M}) = \alpha_{\beta} \beta^{s}_{c,y,x}(\mathbf{S}) + (1 - \alpha_{\beta}) \beta^{o}_{c,y,x}(\mathbf{M}) \end{cases}$$
(11)

where  $\alpha_{\gamma}$  and  $\alpha_{\beta}$  are the trainable parameters.  $\gamma^s$ ,  $\beta^s$  are obtained from convolution with style vector **S** as input, and  $\gamma^o$ ,  $\beta^o$  are obtained from semantic label map **M** after convolution.



Figure 5: Structure of SEAN ResBlk.

Style encoder. The encoder in Figure 4 takes an input image to extract deep features using three convolution layers with leaky rectified linear unit (LReLU) activation and down-sampling. A transpose convolution layer and a Tanh activation layer are then employed to reconstruct the image. The activated feature map is then region-wise average pooled according to the semantic label map corresponding to the input style image. The output is a set of vectors, each containing style information of its corresponding semantic category, which can be utilized for style transfer.

Generator. The generator in Figure 4 employs the structure of the Pix2PixHD generator (Wang et al., 2018) and replaces the residual module with the SEAN module to achieve style embedding with semantic information. A  $3 \times 3$  convolution is used first, followed by seven SEAN ResBlks (shown in Figure 5), and a Tanh activation layer to obtain the output. Upsampling is performed before each SEAN ResBlk. The style vector extracted by the **E** encoder is injected into the first six SEAN ResBlks. The input and output of the generator **G** are the semantic label map **M** and the synthesized image **R**', respectively (Zhu et al., 2020).

*Discriminator.* To determine whether the high-resolution synthesized image is real or fake, a discriminator with a large receptive field is needed. However, using a larger convolution kernel and a deeper network will increase unnecessary computing costs and may lead to overfitting. Therefore, we employ a multi-scale discriminator with instance normalization (IN) and LReLU activation. The structure of each discriminator is shown in Figure 4. Two discriminators are used, and GAN loss is calculated by referring to PatchGAN (Wang et al., 2018), which improves the discriminator by enlarging the receptive field without increasing the network parameters. The inputs of the discriminators are the concatenation of the image and its corresponding semantic label, and the output is an estimation of the true and fake probability of the generated samples.

#### 3.3.2. Multi-domain losses

In order to solve the problem that image translation cannot retain the structural features of source domain, a frequency domain adaptation approach is proposed (Cai et al., 2021). This method can preserve the high-frequency details of the source domain for better results in several tasks. Although our paper combines image translation with style transfer, the lack of source domain structural features is still inevitable. Thus, we use the frequency domain adaptation method to calculate and extract the frequency map (gradient map) and spectrum map corresponding to the image, and retain the consistency between the synthesized result and the real

image in both pixel space and Fourier spectral space to ensure that the generated result maintains more details.

Furthermore, we observe that façade textures have obvious structural characteristics and artificial rules. The optimization of frequency domain adaptation can only enrich the texture details, and is powerless in terms of structural information. Therefore, our paper proposes a regular optimization method that utilizes corner information in the pixel domain to reflect the regular structure of building façades. This regular optimization method allows us to achieve regular optimization of image synthesis in the training process.

Overall, the multi-domain approach of our proposed method is illustrated in Figure 6. We measure and minimize the distance between synthesized images and input images in both pixel and spectral domains to optimize the network for better details and regularities in the results.



Figure 6: Multi-domain of proposed method. Pixel domain consists of original image and synthesized image, as well as their corresponding frequency maps and corner maps. Spectral domain is the spectrum maps of original image and synthesized image. The frequency maps and spectrum maps are obtained using the frequency domain adaptation method (Cai et al., 2021), while the corner maps are calculated using an improved Harris detector (Harris et al., 1988).

# 1)Pixel domain.

Frequency map. We first perform low-pass filtering on the ground-truth image and synthetic image to obtain correspondence low-frequency information expressed in pixel domain. Then, the original images are converted into grayscale images and make difference with the low-frequency images to obtain the high-frequency information expressed in pixel domain. The network weights are optimized by comparing the difference of frequency information in pixel domain between the synthesized result and the ground-truth image, so that to retain more details in generated image. More formally, this paper adopts Gaussian kernel to do low-pass filtering, the specific form is as follows:

$$k_{\sigma}[i,j] = \frac{1}{2\pi\sigma^2} e^{-\frac{1}{2}\left(\frac{i^2+j^2}{\sigma^2}\right)}$$
(12)

where [i, j] denotes the pixel location,  $\sigma$  is standard deviation of Gaussian function. By using the Gaussian low-pass filter, we can get the frequency maps corresponding to the input image that can express high and low frequency information. The specific method is shown in following equation:

$$\begin{cases} I_L = k \otimes I \\ I_H = \mathcal{G}(I) - I_L \end{cases}$$
(13)

where I,  $I_L$  and  $I_H$  denote the image and its correspondence low and high frequency map respectively. k is Gaussian kernel and  $\otimes$  is convolution operation.  $\mathcal{G}$  is the function that can convert an image from RGB color space to grayscale space.

*Corner map.* Traditional corner detection methods usually first calculate the gradient map of the horizontal and vertical direction and then determine whether the pixel is a corner point by using a threshold value. However, this process is not differentiable, which can pose a serious problem for backpropagation optimization in deep convolutional networks. To address this issue, we propose an optimized Harris detector (Harris et al., 1988) that is differentiable and suitable for use in deep learning methods.

As shown in Equation 14, we use Sobel operator  $S_x$  and  $S_y$  to extract gradient information in horizontal and vertical direction of the image I.

$$\begin{cases}
I_x = I \otimes S_x \\
I_y = I \otimes S_y
\end{cases}$$
(14)

where  $I_x$  and  $I_y$  are the gradient maps in two different directions. The product between the two directional gradients and their squares is then calculated, as shown in the following equation.

$$\begin{cases}
I_x^2 = I_x \circ I_x \\
I_y^2 = I_y \circ I_y \\
I_x I_y = I_x \circ I_y
\end{cases}$$
(15)

where  $\circ$  denotes the pixel-wise product. After that, we calculate the Gaussian weighted sum for  $I_x^2$ ,  $I_y^2$  and  $I_x I_y$ .

$$M = \sum_{(x,y)\in W} w_G(x,y) \begin{bmatrix} \mathbf{I}_x^2 & \mathbf{I}_x \mathbf{I}_y \\ \mathbf{I}_x \mathbf{I}_y & \mathbf{I}_y^2 \end{bmatrix}$$
(16)

where  $w_G$  is the window function, which is set to Gaussian kernel function in this paper. W denotes the current sliding window being processed. Subsquently, we can get corner response matrix R by Equation 17.

$$R = \det M - k(\operatorname{trace}M) \tag{17}$$

where det and trace denote the determinant and trace of matrix M. According to Harris et al. (1988), pixels with R values greater than, equal to, and less than zero are considered corner points, flat areas, and edges, respectively. In this paper, we use a differentiable rectified linear unit (ReLU) function to remove negative and zero values that are not relevant, and obtain an equivalent expression of corner information through scaling. The specific form is shown in Equation 18.

$$R^* = \omega \cdot \text{ReLU}(R) \tag{18}$$

where  $\omega = 100000$  is the scaling factor, and as shown in Figure 7,  $R^*$  is the final value matrix that represents the corner information in the pixel domain.



(a) Façade image

(b) Visualized  $R^*$ 

Figure 7: Façade images and their corresponding visualized  $R^*$ .

## 2)Spectral domain.

Spectrum map. The pixel domain information is presented in the form of coordinates, and global information cannot be considered. The calculation of each position in the Fourier spatial spectrum map needs to use the information of all pixels, so it can reflect the global features of the image. Therefore, in addition to pixel domain optimization, this paper also uses Fast Fourier Transform (FFT) to convert image I from pixel domain to Fourier spectral domain. The neural network weights are optimized by comparing the spectrum difference between the synthesized result and ground-truth façade image, thus further retaining the consistency globally. The Discrete Fourier Transform (DFT) applied to a single image I can be formally expressed as follows:

$$\mathcal{F}(\mathbf{I})(x,y) = \frac{1}{HW} \sum_{h=0}^{H-1} \sum_{w=0}^{W-1} e^{-2\pi i \cdot \frac{hx}{H}} e^{-2\pi i \cdot \frac{wy}{W}} \cdot \mathbf{I}(h,w)$$
(19)

where (x, y) denotes the pixel location in image  $I \in \mathbb{R}^{H \times W}$ . Since the values of  $\mathcal{F}(I)$  are complex numbers, we use Equation 20 to transform them into a field of real numbers.

$$\mathcal{F}^{R}(\mathbf{I})(x,y) = \log\left(1 + \sqrt{\left[\mathcal{F}_{R}(\mathbf{I})(x,y)\right]^{2}} + \sqrt{\left[\mathcal{F}_{I}(\mathbf{I})(x,y)\right]^{2}} + \epsilon\right)$$
(20)

where  $\mathcal{F}_R(I)$  and  $\mathcal{F}_I(I)$  denote the real and imaginary parts of  $\mathcal{F}(I)$ .  $\epsilon$  is an additional term of numerical stability, this paper is set to  $1 \times 10^{-8}$  (Cai et al., 2021).  $\mathcal{F}^R(I)$  is the final spectrum map of I to be compared during training.

3) Overall loss. The overall loss function can be formally expressed by following equation:

$$\min_{\mathbf{E},\mathbf{G}} \left( \left( \max_{\mathbf{D}_1,\mathbf{D}_2} \sum_{k=1,2} \mathcal{L}_{GAN} \right) + \lambda_1 \sum_{k=1,2} \mathcal{L}_{FM} + \lambda_2 \mathcal{L}_P + \lambda_3 \mathcal{L}_D + \lambda_4 \mathcal{L}_R \right)$$
(21)

where the five terms represent GAN loss, feature matching loss, perceptual loss, detail loss and regularity loss respectively. These terms are introduced in detail in the following paragraphs. **D**<sub>1</sub> and **D**<sub>2</sub> denote two PatchGAN discriminators. The weights to balance the different losses are set to  $\lambda_1 = \lambda_2 = \lambda_3 = \lambda_4 = 10$  and  $\lambda_5 = 5 \times 10^{-6}$ .

*GAN loss.* The form of the GAN loss function in this paper is shown in Equation 22. The optimization is to minimize the loss of style encoder **E** and generator **G**, and maximize the loss of discriminators  $\mathbf{D_1}, \mathbf{D_2}$ . The main objective is to train generator **G** to generate high-quality

images that can deceive discriminator **D**, making the discriminator unable to distinguish between real and synthesized images.

$$\min_{\mathbf{E},\mathbf{G}} \max_{\mathbf{D}_{1},\mathbf{D}_{2}} \sum_{k=1,2} \mathcal{L}_{GAN}\left(\mathbf{E},\mathbf{G},\mathbf{D}_{k}\right)$$
(22)

The Hinge loss function used for  $\mathcal{L}_{GAN}$  in this paper is defined as follows:

$$\mathcal{L}_{GAN} = \mathbb{E}[\max(0, 1 - \mathbf{D}_{\mathbf{k}}(\mathbf{R}, \mathbf{M}))] + \mathbb{E}[\max(0, 1 - \mathbf{D}_{\mathbf{k}}(\mathbf{G}(\mathbf{S}, \mathbf{M}), \mathbf{M}))]$$
(23)

Refer to section 3.1.2 for the meaning of parameters in the formula.

*Feature matching loss.* The feature matching loss evaluates the difference between the feature maps of the ground-truth image and the synthetic image at different levels, which are extracted from the discriminator network. Optimizing the network weights according to the feature matching loss can achieve better consistency at the feature level, resulting in better results. The definition of the feature matching loss is shown as follows:

$$\mathcal{L}_{FM} = \mathbb{E}\sum_{i=1}^{T} \frac{1}{N_i} \Big[ \|\mathbf{D}_{\mathbf{k}}^{(i)}(\mathbf{R}, \mathbf{M}) - \mathbf{D}_{\mathbf{k}}^{(i)}(\mathbf{G}(\mathbf{S}, \mathbf{M}), \mathbf{M})\|_1 \Big]$$
(24)

where  $D_k^{(i)}$  and  $N_i$  denote the feature map outputted from the i-th layer of the k-th discriminator and the number of corresponding elements. T is the total layer numbers of discriminator (Mirza and Osindero, 2014).

*Perceptual loss.* Perceptual loss measures the difference between two images by comparing the features extracted from the network pre-trained on large datasets, e.g. pre-trained VGG-16 is used in this paper. Specifically, perceptual loss can be calculated using the following equation:

$$\mathcal{L}_P = \mathbb{E}\sum_{i=1}^N \frac{1}{M_i} \left[ \left\| F^{(i)}(\mathbf{R}) - F^{(i)}(\mathbf{G}(\mathbf{S}, \mathbf{M})) \right\|_1 \right]$$
(25)

where N is the layer number of pre-trained network,  $F^{(i)}$  and  $M_i$  denote the feature map extracted from i-th layer of pre-trained network and the number of corresponding elements (Johnson et al., 2016).

Detail loss. Detail loss is shown in Equation 3, and consists of two parts, i.e.,  $\mathcal{L}_{D_{pix}}$  and  $\mathcal{L}_{D_{spectral}}$ .

$$\mathcal{L}_D = \mathcal{L}_{D_{pix}} + \mathcal{L}_{D_{spectral}} \tag{26}$$

where  $\mathcal{L}_{D_{pix}}$  and  $\mathcal{L}_{D_{spectral}}$  denote the different detail losses, they come from frequency information in pixel and spectral domain respectively, they are calculated by following equations:

$$\mathcal{L}_{D_{pix}} = \mathbb{E}\left[\left\|\mathbf{R}_{L} - (\mathbf{G}(\mathbf{S}, \mathbf{M})_{L})\right\|_{1} + \left\|\mathbf{R}_{H} - (\mathbf{G}(\mathbf{S}, \mathbf{M})_{H})\right\|_{1}\right]$$
(27)

where  $\mathbf{R}_L$  and  $\mathbf{R}_H$  denote the low and high frequency map of ground-truth image,  $\mathbf{G}(\mathbf{S}, \mathbf{M})_L$  and  $\mathbf{G}(\mathbf{S}, \mathbf{M})_H$  are the low and high frequency map of generated image. Frequency map is obtained from Equation 13.

$$\mathcal{L}_{D_{spectral}} = \mathbb{E}\left[\left\|\mathcal{F}^{R}(\mathbf{R}) - \mathcal{F}^{R}(\mathbf{G}(\mathbf{S}, \mathbf{M}))\right\|_{1}\right]$$
(28)

where  $\mathcal{F}^R$  refers to the FFT and real number operation using Equation 19 and 20.

*Regularity loss.* The regularity loss is obtained by calculating the difference of corner maps between the generated image and the ground-truth image. The specific form is shown in Equation 29.

$$\mathcal{L}_{\mathbf{R}} = \mathbb{E}\left[\|\mathcal{C}(\mathbf{R}) - \mathcal{C}(\mathbf{G}(\mathbf{S}, \mathbf{M}))\|_{1}\right]$$
(29)

where  $\mathcal{C}(\cdot)$  is the corner map calculation introduced in paragraph 3.3.2.

## 3.4. Image quilting

Because deep learning-based image synthesis methods rely on massive training data, they may not be able to generalize well to data that has not been seen. While embedding image styles can alleviate this problem, there is still a significant gap between the generated results and actual images, particularly for highly structured repetitive textures. Therefore, when the generation ability of deep learning methods is insufficient, this paper applies a sample-based image quilting algorithm to the wall region to compensate for this limitation. This process can be expressed as follows:

$$\begin{cases} \mathbf{R}' = \mathbf{R}' & \mathcal{V}(\mathbf{R}', \mathbf{R}) \leq threshold \\ \mathbf{R}' = \mathbf{R}' \otimes \neg \mathbf{m}' \oplus \mathcal{Q}(\mathbf{R}') \otimes \mathbf{m}' & \mathcal{V}(\mathbf{R}', \mathbf{R}) > threshold \end{cases}$$
(30)

 $\mathbf{R}'$  denotes the synthetic image of generator  $\mathbf{G}$  on the right side of the equation, and the final result on the left side of the equation.  $\mathbf{m}' \in \mathbb{B}^{H \times W}$  is a subset of  $\mathbf{M}'$  which specifies the wall region in synthetic image.  $\mathcal{Q}$  denotes image quilting which generates repetitive textures by stitching together small patches of un-occluded regions (Efros and Freeman, 2001). Next, we will provide a detailed introduction to the image quilting algorithm  $\mathcal{Q}$ .

We begin the image quilting process by selecting a source pixel patch S from the synthetic image  $\mathbf{R}'$ . During image quilting, all pixels are sampled from S. Firstly, we sample all  $N \times N$ pixel patches  $S(\mathbf{p})$  from S to form the set  $S_B$ . Next, we randomly select a patch  $\mathcal{R}(S_B)$  and define this patch as  $B_1$ . Subsequently, the algorithm iterates by a sequence Q which contains the center position  $\mathbf{p}$  of all patches in raster scanning order. The moving step during scanning is  $\mathbf{v}$ , which denotes the overlapping width. In this paper, the offset  $\mathbf{v}$  is set to (5,0) in the first row, (0,5) in the first column, and (5,5) for the rest of the positions. During iteration, we first find pixel patch  $B_2$  from  $S_B$  by measuring its similarity with  $B_1$ . The overlapping regions in  $B_1$  and  $B_2$  are represented by  $B_{ov}^1$  and  $B_{ov}^2$  respectively. Next, we find a boundary  $\mathbf{l}$  around the center line of the overlapping region by using the Dijkstra algorithm  $\mathcal{D}$ , based on the error  $\mathcal{E}(B_{ov}^1 - B_{ov}^2)^2 < \mathbf{t}$ , where  $\mathbf{t} = 0.1$  denotes the minimum tolerance. After that, we update the pixel patches  $\mathcal{T}(\mathbf{p})$  and  $\mathcal{T}(\mathbf{p} + \mathbf{v})$  by stitching  $B_1$  and  $B_2$  using the minimum error boundary  $\mathbf{l}$ . After the iteration process, we obtain the texture synthesized result  $\mathcal{T}$  of image quilting.

orithm 1 Image quilting algorithm.	
procedure ImageQuilting $(\mathcal{S}, \mathcal{T}, Q, \boldsymbol{t}, \boldsymbol{v})$	
$\mathcal{S}_B \leftarrow \{\mathcal{S}(oldsymbol{p}), orall oldsymbol{p} \in \mathcal{S}\}$	$\triangleright$ Random Initialization
$B_1 \leftarrow \mathcal{R}(\mathcal{S}_B)$	
$\mathbf{for}\; \boldsymbol{p} \in Q\; \mathbf{do}$	▷ Raster Scanning
$B_2 \leftarrow \mathcal{V}(B_1, \mathcal{S}_B)$	
$B_{ov}^1, B_{ov}^2 \leftarrow B_1 \wedge B_2$	
while $\mathcal{E}(B_{ov}^1 - B_{ov}^2)^2 < t$ do	$\triangleright$ Minimum Error Boundary Cut
$oldsymbol{l} \leftarrow \mathcal{D}(B^1_{ov}, B^2_{ov})$	
end while	
$\mathcal{T}(oldsymbol{p}), \mathcal{T}(oldsymbol{p}+oldsymbol{v}) \leftarrow \mathcal{Q}(B_1, B_2, oldsymbol{l})$	$\triangleright$ Pixel Patch Quilting
$B_1 \leftarrow B_2$	
end for	
nd procedure	
	$ \begin{array}{l} \begin{array}{l} \mbox{rithm 1 Image quilting algorithm.} \\ \hline \mbox{rocedure IMAGEQUILTING}(\mathcal{S},\mathcal{T},Q,t,v) \\ \mathcal{S}_B \leftarrow \{\mathcal{S}(p), \forall p \in \mathcal{S}\} \\ B_1 \leftarrow \mathcal{R}(\mathcal{S}_B) \\ \mbox{for } p \in Q \ \mbox{do} \\ B_2 \leftarrow \mathcal{V}(B_1,\mathcal{S}_B) \\ B_{ov}^1, B_{ov}^2 \leftarrow B_1 \wedge B_2 \\ \mbox{while } \mathcal{E}(B_{ov}^1 - B_{ov}^2)^2 < t \ \mbox{do} \\ l \leftarrow \mathcal{D}(B_{ov}^1, B_{ov}^2) \\ \mbox{end while} \\ \mathcal{T}(p), \mathcal{T}(p+v) \leftarrow \mathcal{Q}(B_1, B_2, l) \\ B_1 \leftarrow B_2 \\ \mbox{end for } \\ \mbox{nd procedure} \end{array} $

## 4. Experimental evaluation and analysis

# 4.1. Dataset description

To verify the effectiveness of the proposed approach in this paper, we used two public urban textured 3D building datasets in .skp format from different countries as shown in Figure 8. The first dataset was collected from central London, including several landmarks such as London Bridge. We selected some buildings with missing texture in this dataset for the experiment. The second dataset is from Pipitea South of Wellington, New Zealand, and we selected a typical block (Wellington train station) consisting of a series of buildings with occluded façade textures to evaluate our approach. Moreover, we trained our proposed network using the Large Scale Architectural Asset (LSAA) dataset for performance evaluation. The LSAA dataset contains 199,723 façade images of different styles extracted from large-scale rectified panoramas (Zhu et al., 2020).



(a) London dataset

(b) Wellington dataset

# 4.2. Results

#### 4.2.1. Semantics completion of building façade

To solve the problem caused by frequent occlusions of building façade in built-up areas, we first recover the semantics of occluded regions. Figure 9 and 10 show two different experiments to validate our semantics completion method. As shown in Figure 9, we first consider non-occluded

Figure 8: Datasets for experiment. London dataset is rendered in SketchUp with its own base map, and Wellington dataset is rendered in ArcGIS with public world imagery base map. They both are skp format which is modeled by SketchUp. The textures and base map of London dataset are the manual production of model creator. Due to the mutual occlusions between buildings, some building façades have no texture. Unlike London dataset, the textures of Wellington dataset come from realistic scene, that contain several inevitable occlusions.

façade images and their correspondence semantics as reference, then manually or randomly remove some regions of semantic label maps, finally adopt proposed method to recover the missing semantics. It can be seen that, except for the first group of experiments in the third row, the other experiments have obtained desirable results compared to the reference.



(a) Ground-truth

(b) Incomplete semantic map & semantics completion results

Figure 9: Hypothetical experiment of proposed semantics completion method. (a) Ground-truth non-occluded façde images and their correspondence semantic label maps. (b) The hypothetical incomplete images and semantics completion results recovered by proposed method. Different colors in (b) denote different semantics, i.e. window, wall and door.

Figure 10 shows experiments on building façades with occluded realistic textures, and reasonable results are obtained except for Figure 10(b). By combining the experiments shown in Figure 9 and 10, the proposed method demonstrates good performance on small occlusions. However, when the occlusion area is too large, achieving satisfactory recovery results may be challenging. Fortunately, by using the approach proposed in this paper, people can also manually annotate the semantics of building façades for final texture synthesis. Moreover, the colors of different semantics are predetermined, so obtaining the label map with standard colors through simple correction after image completion is possible.

![](_page_17_Figure_6.jpeg)

(c) Building façade in London

(d) Building façade in Australia

Figure 10: Semantics completion experiment on building façades with occluded realistic textures. In each subfigure, the occlusion texture and its corresponding semantic label image, mask image and recovery result are displayed from left to right. Different colors denote different semantics in label map, including window, wall, door, vegetation and cornice.

## 4.2.2. Façade texture repair of 3D building

Figure 11 showcases the building models before and after applying the proposed repair method. The details of the repaired building façades are presented in Figure 12 and 13. For instance, the Wellington train station building façade textures have numerous occlusions, resulting in the model having visible defects. To tackle this problem, we first annotate the semantics of the occluded regions and employ the proposed semantics completion method to repair the semantics of those regions. We then manually annotate the semantics of façades with large occlusions, which can hardly yield satisfactory semantics completion results. Afterwards, we feed the labels recovered manually and automatically and style images with their corresponding semantic labels (12 (a) set in this paper) to the proposed GAN network to synthesize plausible textures. Finally, we map the synthetic images to the respective façades to obtain the final repaired model. Moreover, Figure 13 demonstrates the effectiveness of our approach in handling missing textures. We select several buildings in the London dataset with incomplete fadade textures and apply the same process as the Wellington experiment. This experiment also highlights the repair ability of the proposed approach on different architectural styles.

![](_page_18_Picture_2.jpeg)

(c) Unrepaired London model

(d) Repaired London model

Figure 11: Repaired results of Wellington and London model. (a) Original model of Wellington train station, it has many occlusions on several building façades. (b) The Wellington train station model after repaired, we de-occluded most occlusions by proposed approach. (c) Original building model which has many missing textures of London. (d) Repaired London model with complete textures synthesized by proposed approach.

![](_page_19_Figure_0.jpeg)

Figure 12: Details of Wellington train station model experiment. All sub-experiments consider (a) as style input. Besides, (a) and (e) use semantics completion results as input labels, others use manual annotated labels as inputs.

![](_page_19_Figure_2.jpeg)

Figure 13: Details of London model experiment.

# 4.3. Comparison of image translation

# 4.3.1. Qualitative comparison

To evaluate the proposed image translation method, we adopt SPADE, SEAN and our method on LSAA dataset to do a series of experiments. Figure 14 shows some synthesized results

of different methods. The first row in every sub-figure denotes the inputs, that specify style and content to synthetic images. The next three rows are the synthesized results of different methods mentioned above, and the red rectangles indicate the obviously differences between different methods. SPADE only considers semantic label map as input and cannot specify the style of result, consequently it is far different from ground-truth. Comparing the results of SEAN and proposed method, it can be seen from Figure 14 (a), (b), (c), (d) that the proposed method has better texture details than SEAN, especially in Figure 14 (c) and (d), the results of proposed method have much more clear and regular details on balustrade areas highlighted by red rectangles. In addition to this, our method has better regular structures generation ability. More specifically, proposed method can synthesizes windows with clear regular structures in Figure 14 (e), (g), (j), (k) and (l), on the contrary, SEAN even cannot generate any structures of windows in the experiments of Figure 14 (j), (k), (l). Besides, it also can be seen from Figure 14 (f), (h) and (i), that the synthesized façade components of proposed methods are more horizontal and vertical, which is more similar to the ground-truth.

![](_page_21_Figure_0.jpeg)

Figure 14: Image translation results of different methods. For each sub-figure, from top to bottom, it is the ground-truth, results of SPADE, results of SEAN and results of proposed method. The red rectangles indicate the areas with differences for synthetic images of different methods. 22

As shown in Figure 15, proposed image translation method still doesn't have enough ability on image synthesis of unseen highly structured styles. Thus, we use image quilting algorithm to eliminate the impact of this deficiency on the results. Figure 15 (d) shows the finally result after the replacement of image quilting texture.

![](_page_22_Figure_1.jpeg)

Figure 15: Result after image quilting. (a) Inputs of GAN network, including an ground-truth image with its semantic label map, the red rectangle is the selected input of image quilting algorithm. (b) Synthesized façade image of proposed network. (c) Synthesized texture of wall region from the pixel patch highlighted in (a). (d) Finally result after replacing wall region with quilted texture.

## 4.3.2. Quantitative comparison

Peak signal-to-noise ratio (PSNR) and Structural similarity (SSIM) are wildly used indexes for image quality evaluation. PSNR can be formally expressed as Equation 31.

$$\begin{cases} \text{MSE} = \|\mathbf{R}' - \mathbf{R}\|_2^2 / N \\ \text{PSNR} = 10 \log_{10} \left( M^2 / \text{ MSE} \right) \end{cases}$$
(31)

where MSE denotes mean square error, M is the range of the data type of **R** and **R'**. SSIM can be calculated by Equation 32.

$$SSIM = \frac{1}{N} \sum_{i=1}^{N} \frac{(2\mu_{\mathbf{R}'}\mu_{\mathbf{R}} + C_1) (2\sigma_{\mathbf{R'R}} + C_2)}{(\mu_{\mathbf{R'}}^2 + \mu_{\mathbf{R}}^2 + C_1) (\sigma_{\mathbf{R'}}^2 + \sigma_{\mathbf{R}}^2 + C_2)}$$
(32)

where N is the total number of pixels of images  $\mathbf{R}'$  and  $\mathbf{R}$ .  $\mu_{\mathbf{R}'}$ ,  $\mu_{\mathbf{R}}$ ,  $\sigma_{\mathbf{R}'}$ ,  $\sigma_{\mathbf{R}}$  and  $\sigma_{\mathbf{R}'\mathbf{R}}$  are the local means, standard deviations, and cross-covariance for  $\mathbf{R}'$  and  $\mathbf{R}$ . And in this paper,  $C_1 = (0.01 \times 255)^2$ ,  $C_2 = (0.03 \times 255)^2$  and  $C_3 = (0.03 \times 255)^2/2$ . SSIM evaluates similarities between two images from luminance, contrast and structure.

PSNR and SSIM evaluate the differences between synthesized image and ground-truth image from different perspectives. The higher their values are, more similar the synthesized image is to the ground-truth image, which also means that the method is better.

Learned Perceptual Image Patch Similarity (LPIPS), also named perceptual loss, is a metric that measure the differences between two images on feature level by using pre-trained network, e.g. pre-trained AlexNet in this paper. It is more similar to human perception than PSNR and SSIM, and the lower value denotes better synthesis results.

Method Metrics	SPADE	SEAN	proposed
PSNR ↑ SSIM ↑ LPIPS ↓	13.77 0.291 0.644	$16.76 \\ 0.489 \\ 0.555$	$17.13 \\ 0.502 \\ 0.550$

Table 1: Quantitative comparison of 3 different methods on LSAA. The best results are highlighted in bold.

We trained SPADE, SEAN and proposed method on LSAA training dataset (28494 façade images with semantic labels) with  $4 \times \text{NVIDIA} 3090$ . These trained models are applied to LSAA testing dataset, which contains 1000 façade images with semantic labels, that never be seen during training. The quantitative results of the comparison are presented in Table 1, where the proposed method outperformed SPADE and SEAN in all three quality assessment metrics on the LSAA dataset.

# 4.4. Analysis of detail and regularity losses

The proposed image translation method enriches details and regularizes the structures of synthesized results by multi-domain losses. In order to prove the availability of our innovation, we conducted ablation experiment on LSAA testing dataset by setting different losses during training. Figure 16 shows several typical results synthesized by different settings of proposed method.

Focusing on the red rectangles in Figure 16, it can be concluded that both detail loss and regularity loss proposed by us are effective for the improvement of synthesis results. As we can see the vegetation regions in Figure 16 (a) and (e), detail loss can make the synthesized textures be more natural with vivid details. Comparing the balustrade areas in Figure 16 (b) and (c), we can find out that the results with detail loss have more clear structures in details. Unfortunately, its regularity maintenance capacity is still not enough for some highly structured façade components synthesis, and the regularity loss can reinforces it to a certain degree. The window frames in Figure 16 (d), (e) and (f) are not plausible enough when synthesized by the model without regularity and detail losses. However, as shown in Figure 16 (d), when detail loss is added, the results have distinct difference in window areas. In addition, Figure 16 (e), (f) also have significant improvement when combining regularity loss with detail loss.

![](_page_24_Figure_0.jpeg)

Figure 16: Ablation experiments on LSAA dataset, and the baseline is SEAN. From top to bottom in every sub-figure, it is ground-truth, w/o. both, w/o. regularity loss, w/o. detail loss, and proposed. Red rectangles indicate regions that have differences with each others.

Table 2 shows the evaluation of synthesized results under different loss settings by quality assessment indexes. It is demonstrated that detail loss and regularity loss are valid for façade images synthesis which are always highly structured, and their combination will make the results better and more realistic than their own.

Table 2: Quantitative evaluation of results under different settings on proposed method, baseline is SEAN (Zhu et al., 2020). The best results are highlighted in bold.

Method Metrics	w/o. both	w/o. detail loss	w/o. regularity loss	proposed
$\mathrm{PSNR}\uparrow$	16.76	17.02	16.96	17.13
SSIM $\uparrow$	0.489	0.500	0.496	0.502
$\mathrm{LPIPS}\downarrow$	0.555	0.551	0.553	0.550

#### 4.5. Discussion and limitations

The proposed approach for repairing defective building façade textures yields plausible results, addressing the problem of texture occlusion or missing in 3D realistic building models. The repaired models exhibit reasonable textures and can be used for model exhibition. The proposed semantics completion method also produces desirable results for occlusions, improving the automation of the texture repair approach. Moreover, the regularity-aware multi-domain universal image translation method has demonstrated the ability to synthesize more detailed and structured building façade textures through qualitative experiments. Quantitative comparisons also show the superiority of the proposed image translation method compared to others.

Despite the better results of façade texture repair and image translation, there are still limitations. The resolution and definition of style images affect the final synthesis results. For higher resolution, we can resolve this problem by simple down-sampling. However, for lower resolution, it is hard to realize a plausible result by up-sampling. In addition to this, GAN based methods are hardly to train and too computation-consuming to synthesis higher resolution images. Future research will explore the combination of image translation and super-resolution techniques.

# 5. CONCLUSION

3D building model is the fundamental of digital city, live navigation and smart driving, its realism usually comes from photogrammetric textures. However, there are inevitable occlusions in built-up areas, the vegetation near buildings also makes it difficult to acquire un-occluded façade textures by handheld camera. The above problems eventually lead to the inaccessible areas of photogrammetry, which is cannot deal with in traditional texture mapping pipeline. To solve this problem, this paper proposed a deep learning based approach to repair façade defect textures from easily accessible semantic label map. Specifically, we proposed a semantics recovery method by using image completion algorithm to improve automation of de-occlusion. A regularity-aware multi-domain universal image translation method is used to synthesize building façade textures of arbitrary styles. This method achieves better results by enriching the details and improving the regularity of synthesized images. Overall, the proposed texture repair approach not only can de-occlusion, but also can generate realistic façade textures which have actual architectural style from nothing. Future directions on the building models processing may include: (1) deocclusion of other regions by proposed approach; (2) high-resolution texture synthesis under limited computation resources; (3) deep learning based geometry generation of 3D model.

#### Acknowledgments

This work was supported in part by the National Natural Science Foundation of China (Project No. 42230102, 42071355, 41871291) and the National Key Research and Development Program of China (Project No. 2022YFF0904400).

#### References

- Arjovsky, M., Chintala, S., Bottou, L., 2017. Wasserstein generative adversarial networks, in: International conference on machine learning, PMLR. pp. 214–223.
- Barnes, C., Shechtman, E., Finkelstein, A., Goldman, D.B., 2009. Patchmatch: A randomized correspondence algorithm for structural image editing. ACM Trans. Graph. 28, 24.
- Bertalmio, M., Sapiro, G., Caselles, V., Ballester, C., 2000. Image inpainting, in: Proceedings of the 27th annual conference on Computer graphics and interactive techniques, pp. 417–424.

- Buyukdemircioglu, M., Kocaman, S., 2020. Reconstruction and efficient visualization of heterogeneous 3d city models. Remote Sensing 12, 2128.
- Cai, M., Zhang, H., Huang, H., Geng, Q., Li, Y., Huang, G., 2021. Frequency domain image translation: More photo-realistic, better identity-preserving, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 13930–13940.
- Callieri, M., Cignoni, P., Corsini, M., Scopigno, R., 2008. Masked photo blending: Mapping dense photographic data set on high-resolution sampled 3d models. Computers & Graphics 32, 464–473.
- Cao, J., Metzmacher, H., O'Donnell, J., Frisch, J., Bazjanac, V., Kobbelt, L., van Treeck, C., 2017. Facade geometry generation from low-resolution aerial photographs for building energy modeling. Building and Environment 123, 601–624.
- Catmull, E.E., 1974. A subdivision algorithm for computer display of curved surfaces. The University of Utah.
- Chen, D., Yuan, L., Liao, J., Yu, N., Hua, G., 2017. Stylebank: An explicit representation for neural image style transfer, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 1897–1906.
- Chen, S., Zhang, W., Wong, N.H., Ignatius, M., 2020. Combining citygml files and data-driven models for microclimate simulations in a tropical city. Building and Environment 185, 107314.
- Cohen, A., Oswald, M.R., Liu, Y., Pollefeys, M., 2017. Symmetry-aware facade parsing with occlusions, in: 2017 International Conference on 3D Vision (3DV), IEEE. pp. 393–401.
- Criminisi, A., Pérez, P., Toyama, K., 2004. Region filling and object removal by exemplar-based image inpainting. IEEE Transactions on image processing 13, 1200–1212.
- Dehbi, Y., Hadiji, F., Gröger, G., Kersting, K., Plümer, L., 2017. Statistical relational learning of grammar rules for 3d building reconstruction. Transactions in GIS 21, 134–150.
- Efros, A.A., Freeman, W.T., 2001. Image quilting for texture synthesis and transfer, in: Proceedings of the 28th annual conference on Computer graphics and interactive techniques, pp. 341–346.
- Fan, L., Musialski, P., Liu, L., Wonka, P., 2014. Structure completion for facade layouts. ACM Trans. Graph. 33, 210–1.
- Friedman, S., Stamos, I., 2012. Online facade reconstruction from dominant frequencies in structured point clouds, in: 2012 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops, IEEE. pp. 1–8.
- Gadde, R., Jampani, V., Marlet, R., Gehler, P.V., 2017. Efficient 2d and 3d facade segmentation using auto-context. IEEE transactions on pattern analysis and machine intelligence 40, 1273– 1280.
- Gal, R., Wexler, Y., Ofek, E., Hoppe, H., Cohen-Or, D., 2010. Seamless montage for texturing models, in: Computer Graphics Forum, Wiley Online Library. pp. 479–486.
- Gatys, L.A., Ecker, A.S., Bethge, M., 2016. Image style transfer using convolutional neural networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2414–2423.

- Girshick, R., Donahue, J., Darrell, T., Malik, J., 2014. Rich feature hierarchies for accurate object detection and semantic segmentation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 580–587.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y., 2020. Generative adversarial networks. Communications of the ACM 63, 139–144.
- Grammatikopoulos, L., Kalisperakis, I., Karras, G., Petsa, E., 2007. Automatic multi-view texture mapping of 3d surface projections, in: Proceedings of the 2nd ISPRS International Workshop 3D-ARCH, Citeseer. pp. 1–6.
- Harris, C., Stephens, M., et al., 1988. A combined corner and edge detector, in: Alvey vision conference, Citeseer. pp. 10–5244.
- He, K., Gkioxari, G., Dollár, P., Girshick, R., 2017. Mask r-cnn, in: Proceedings of the IEEE international conference on computer vision, pp. 2961–2969.
- He, K., Sun, J., 2012. Statistics of patch offsets for image completion, in: European conference on computer vision, Springer. pp. 16–29.
- Hirschmuller, H., 2007. Stereo processing by semiglobal matching and mutual information. IEEE Transactions on pattern analysis and machine intelligence 30, 328–341.
- Hu, H., Ding, Y., Zhu, Q., Wu, B., Xie, L., Chen, M., 2016. Stable least-squares matching for oblique images using bound constrained optimization and a robust loss function. ISPRS journal of photogrammetry and remote sensing 118, 53–67.
- Hu, H., Wang, L., Zhang, M., Ding, Y., Zhu, Q., 2020. Fast and regularized reconstruction of building fa\c {c} ades from street-view images using binary integer programming. arXiv preprint arXiv:2002.08549.
- Huang, J.B., Kang, S.B., Ahuja, N., Kopf, J., 2014. Image completion using planar structure guidance. ACM Transactions on graphics (TOG) 33, 1–10.
- Huang, X., Belongie, S., 2017. Arbitrary style transfer in real-time with adaptive instance normalization, in: Proceedings of the IEEE international conference on computer vision, pp. 1501–1510.
- Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A., 2017. Image-to-image translation with conditional adversarial networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 1125–1134.
- Jing, Y., Yang, Y., Feng, Z., Ye, J., Yu, Y., Song, M., 2019. Neural style transfer: A review. IEEE transactions on visualization and computer graphics 26, 3365–3385.
- Johnson, J., Alahi, A., Fei-Fei, L., 2016. Perceptual losses for real-time style transfer and superresolution, in: European conference on computer vision, Springer. pp. 694–711.
- Karras, T., Aila, T., Laine, S., Lehtinen, J., 2017. Progressive growing of gans for improved quality, stability, and variation. arXiv preprint arXiv:1710.10196.
- Karras, T., Laine, S., Aila, T., 2019. A style-based generator architecture for generative adversarial networks, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 4401–4410.

- Kelly, T., Femiani, J., Wonka, P., Mitra, N.J., 2017. Bigsur: large-scale structured urban reconstruction. ACM Transactions on Graphics 36.
- Kelly, T., Guerrero, P., Steed, A., Wonka, P., Mitra, N.J., 2018. Frankengan: guided detail synthesis for building mass-models using style-synchonized gans. arXiv preprint arXiv:1806.07179
- Kelly, T., Wonka, P., 2011. Interactive architectural modeling with procedural extrusions. ACM Transactions on Graphics (TOG) 30, 1–15.
- Koutsourakis, P., Simon, L., Teboul, O., Tziritas, G., Paragios, N., 2009. Single view reconstruction using shape grammars for urban environments, in: 2009 IEEE 12th international conference on computer vision, IEEE. pp. 1795–1802.
- Lempitsky, V., Ivanov, D., 2007. Seamless mosaicing of image-based texture maps, in: 2007 IEEE conference on computer vision and pattern recognition, IEEE. pp. 1–6.
- Li, Q., Huang, H., Yu, W., Jiang, S., 2023. Optimized views photogrammetry: Precision analysis and a large-scale case study in qingdao. IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing 16, 1144–1159.
- Li, Y., Liu, M.Y., Li, X., Yang, M.H., Kautz, J., 2018. A closed-form solution to photorealistic image stylization, in: Proceedings of the European Conference on Computer Vision (ECCV), pp. 453–468.
- Lin, D., Jarzabek-Rychard, M., Tong, X., Maas, H.G., 2019. Fusion of thermal imagery with point clouds for building façade thermal attribute mapping. ISPRS Journal of Photogrammetry and Remote Sensing 151, 162–175.
- Lin, H., Chen, M., Lu, G., Zhu, Q., Gong, J., You, X., Wen, Y., Xu, B., Hu, M., 2013. Virtual geographic environments (vges): A new generation of geographic analysis tool. Earth-Science Reviews 126, 74–84.
- Mao, X., Li, Q., Xie, H., Lau, R.Y., Wang, Z., Paul Smolley, S., 2017. Least squares generative adversarial networks, in: Proceedings of the IEEE international conference on computer vision, pp. 2794–2802.
- Martinovic, A., Van Gool, L., 2013. Bayesian grammar learning for inverse procedural modeling, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 201–208.
- Mirza, M., Osindero, S., 2014. Conditional generative adversarial nets. arXiv preprint arXiv:1411.1784.
- Müller, P., Zeng, G., Wonka, P., Van Gool, L., 2007. Image-based procedural modeling of facades. ACM Trans. Graph. 26, 85.
- Park, T., Liu, M.Y., Wang, T.C., Zhu, J.Y., 2019. Semantic image synthesis with spatiallyadaptive normalization, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 2337–2346.
- Pu, Y., Gan, Z., Henao, R., Yuan, X., Li, C., Stevens, A., Carin, L., 2016. Variational autoencoder for deep learning of images, labels and captions. Advances in neural information processing systems 29.

- Radford, A., Metz, L., Chintala, S., 2015. Unsupervised representation learning with deep convolutional generative adversarial networks. arXiv preprint arXiv:1511.06434.
- Remondino, F., Gerke, M., 2015. Oblique aerial imagery-a review, in: Photogrammetric week, pp. 75–81.
- Ren, S., He, K., Girshick, R., Sun, J., 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. Advances in neural information processing systems 28.
- Ripperda, N., Brenner, C., 2009. Application of a formal grammar to facade reconstruction in semiautomatic and automatic environments, in: Proc. of the 12th AGILE Conference on GIScience, Citeseer. pp. 1–12.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks for biomedical image segmentation, in: International Conference on Medical image computing and computerassisted intervention, Springer. pp. 234–241.
- Sinha, S.N., Steedly, D., Szeliski, R., Agrawala, M., Pollefeys, M., 2008. Interactive 3d architectural modeling from unordered photo collections. ACM Transactions on Graphics (TOG) 27, 1–10.
- Stiny, G.N., 1975. Pictorial and formal aspects of shape and shape grammars and aesthetic systems. University of California, Los Angeles.
- Tan, Y., Kwoh, L., Ong, S., 2008. Large scale texture mapping of building facades. IAPRS B37 , 687–691.
- Tao, F., Qi, Q., 2019. Make more digital twins. Nature 573, 490–491.
- Teboul, O., Kokkinos, I., Simon, L., Koutsourakis, P., Paragios, N., 2011. Shape grammar parsing via reinforcement learning, in: CVPR 2011, IEEE. pp. 2273–2280.
- Ulyanov, D., Vedaldi, A., Lempitsky, V., 2016. Instance normalization: The missing ingredient for fast stylization. arXiv preprint arXiv:1607.08022 .
- Vanegas, C.A., Kelly, T., Weber, B., Halatsch, J., Aliaga, D.G., Müller, P., 2012. Procedural generation of parcels in urban modeling, in: Computer graphics forum, Wiley Online Library. pp. 681–690.
- Verykokou, S., Ioannidis, C., 2018. Oblique aerial images: a review focusing on georeferencing procedures. International journal of remote sensing 39, 3452–3496.
- Waechter, M., Moehrle, N., Goesele, M., 2014. Let there be color! large-scale texturing of 3d reconstructions, in: European conference on computer vision, Springer. pp. 836–850.
- Wang, T.C., Liu, M.Y., Zhu, J.Y., Tao, A., Kautz, J., Catanzaro, B., 2018. High-resolution image synthesis and semantic manipulation with conditional gans, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 8798–8807.
- Yang, C., Zhang, F., Gao, Y., Mao, Z., Li, L., Huang, X., 2021. Moving car recognition and removal for 3d urban modelling using oblique images. Remote Sensing 13, 3458.
- Yu, J., Lin, Z., Yang, J., Shen, X., Lu, X., Huang, T.S., 2019. Free-form image inpainting with gated convolution, in: Proceedings of the IEEE/CVF international conference on computer vision, pp. 4471–4480.

- Zhang, H., Xu, K., Jiang, W., Lin, J., Cohen-Or, D., Chen, B., 2013. Layered analysis of irregular facades via symmetry maximization. ACM Trans. Graph. 32, 121–1.
- Zhang, H., Yao, Y., Xie, K., Fu, C.W., Zhang, H., Huang, H., 2021. Continuous aerial path planning for 3d urban scene reconstruction. ACM Trans. Graph. 40, 225–1.
- Zhou, G., Bao, X., Ye, S., Wang, H., Yan, H., 2020. Selection of optimal building facade texture images from uav-based multiple oblique image flows. IEEE Transactions on Geoscience and Remote Sensing 59, 1534–1552.
- Zhu, P., Abdal, R., Qin, Y., Wonka, P., 2020. Sean: Image synthesis with semantic regionadaptive normalization, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5104–5113.
- Zhu, P., Para, W.R., Fruehstueck, A., Femiani, J., Wonka, P., 2020. Large scale architectural asset extraction from panoramic imagery. IEEE Transactions on Visualization and Computer Graphics, 1–1.
- Zhu, P., Para, W.R., Frühstück, A., Femiani, J., Wonka, P., 2020a. Large scale architectural asset extraction from panoramic imagery. IEEE Transactions on Visualization and Computer Graphics.
- Zhu, Q., Huang, S., Hu, H., Li, H., Chen, M., Zhong, R., 2021a. Depth-enhanced feature pyramid network for occlusion-aware verification of buildings from oblique images. ISPRS Journal of Photogrammetry and Remote Sensing 174, 105–116.
- Zhu, Q., Shang, Q., Hu, H., Yu, H., Zhong, R., 2021b. Structure-aware completion of photogrammetric meshes in urban road environment. ISPRS Journal of Photogrammetry and Remote Sensing 175, 56–70.
- Zhu, Q., Wang, Z., Hu, H., Xie, L., Ge, X., Zhang, Y., 2020b. Leveraging photogrammetric mesh models for aerial-ground feature point matching toward integrated 3d reconstruction. ISPRS Journal of Photogrammetry and Remote Sensing 166, 26–40.