

# **Paper Review**

## **Old Photo Restoration via Deep Latent Space Translation**

**YeongHyeon Park**

**Department of Electrical and Computer Engineering**

**SungKyunKwan University**

# Old Photo Restoration via Deep Latent Space Translation

Ziyu Wan, Bo Zhang, Dongdong Chen, Pan Zhang, Dong Chen, Fang Wen, Jing Liao

**Abstract**—We propose to restore old photos that suffer from severe degradation through a deep learning approach. Unlike conventional restoration tasks that can be solved through supervised learning, the degradation in real photos is complex and the domain gap between synthetic images and real old photos makes the network fail to generalize. Therefore, we propose a novel triplet domain translation network by leveraging real photos along with massive synthetic image pairs. Specifically, we train two variational autoencoders (VAEs) to respectively transform old photos and clean photos into two latent spaces. And the translation between these two latent spaces is learned with synthetic paired data. This translation generalizes well to real photos because the domain gap is closed in the compact latent space. Besides, to address multiple degradations mixed in one old photo, we design a global branch with a partial nonlocal block targeting the structured defects, such as scratches and dust spots, and a local branch targeting the unstructured defects, such as noises and blurriness. We also extend the global branch with a more memory-efficient scheme, named multi-scale patch-based attention to processing high-resolution photos. Two branches are fused in the latent space, leading to improved capability to restore old photos from multiple defects. Furthermore, we apply another face refinement network to recover fine details of faces in the old photos, thus ultimately generating photos with enhanced perceptual quality. With comprehensive experiments, the proposed pipeline demonstrates superior performance over state-of-the-art methods as well as existing commercial tools in terms of visual quality for old photos restoration. Both code and models could be found [here](#).

**Index Terms**—Image Restoration, Image Generation, Latent Space Translation, Mixed degradation

IEEE TRANSACTIONS ON  
PATTERN ANALYSIS AND  
MACHINE INTELLIGENCE



# Main purpose

- **Unstructured defects:** noise, blurriness, color fading, and low resolution
- **Structured defects:** holes, scratches, and spots



Fig. 1: **Old photo restoration results produced by our method.** Our method can handle the complex degradation mixed with both unstructured and structured defects in real old photos. In particular, we recover high-frequency details for face regions, further improving the perceptual quality for portraits. For each image pair, left is the input while the retouched output is shown on the right.

**Bring old photos back to life!**

# Prior Work

**CVPR VIRTUAL**

## Bringing Old Photos Back to Life

Ziyu Wan<sup>1\*</sup>, Bo Zhang<sup>2</sup>, Dongdong Chen<sup>3</sup>, Pan Zhang<sup>4</sup>, Dong Chen<sup>2</sup>, Jing Liao<sup>1†</sup>, Fang Wen<sup>2</sup>

<sup>1</sup>City University of Hong Kong    <sup>2</sup>Microsoft Research Asia    <sup>3</sup>Microsoft Cloud + AI

<sup>4</sup>University of Science and Technology of China



Figure 1: **Old image restoration results produced by our method.** Our method can handle the complex degradation mixed by both unstructured and structured defects in real old photos.

# About First Author



Ziyu Wan

FOLLOW

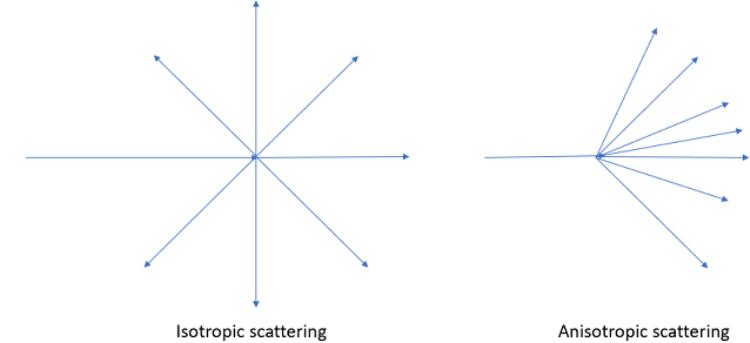
[City University of Hong Kong](#)

Verified email at my.cityu.edu.hk - [Homepage](#)

Computer Vision Computational Photography Computer Graphics

TITLE	CITED BY	YEAR
<a href="#">Bringing Old Photos Back to Life</a> Z Wan, B Zhang, D Chen, P Zhang, D Chen, J Liao, F Wen IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2747-2757	76	2020
<a href="#">Transductive Zero-Shot Learning with Visual Structure Constraint</a> Z Wan, D Chen, Y Li, X Yan, J Zhang, Y Yu, J Liao Advances in Neural Information Processing Systems (NeurIPS)	56	2019
<a href="#">PD-GAN: Probabilistic Diverse GAN for Image Inpainting</a> H Liu, Z Wan, W Huang, Y Song, X Han, J Liao IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)	55	2021
<a href="#">High-Fidelity Pluralistic Image Completion with Transformers</a> Z Wan, J Zhang, D Chen, J Liao IEEE/CVF International Conference on Computer Vision (ICCV)	44	2021
<a href="#">Meta-PU: An arbitrary-scale upsampling network for point cloud</a> S Ye, D Chen, S Han, Z Wan, J Liao IEEE Transactions on Visualization and Computer Graphics (TVCG)	19	2021
<a href="#">Old Photo Restoration via Deep Latent Space Translation</a> Z Wan, B Zhang, D Chen, P Zhang, D Chen, J Liao, F Wen IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)	16	2022
<a href="#">DeFLOCNet: Deep Image Editing via Flexible Low-level Controls</a> H Liu, Z Wan, W Huang, Y Song, X Han, J Liao, B Jiang, W Liu IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)	8	2021
<a href="#">Transductive Zero-Shot Learning via Visual Center Adaptation</a> Z Wan, Y Li, M Yang, J Zhang Proceedings of the AAAI Conference on Artificial Intelligence 33 (01), 10059 ...	3	2019
<a href="#">Visual Structure Constraint for Transductive Zero-Shot Learning in the Wild</a> Z Wan, D Chen, J Liao International Journal of Computer Vision (IJCV) 129 (6), 1893-1909	2	2021
<a href="#">Bringing Old Films Back to Life</a> Z Wan, B Zhang, D Chen, J Liao IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)	1	2022
<a href="#">FDNeRF: Few-shot Dynamic Neural Radiance Fields for Face Reconstruction and Expression Editing</a> J Zhang, X Li, Z Wan, C Wang, J Liao		2022
<a href="#">Adaptive Joint Optimization for 3D Reconstruction with Differentiable Rendering</a> J Zhang, Z Wan, J Liao IEEE Transactions on Visualization and Computer Graphics (TVCG)		2022
<a href="#">Supplementary Material-Bringing Old Films Back to Life</a> Z Wan, B Zhang, D Chen, J Liao		

# Summary



## Purpose

- Unpaired old (degraded) photo restoration

## Contributions

- Domain alignment: well generalization on real photos by training synthetic data
- Partial nonlocal block: reconstruction of blind spots / corruption area
- Multi-scale patch-based fusion: reduction of memory cost and better performance
- Coarse-to-fine generator: reconstruction of the high-resolution face

## Limitations

- Anisotropic pollution is not totally restored.
- Shading along the folds of the photo hindrance the restoration process.

# Problem Definition

# Problems

- **Unstructured defects:** noise, blurriness, color fading, and low resolution
- **Structured defects:** holes, scratches, and spots

Unstructured distortion



Structured distortion



## Single / Mixed degradation image restoration

- Single degradation = unstructured or structured defects
- Mixed degradation = unstructured and structured defects
- Classical methods have **over smoothness** issue.
- Deep learning methods have a **limitation of generalization**.
  - The deep neural network learns synthetic data due to unpaired dataset.
  - So, it has a limitation to restore the real photo.

Noise, blurring, color fading, sepia issue, ...

Scratches, dust, pollution, ....

## Face restoration

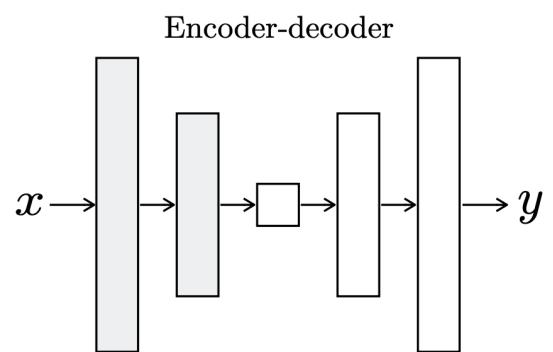
- Wild-portraits have **various posed faces**.
- Face identity should be preserved while the face restoration task.

## Old photo restoration

- Prior studies focus on inpainting task to remove structured defects.
- Also, those are not consider to restore **unstructured defects**.

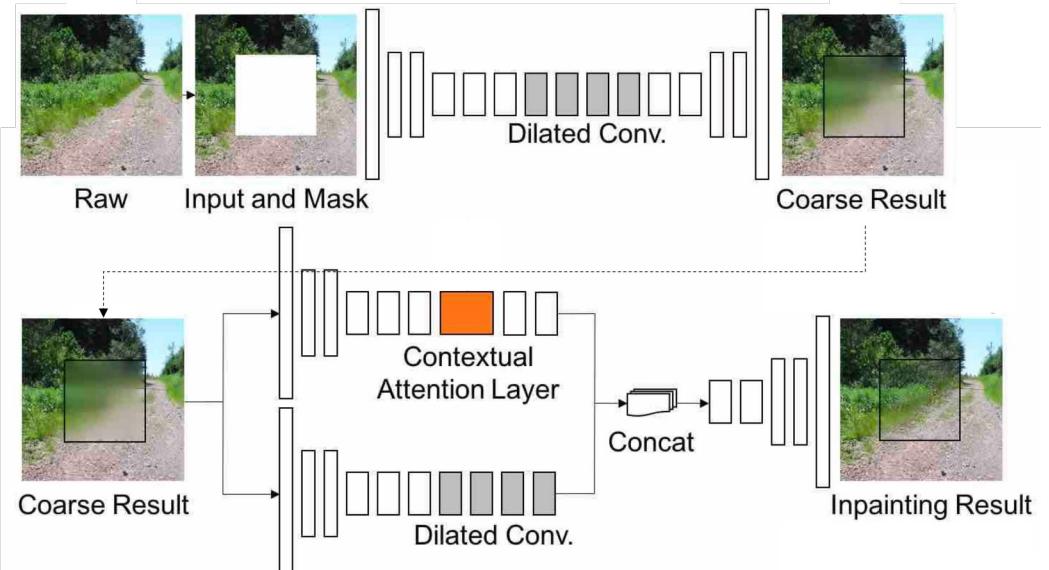
# Single degradation image restoration

**Pix2Pix**



Pix2Pix cannot fill hole.

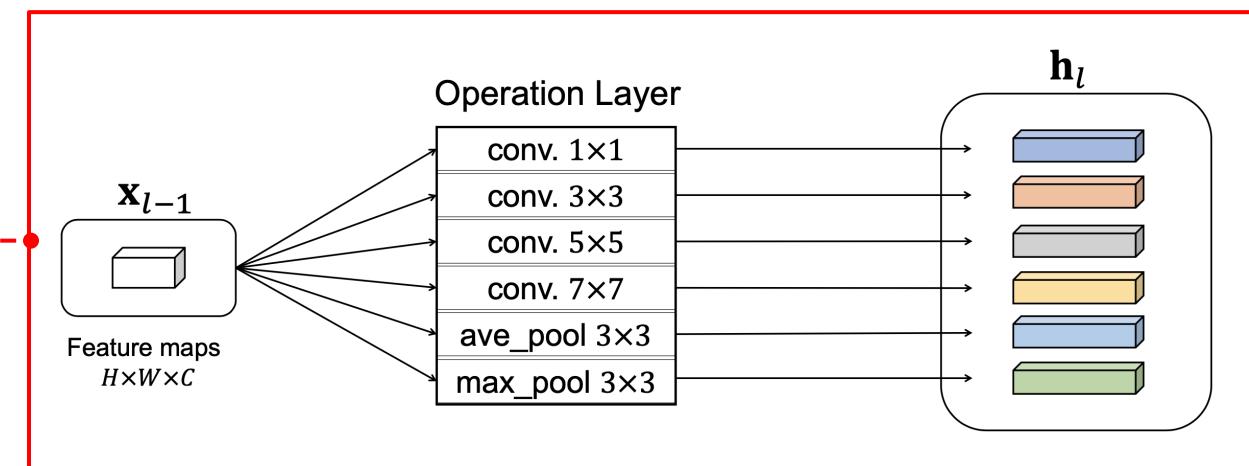
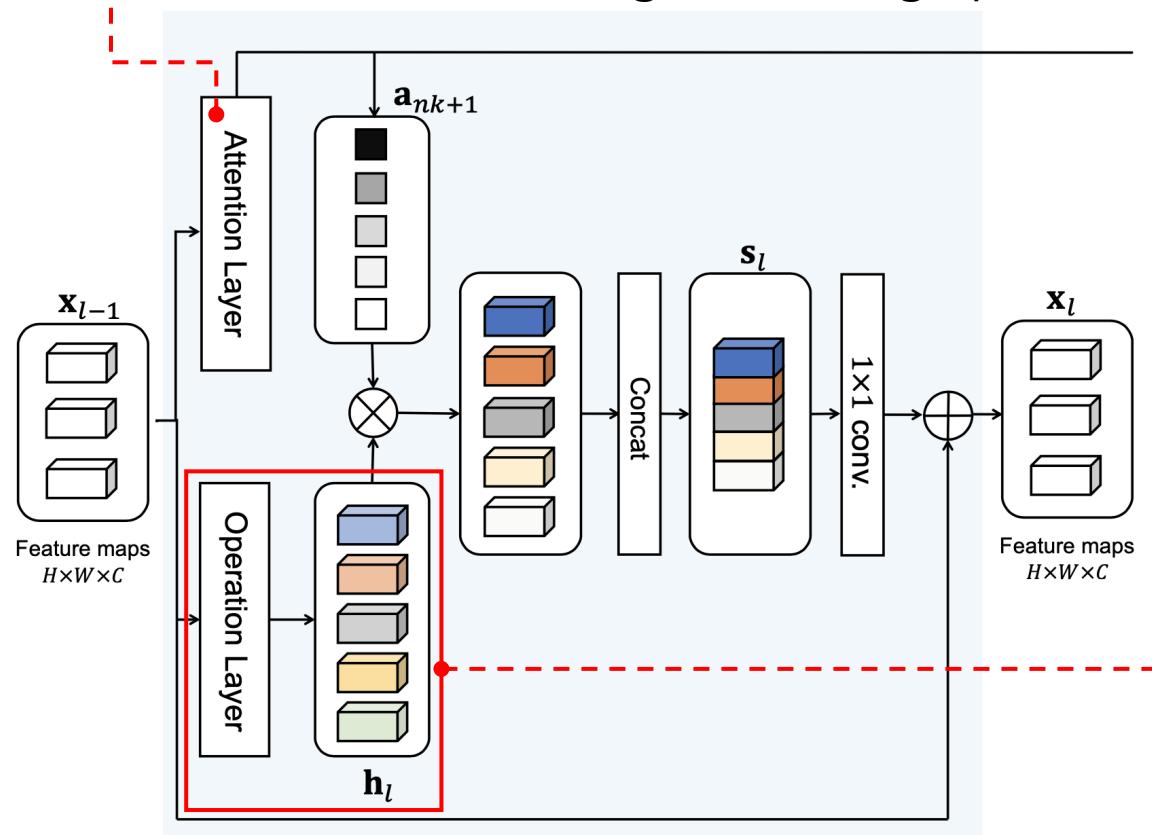
**Contextual Attention**



# Mixed degradation image restoration [1/2]

Operation-wise attention

- 1 x 1 convolution & global average pool



- $k \times k$  convolution ( $k=1, 3, 5, 7$ )
- $k \times k$  dilated convolution
- $3 \times 3$  average pool
- $3 \times 3$  max pool

## Mixed degradation image restoration [2/2]

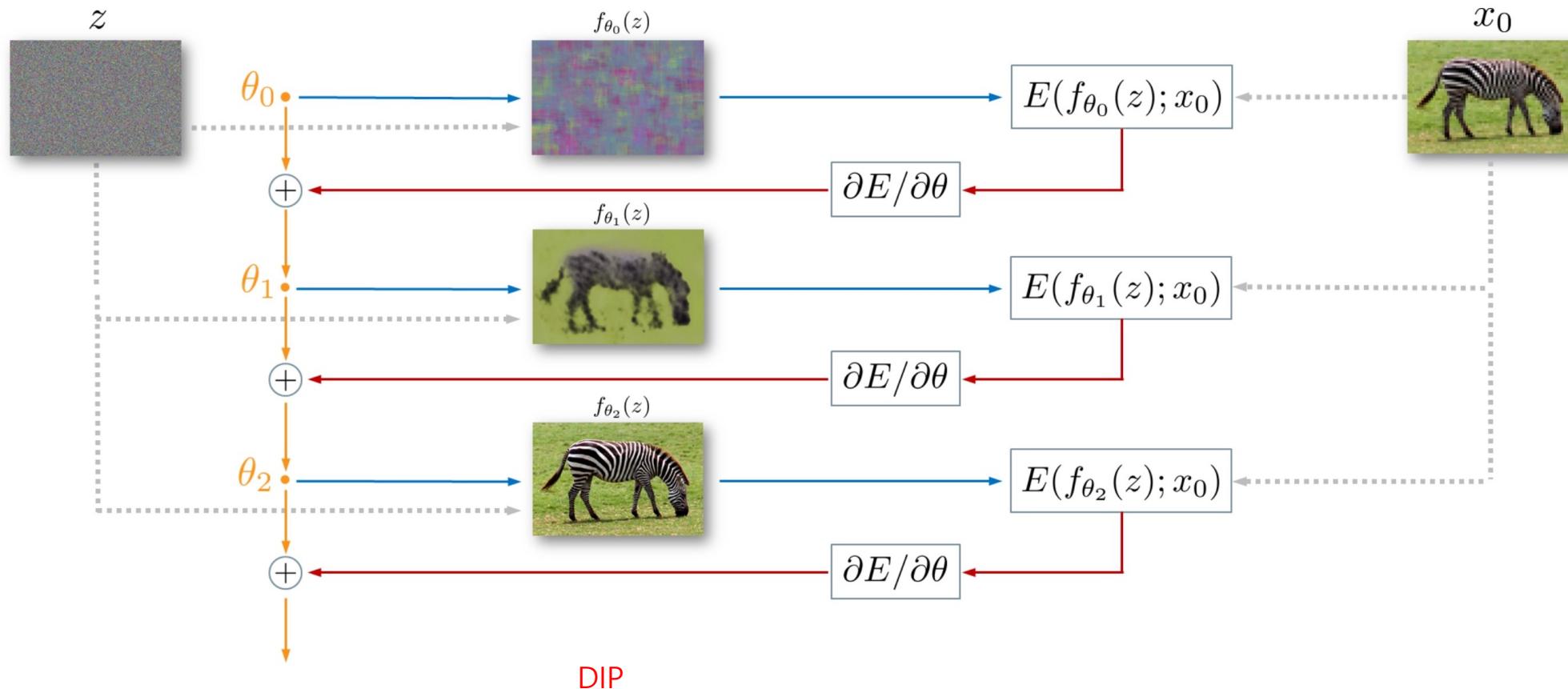
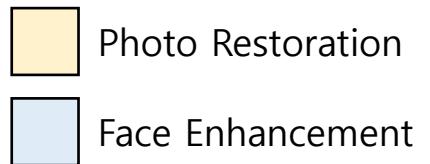
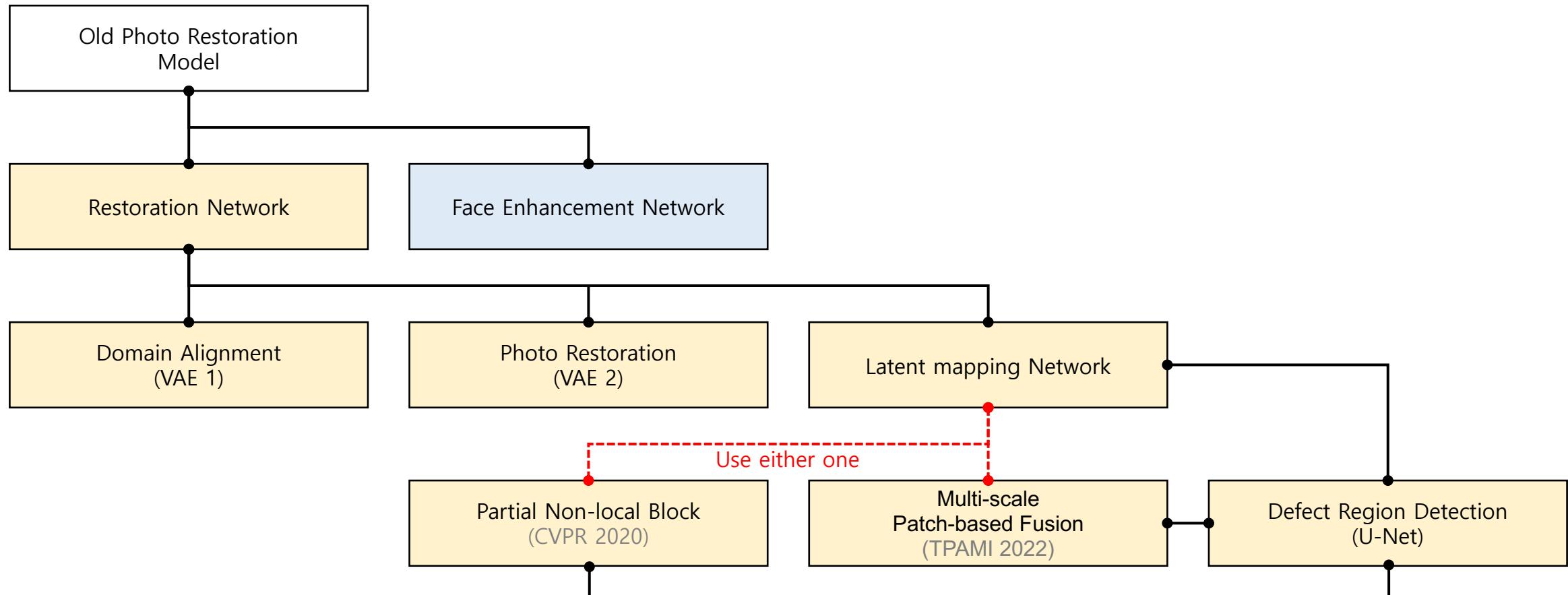


Fig. 2: **Image restoration using the deep image prior.** Starting from a random weights  $\theta_0$ , we iteratively update them in order to minimize the data term eq. (2). At every iteration the weights  $\theta$  are mapped to an image  $x = f_{\theta}(z)$ , where  $z$  is a fixed tensor and the mapping  $f$  is a neural network with parameters  $\theta$ . The image  $x$  is used to compute the task-dependent loss  $E(x, x_0)$ . The gradient of the loss w.r.t. the weights  $\theta$  is then computed and used to update the parameters.

# Old Photo Restoration Model



# Overview



# Photo Restoration [1/11]

Latent Space Translation (Intersection Maximization)

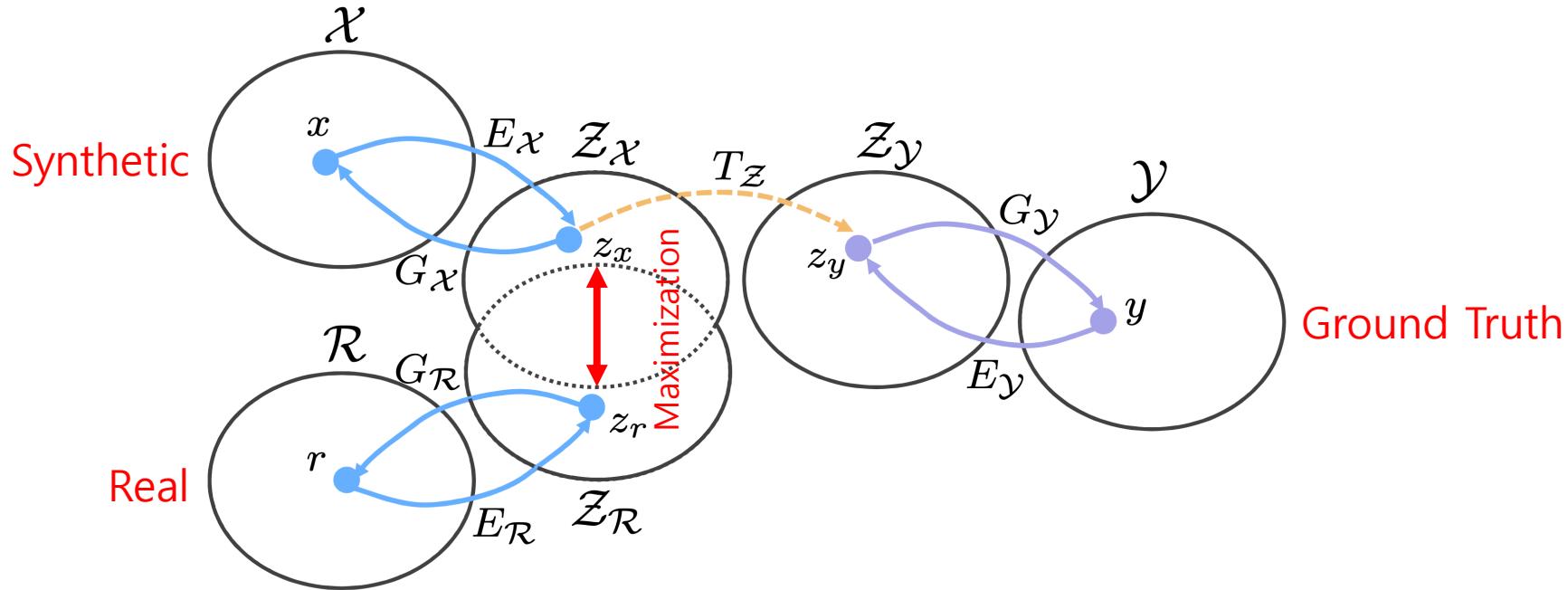


Fig. 2: Illustration of our translation method with three domains. The domain gap between  $\mathcal{Z}_x$  and  $\mathcal{Z}_{\mathcal{R}}$  will be reduced in the shared latent space.

# Photo Restoration [2/11]

## Out-of-Distribution (Intersection Minimization)

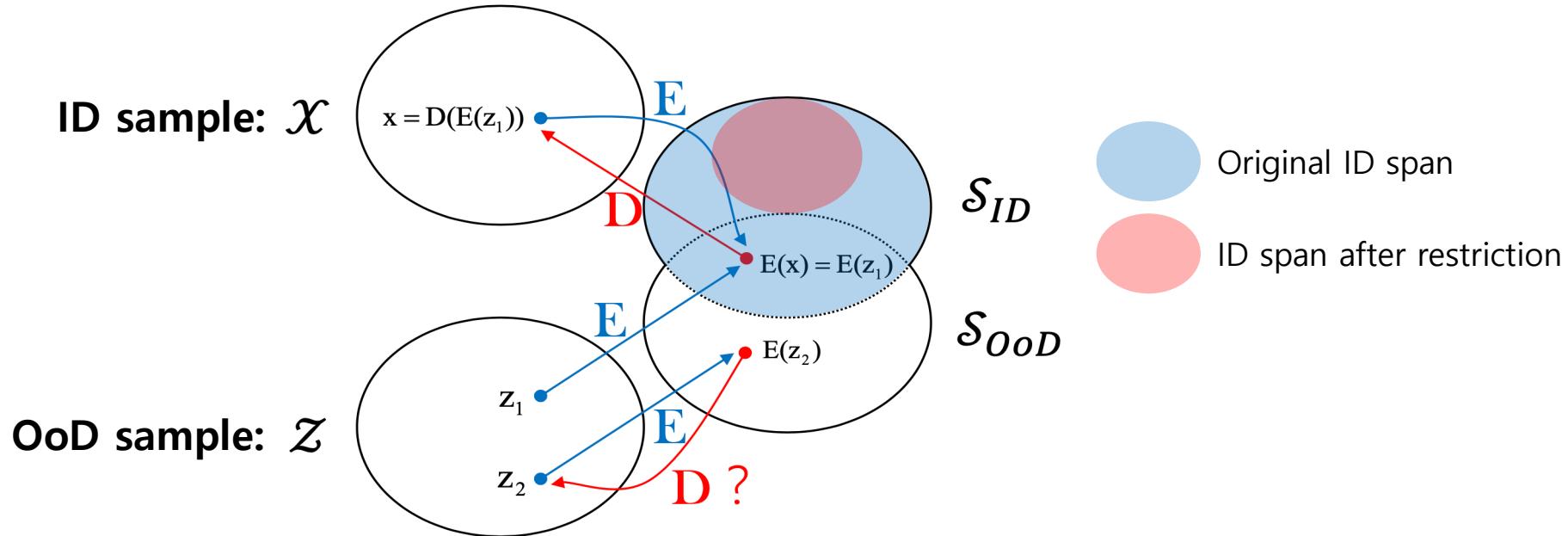
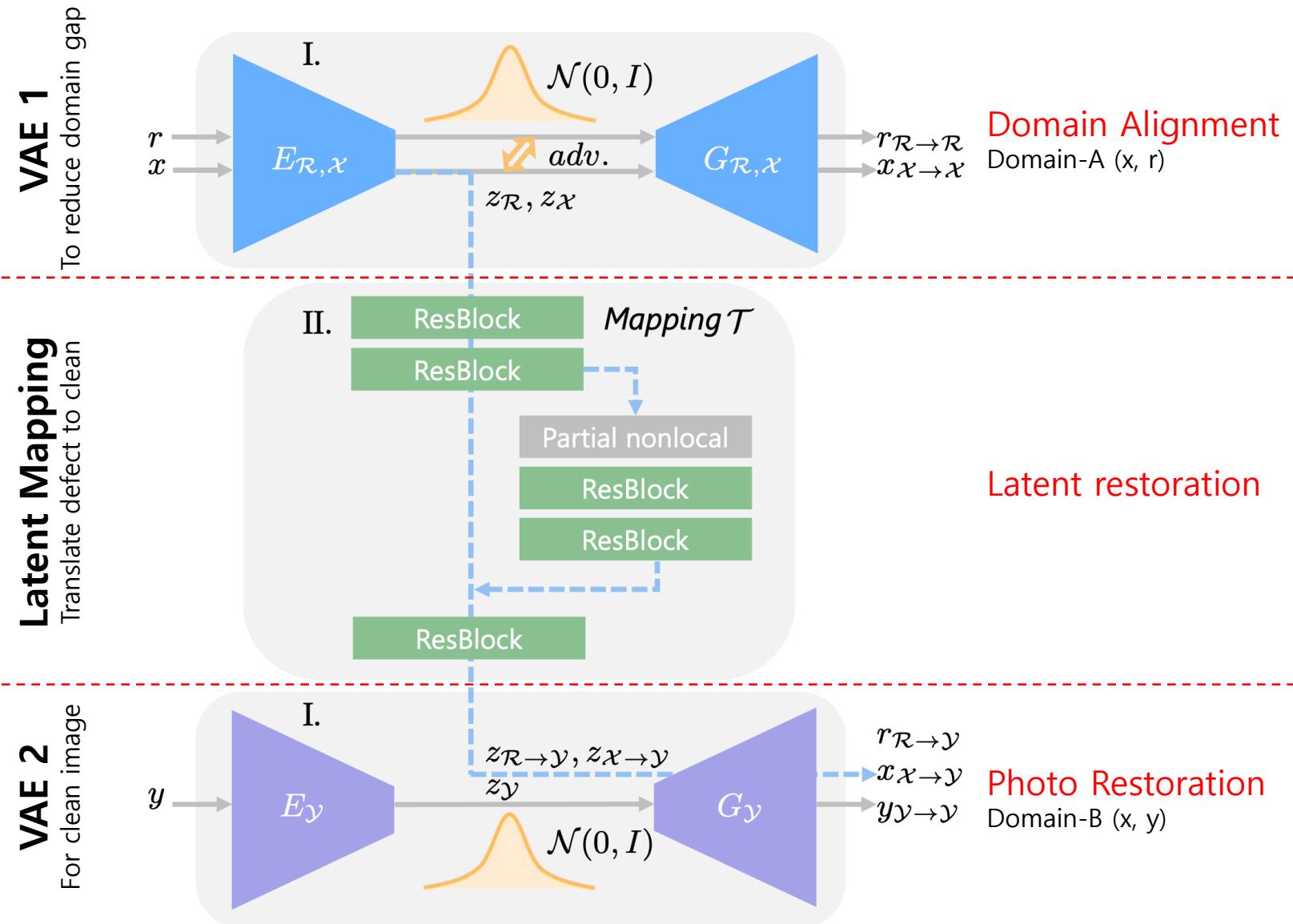


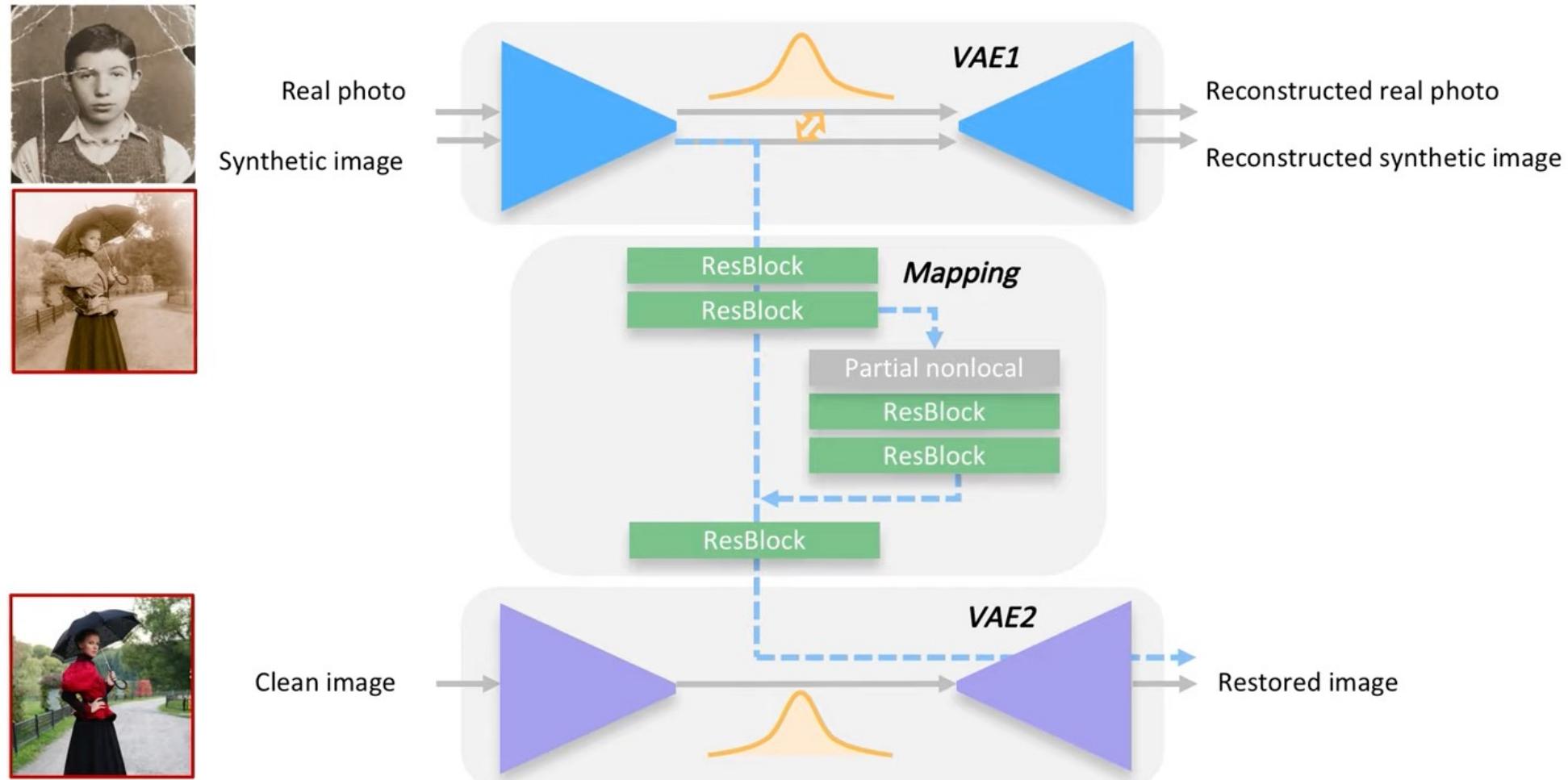
Figure 1. Illustration of the described quadruplet domain translation. For an OoD sample  $z_1$  encoded into  $\mathcal{S}_{ID} \cap \mathcal{S}_{OoD}$ , its latent representation  $E(z_1)$  is equal potentially to that of an ID sample  $x$ . Therefore,  $E(z_1)$  can be decoded to a different sample  $x$  within  $\mathcal{X}$ , resulting in a large reconstruction error. However, for an OoD sample  $z_2$  with latent representation  $E(z_2)$  lying outside  $\mathcal{S}_{ID}$ , it offers no guarantee that it could not be reconstructed well.

# Photo Restoration [3/11]

## Networks

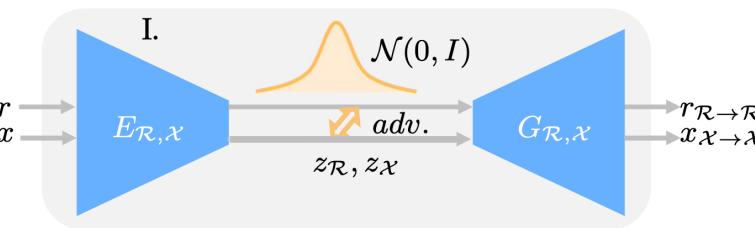


# Photo Restoration [4/11]

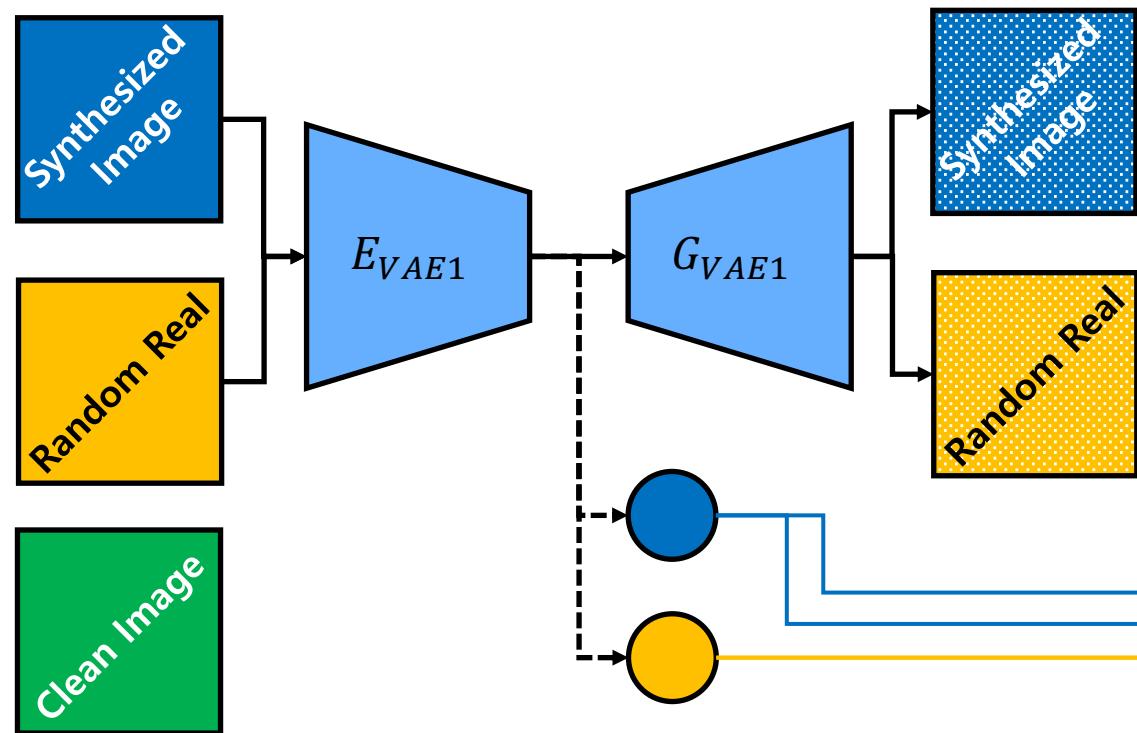


# Photo Restoration [5/11]

Training – VAE 1



$$\begin{aligned} \min_D V_{\text{LSGAN}}(D) &= \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [(D(\mathbf{x}) - 1)^2] \\ &\quad + \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [(D(G(\mathbf{z})) + 1)^2] \quad (8) \\ \min_G V_{\text{LSGAN}}(G) &= \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [(D(G(\mathbf{z})))^2]. \end{aligned}$$



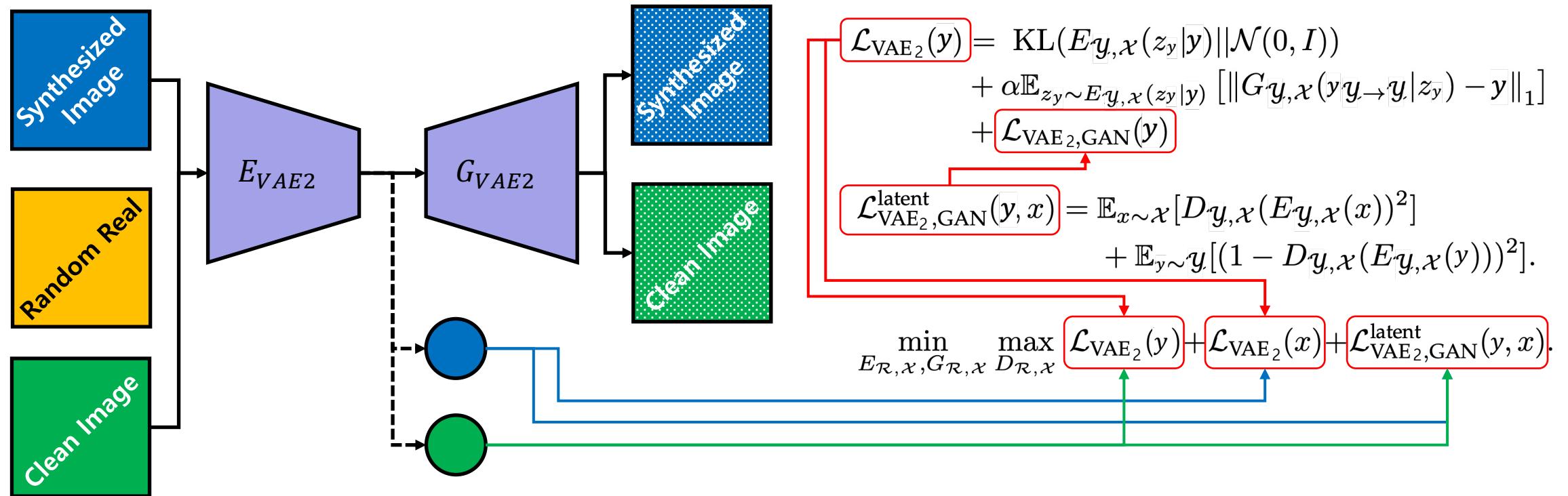
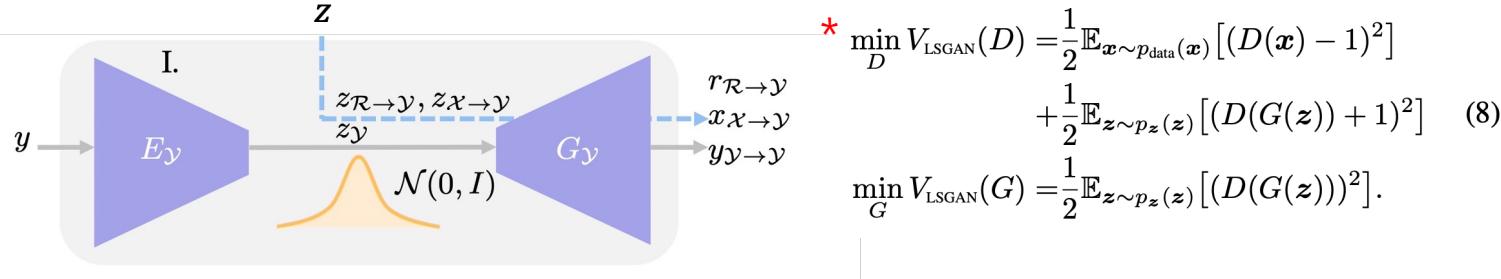
$$\begin{aligned} \mathcal{L}_{\text{VAE}_1}(r) &= \text{KL}(E_{\mathcal{R}, \mathcal{X}}(z_r | r) || \mathcal{N}(0, I)) \\ &\quad + \alpha \mathbb{E}_{z_r \sim E_{\mathcal{R}, \mathcal{X}}(z_r | r)} [\|G_{\mathcal{R}, \mathcal{X}}(r_{\mathcal{R} \rightarrow \mathcal{R}} | z_r) - r\|_1] \quad (2) \\ &\quad + \mathcal{L}_{\text{VAE}_1, \text{GAN}}(r) \end{aligned}$$

$$\begin{aligned} \mathcal{L}_{\text{VAE}_1, \text{GAN}}^{\text{latent}}(r, x) &= \mathbb{E}_{x \sim \mathcal{X}} [D_{\mathcal{R}, \mathcal{X}}(E_{\mathcal{R}, \mathcal{X}}(x))^2] \\ &\quad + \mathbb{E}_{r \sim \mathcal{R}} [(1 - D_{\mathcal{R}, \mathcal{X}}(E_{\mathcal{R}, \mathcal{X}}(r)))^2]. \quad (3) \end{aligned}$$

$$\min_{E_{\mathcal{R}, \mathcal{X}}, G_{\mathcal{R}, \mathcal{X}}} \max_{D_{\mathcal{R}, \mathcal{X}}} \mathcal{L}_{\text{VAE}_1}(r) + \mathcal{L}_{\text{VAE}_1}(x) + \mathcal{L}_{\text{VAE}_1, \text{GAN}}^{\text{latent}}(r, x). \quad (4)$$

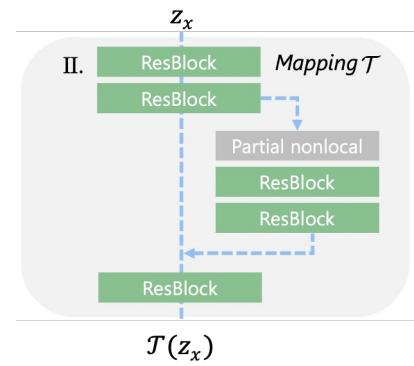
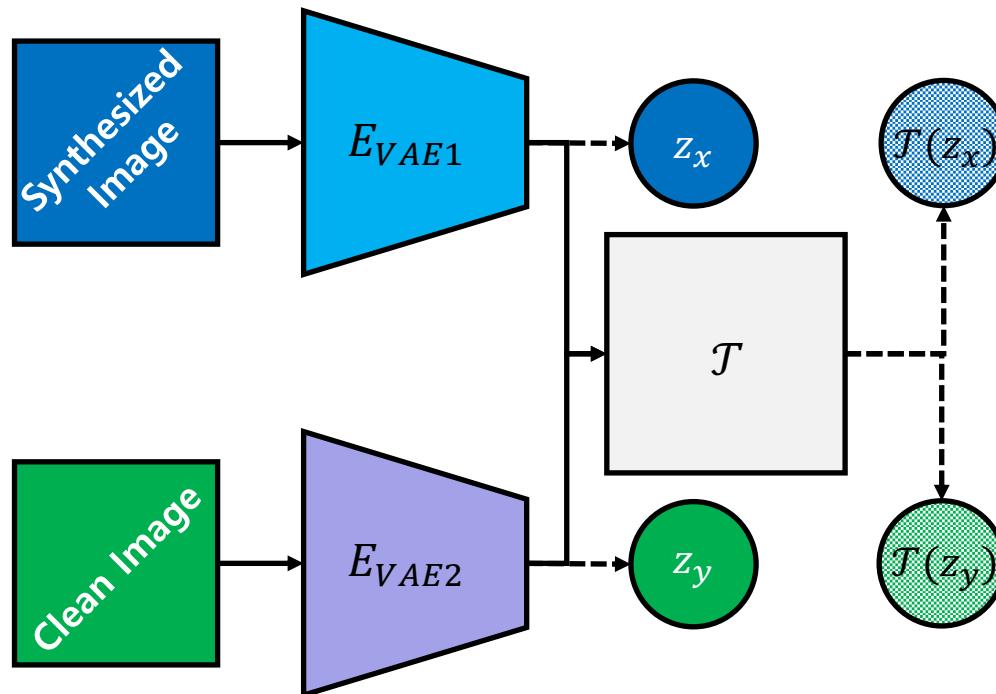
# Photo Restoration [6/11]

Training – VAE 2



# Photo Restoration [7/11]

## Training – Mapping Network



$$\mathcal{L}_{\mathcal{T}}(x, y) = \lambda_1 \mathcal{L}_{\mathcal{T}, \ell_1} + \mathcal{L}_{\mathcal{T}, GAN} + \lambda_2 \mathcal{L}_{FM} \quad (5)$$

$$\mathcal{L}_{\mathcal{T}, \ell_1} = \mathbb{E} \|\mathcal{T}(z_x) - z_y\|_1$$

$$\mathcal{L}_{VAE_1, GAN}^{\text{latent}}(r, x) = \mathbb{E}_{x \sim \mathcal{X}} [D_{\mathcal{T}}(\mathcal{T}(z_x))^2] + \mathbb{E}_{r \sim \mathcal{R}} [(1 - D(r))^2] \quad \text{LPIPS with L1}$$

$$\mathcal{L}_{FM} = \mathbb{E} \left[ \sum_i \frac{1}{n_{D_{\mathcal{T}}}^i} \|\phi_{D_{\mathcal{T}}}^i(x_{\mathcal{X} \rightarrow \mathcal{Y}}) - \phi_{D_{\mathcal{T}}}^i(y_{\mathcal{Y} \rightarrow \mathcal{Y}})\|_1 + \sum_i \frac{1}{n_{VGG}^i} \|\phi_{VGG}^i(x_{\mathcal{X} \rightarrow \mathcal{Y}}) - \phi_{VGG}^i(y_{\mathcal{Y} \rightarrow \mathcal{Y}})\|_1 \right], \quad (6)$$

$\phi_{D_{\mathcal{T}}}^i$  ( $\phi_{VGG}^i$ ) : feature map of the i-th layer

$n_{D_{\mathcal{T}}}^i$  ( $n_{VGG}^i$ ) : number of activation in i-th layer

# Photo Restoration [8/11]

## Non-local Block (Inpainting with global and local features)

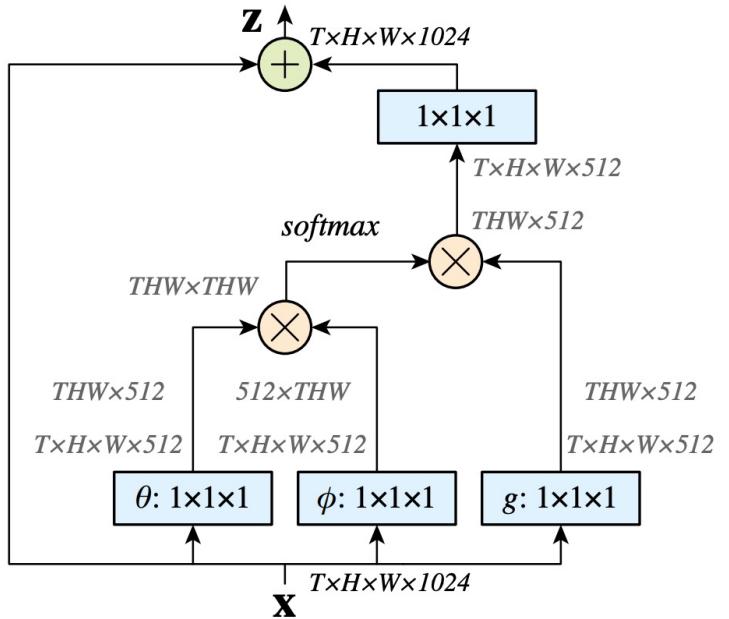


Figure 2. A spacetime **non-local block**. The feature maps are shown as the shape of their tensors, e.g.,  $T \times H \times W \times 1024$  for 1024 channels (proper reshaping is performed when noted). “ $\otimes$ ” denotes matrix multiplication, and “ $\oplus$ ” denotes element-wise sum. The softmax operation is performed on each row. The blue boxes denote  $1 \times 1 \times 1$  convolutions. Here we show the embedded Gaussian version, with a bottleneck of 512 channels. The vanilla Gaussian version can be done by removing  $\theta$  and  $\phi$ , and the dot-product version can be done by replacing softmax with scaling by  $1/N$ .

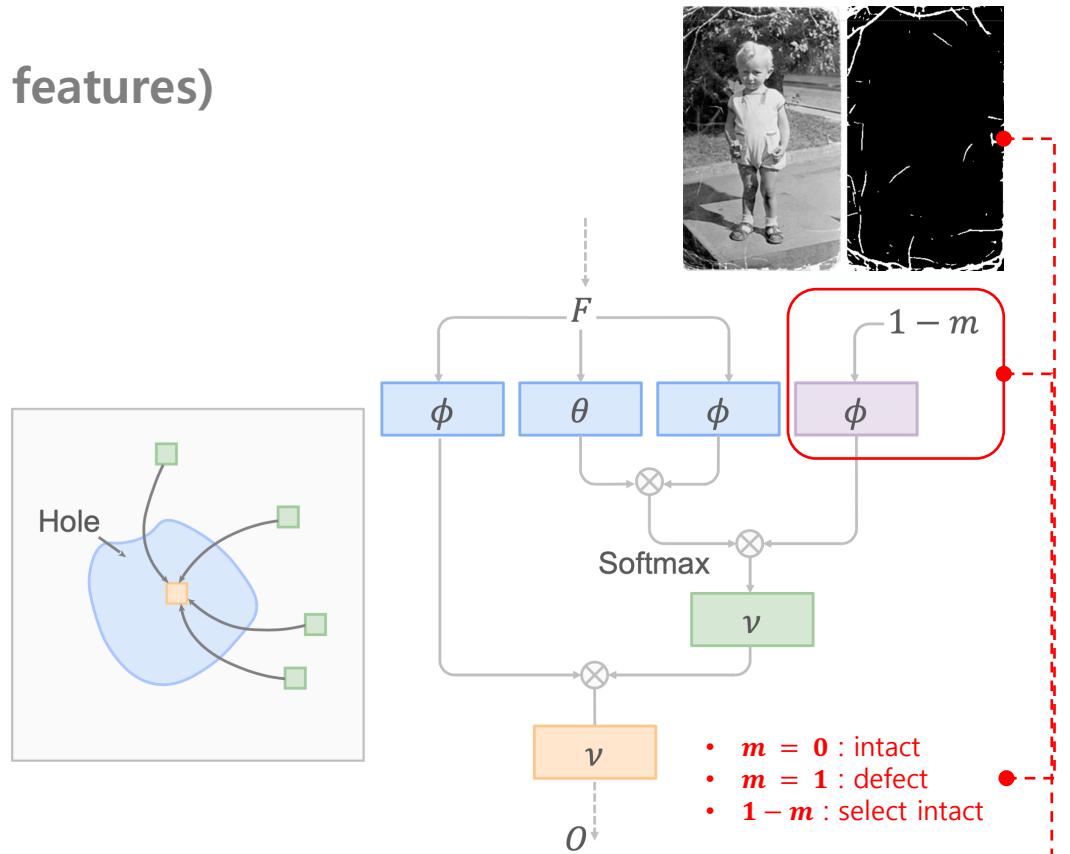
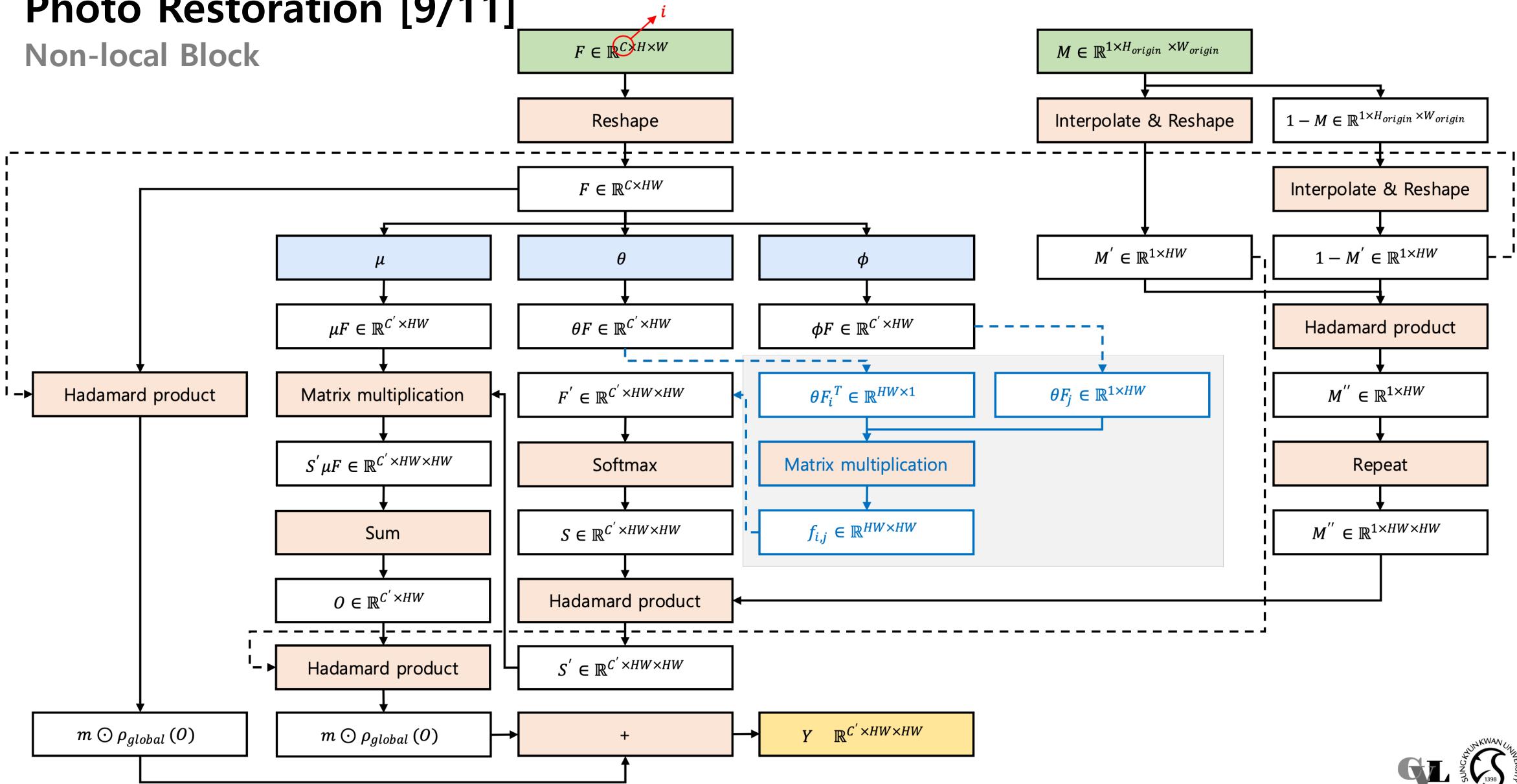


Fig. 4: **Partial nonlocal block**. Left shows the principle. The pixels within the hole areas are inpainted by the context pixels outside the corrupted region. Right shows the detailed implementation.

Suppress defective features for inpainting

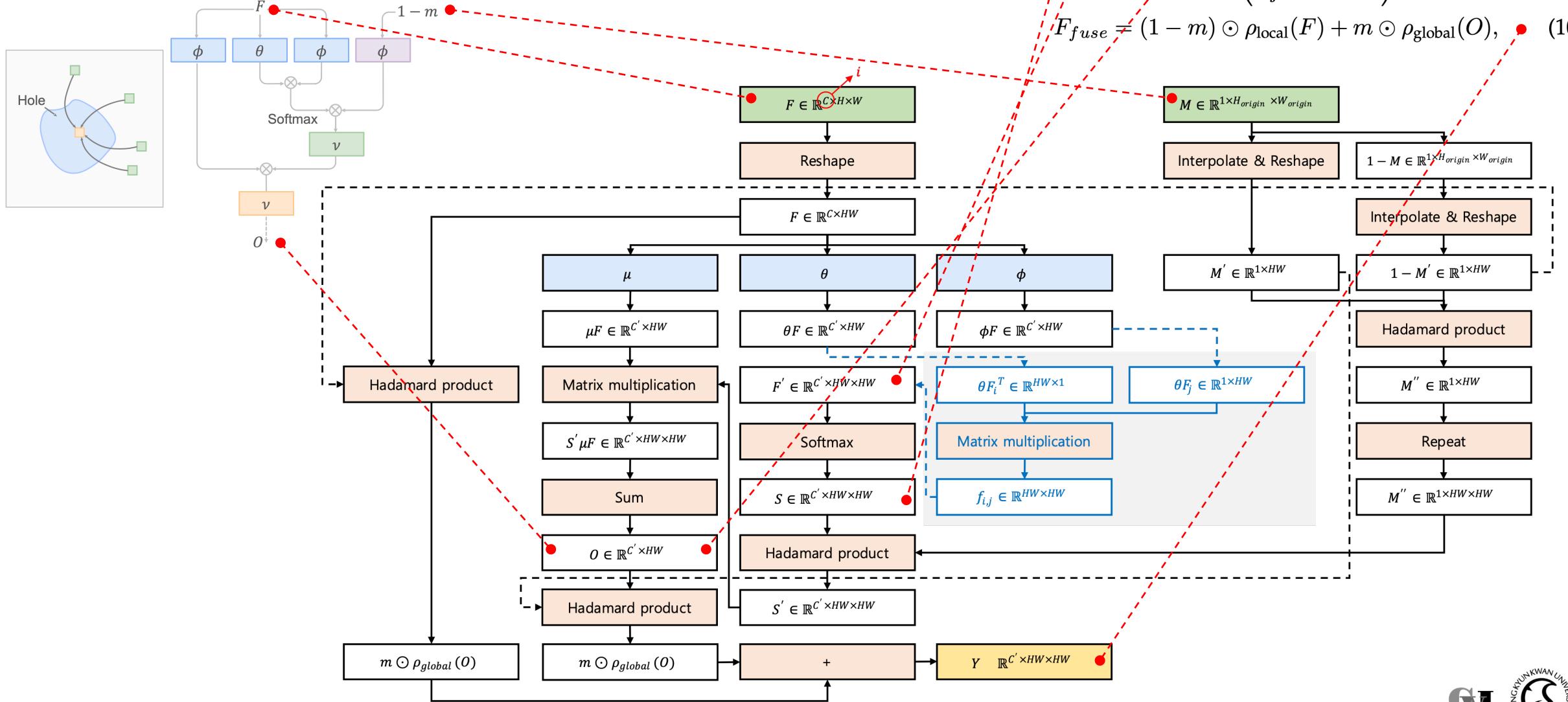
## Photo Restoration [9/11]

## Non-local Block



# Photo Restoration [10/11]

## Non-local Block



$$s_{i,j} = (1 - m_j) f_{i,j} / \sum_{k \neq i} (1 - m_k) f_{i,k}, \quad (7)$$

$$f_{i,j} = \exp(\theta(F_i)^T \cdot \phi(F_j)) \quad (8)$$

$$O_i = \nu \left( \sum_{\forall j} s_{i,j} \mu(F_j) \right), \quad (9)$$

$$F_{fuse} = (1 - m) \odot \rho_{local}(F) + m \odot \rho_{global}(O), \quad (10)$$

# Photo Restoration [11/11]

## High-Resolution Processing

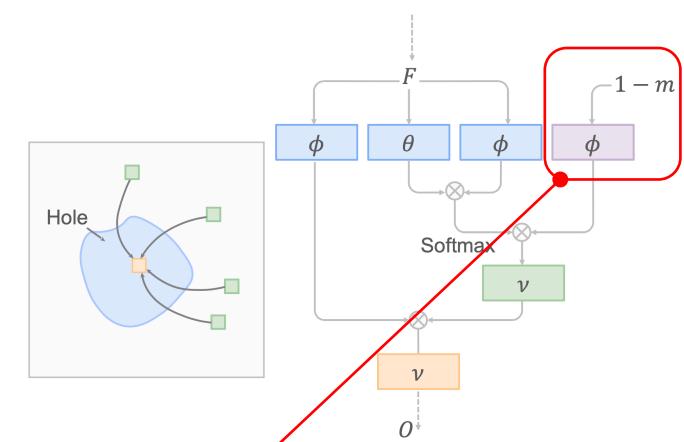
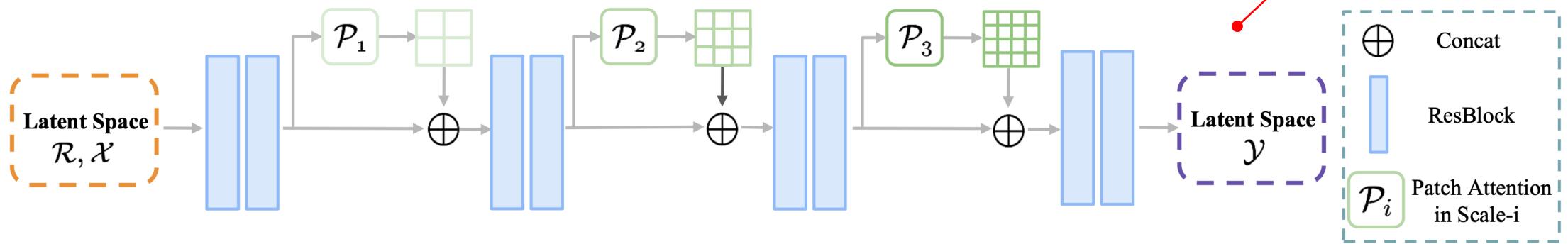


Fig. 5: **Multi-scale patch-based fusion.** The new proposed latent space mapping network vastly reduces the memory cost of partial nonlocal block and enable the processing of high-resolution photos.

# Face Enhancement

To restore the blind spot

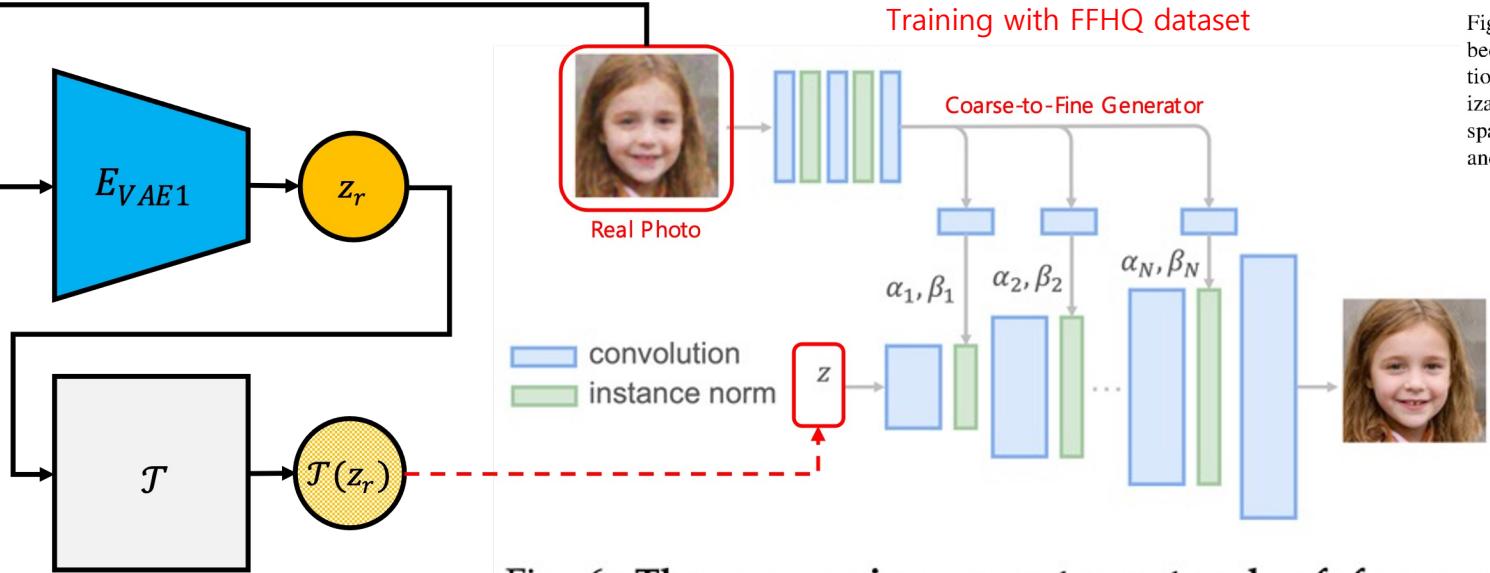


Fig. 6: The progressive generator network of face enhancement. Starting from a latent vector  $z$ , the network up-samples the feature map by deconvolution progressively. The degraded face will be injected into different resolutions in a spatial condition manner.

$$\mathcal{L}_{\text{perc}}^{\text{face}} = \mathbb{E} \left[ \sum_i \frac{1}{n_{\text{VGG}}^i} \|\phi_{\text{VGG}}^i(G_f(z, r_f)) - \phi_{\text{VGG}}^i(r_c)\|_1 \right], \quad (13)$$

$r_f$ : degraded face  
 $r_c$ : ground truth

$$\begin{aligned} \mathcal{L}_{\text{GAN}}^{\text{face}}(z, r_f, r_c) = & \mathbb{E}_{z \sim \mathcal{Z}, r_f \sim \mathcal{R}_f} [D_f(G_f(z, r_f))^2] \\ & + \mathbb{E}_{r_c \sim \mathcal{R}_c} [(1 - D_f(r_c))^2]. \end{aligned} \quad (14)$$

$$\gamma_{x,y,c}(r_f^i) \frac{h_{x,y,c} - \mu_c}{\sqrt{\sigma_c^2 + \epsilon}} + \beta_{x,y,c}(r_f^i), \quad (12)$$

Learnable parameter

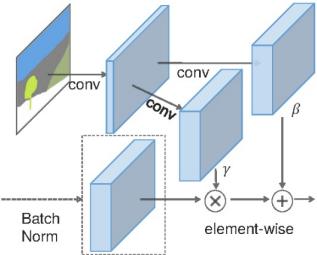


Figure 2: In SPADE, the mask is first projected onto an embedding space, and then convolved to produce the modulation parameters  $\gamma$  and  $\beta$ . Unlike prior conditional normalization methods,  $\gamma$  and  $\beta$  are not vectors, but tensors with spatial dimensions. The produced  $\gamma$  and  $\beta$  are multiplied and added to the normalized activation element-wise.

# Experiments

# Data Preparation



Pascal VOC



Random Crop



Unstructured Degradation

- Gaussian white noise
- Gaussian blur
- JPEG compression
- Downsampling and upsampling

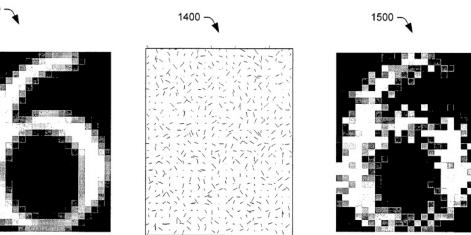


FIG. 13

FIG. 14

FIG. 15



Structured Degradation  
(using Elastic Distortions)

- 62 scratchtexture images
- 55 paper texture images
- Random holes
- Film grain noise

+ 783 manual annotated samples

# Defect Region Detection

For mask generation



Fig. 9: Defect region detection results on real photos.

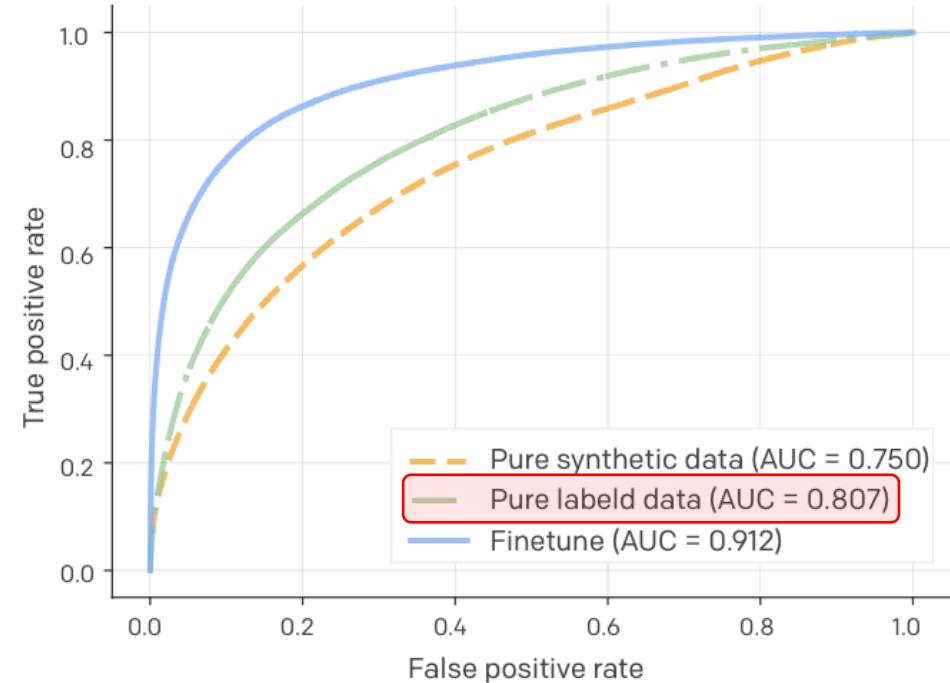


Fig. 8: ROC curve for scratch detection of different data settings. Combining both synthetic structured degradations and a small amount of labeled data, the scratch detection network could achieve great results.

Pure synthetic: synthetic data from Pascal VOC

Pure labeled: manually annotated real degraded photo

Finetune: finetune with manually annotated data on pretrained model (w/ synthetic data)

# Quantitative Comparison

Method	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
Input	12.92	0.49	0.59	306.80
DIP [53]	22.59	0.57	0.54	194.55
Sequential [27], [70]	22.71	0.60	0.49	191.98
Attention [52]	<b>24.12</b>	<b>0.70</b>	0.33	208.11
Pix2pix [71]	22.18	0.62	<b>0.23</b>	135.14
Ours w/o partial nonlocal	23.14	0.68	0.26	143.62
Ours* [72] (CVPR 2020)	23.33	0.69	0.25	<b>134.35</b>
<b>Ours</b>	<b>23.45</b>	<b>0.71</b>	<b>0.24</b>	<b>130.64</b>

TABLE 2: Quantitative results on the DIV2K dataset.  $\uparrow$  indicates that a higher score denotes a good image quality. We highlight the best two scores for each measure.

- LPIPS: feature distance
- FID: feature distribution distance

$$FID = \|\mu_x - \mu_y\|^2 - \|\sigma_x - \sigma_y\|^2$$

Method	Top 1	Top 2	Top 3	Top 4	Top 5
DIP [53]	2.54	8.49	19.26	39.09	74.22
CycleGAN [48]	4.24	8.21	19.54	28.32	50.42
Sequential [27], [70]	4.81	18.13	47.87	79.60	94.61
Attention [52]	6.79	21.24	49.85	73.08	88.38
Pix2Pix [71]	16.14	60.90	73.65	86.68	94.90
<b>Ours</b>	<b>65.43</b>	<b>83.00</b>	<b>89.80</b>	<b>93.20</b>	<b>97.45</b>

TABLE 3: User study results. The percentage (%) of each method is selected as the top  $K$  ( $K = 1 - 5$ ) by users.

# Qualitative Comparison [1/2]

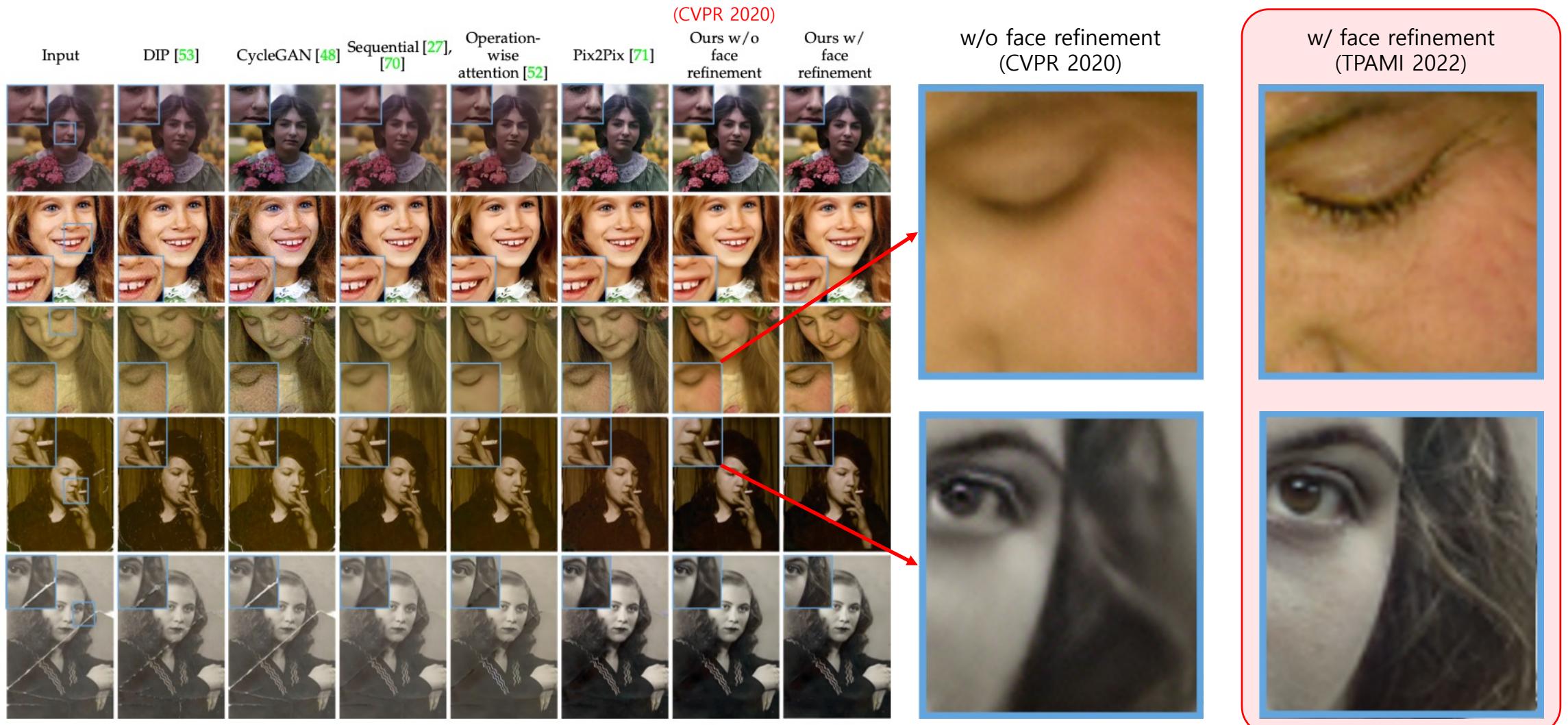


Fig. 10: **Qualitative comparison against state-of-the-art methods.** It shows that our method can restore both unstructured and structured degradation and our recovered results are significantly better than other methods.

## Qualitative Comparison [2/2]

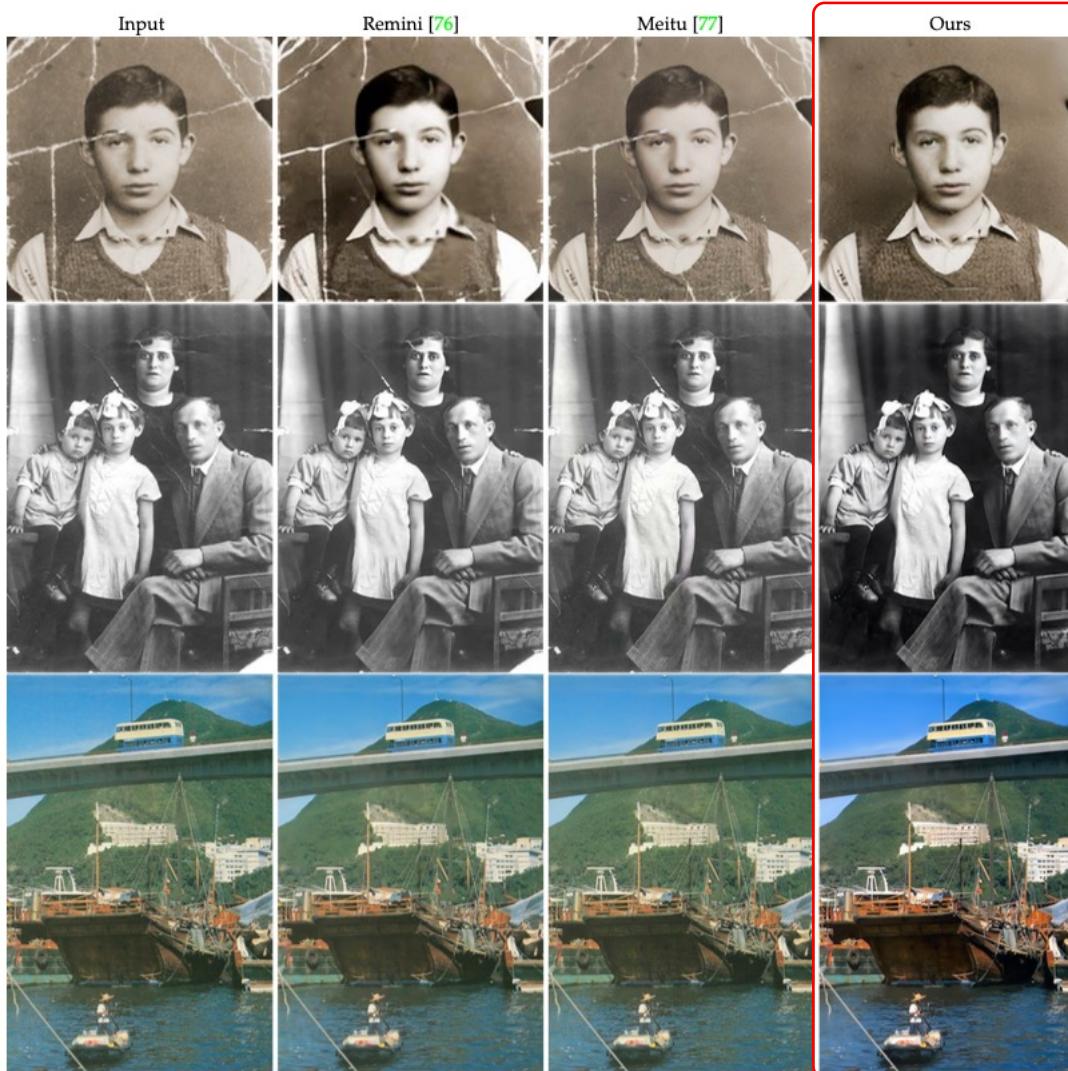


Fig. 11: Qualitative comparisons against commercial tools. Remini Photo Enhancer [76], Meitu [77] and our full pipeline results are included. The last two rows are high-resolution restoration (736x1024 and 640x960).

# Ablation Study

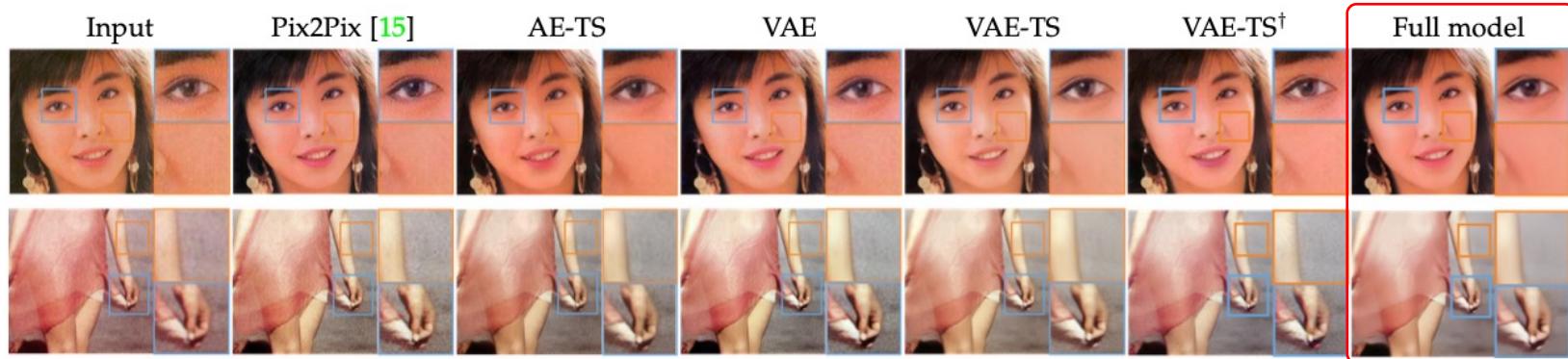


Fig. 12: **Ablation study for latent translation with VAEs.** By involving feature translation, VAE and feature-level adversarial loss, the domain gap between synthetic degradations and real-world defects could be narrowed better, leading to better restoration results step by step. VAE-TS<sup>†</sup>: Remove the real old data for two-stage VAE.

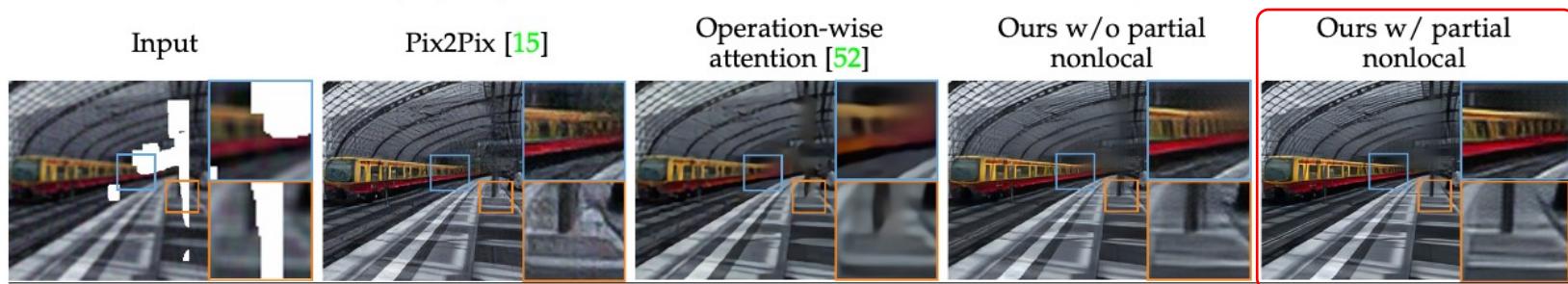


Fig. 13: **Ablation study of the partial nonlocal block.** Partial nonlocal better inpaints the structured defects.



Fig. 14: **Ablation study of the partial nonlocal block.** Partial nonlocal does not touch the non-hole regions as this operation is aware of the corruption area.

Annotation	Neural Network	Training
AE-TS	Two Auto-Encoder	Two-stage
VAE	Two Variational AE	Simultaneously
VAE-TS	Two VAE	Two-stage
VAE-TS <sup>†</sup>	Two VAE	Two-stage
Full model	Two VAE + Face refine	Two-stage

Method	Pix2Pix	AEs-TS	VAEs	VAEs-TS	VAE-TS <sup>†</sup>	full model
Wasserstein ↓	1.837	1.432	1.048	0.765	1.027	<b>0.581</b>
BRISQUE ↓	25.549	24.547	23.949	23.396	23.850	<b>23.016</b>

TABLE 4: **Ablation study of latent translation with VAEs.** We provide some quantitative comparisons here to demonstrate the superior performance of the full model. Our full method achieves the best results on both distribution distance and BRISQUE metric.

# Effect of Multi-Scale Patch-based Fusion

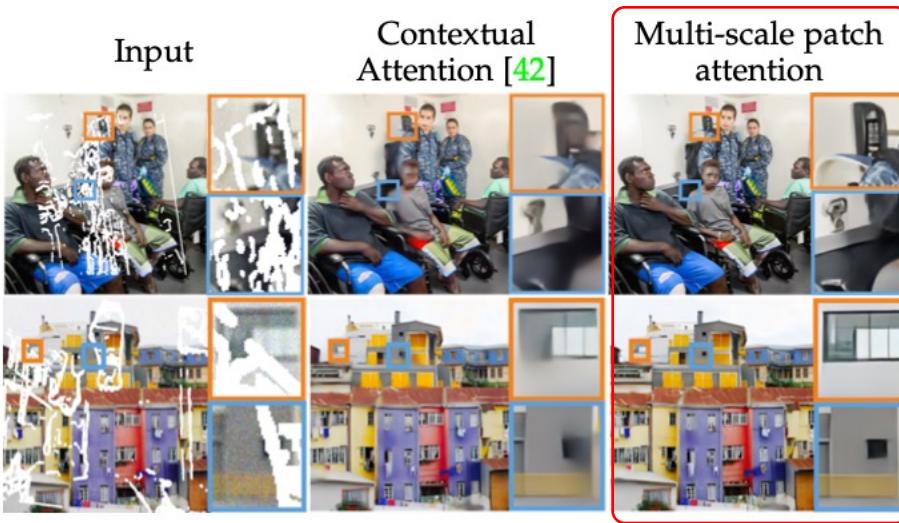


Fig. 15: Qualitative comparisons between contextual attention and ours on high-resolution (1024x1024) photos.

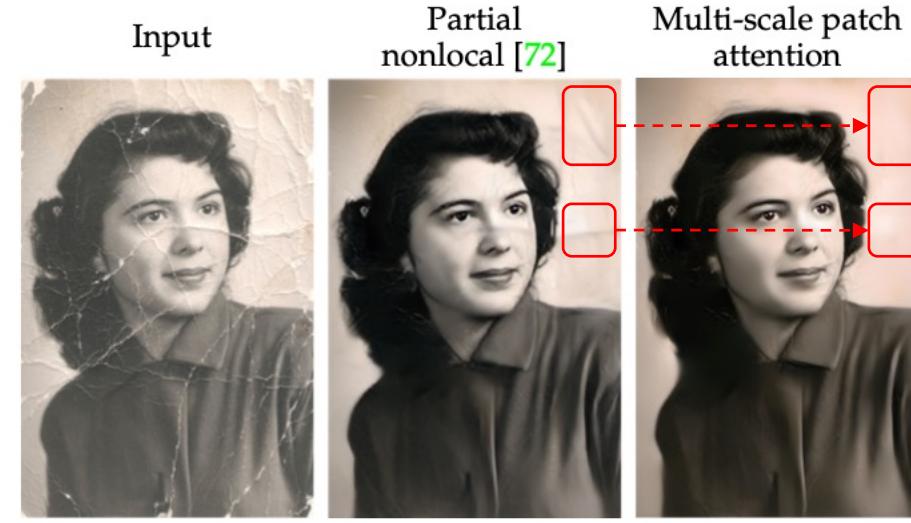


Fig. 16: Qualitative comparisons between two designed global branches. The multi-scale patch attention could better resolve the structured defects.

Method	PSNR ↑	SSIM ↑	LPIPS ↓	FID ↓
Remove Global	23.14	0.68	0.26	143.62
Nonlocal [64]	23.19	0.68	0.26	141.49
Contextual Attention [42]	23.25	<u>0.69</u>	<u>0.25</u>	138.27
Partial Nonlocal (CVPR 2020)	<u>23.33</u>	<u>0.69</u>	<u>0.25</u>	<u>134.35</u>
<b>Multi-scale Patch Fusion</b>	<b>23.45</b>	<b>0.71</b>	<b>0.24</b>	<b>130.64</b>

TABLE 5: Quantitative comparisons on the DIV2K dataset (256x256). **Bold**: Best. Underline: Second best.

# Summary

## Purpose

- Unpaired old (degraded) photo restoration

## Contributions

- Domain alignment: well generalization on real photos by training synthetic data
- Partial nonlocal block: reconstruction of blind spots / corruption area
- Multi-scale patch-based fusion: reduction of memory cost and better performance
- Coarse-to-fine generator: reconstruction of the high-resolution face

## Limitations

- Anisotropic pollution is not totally restored.
- Shading along the folds of the photo hindrance the restoration process.