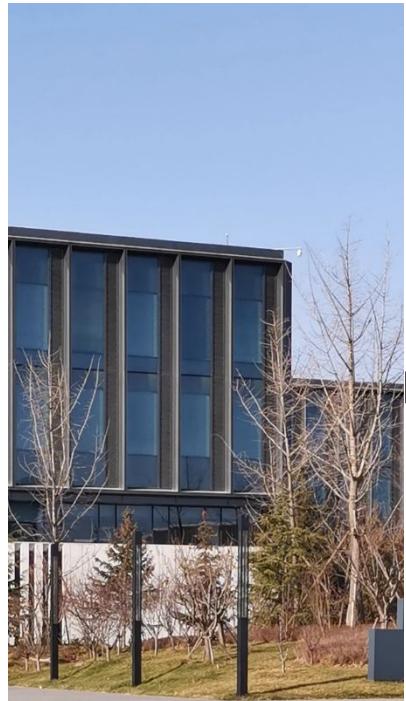


# Compound Quality Mapping on High-resolution Images and Videos

**Danda Pani Paudel & Zhiwu Huang**  
Computer Vision Lab @ ETH Zurich

# Compound Quality Mapping



Color Adjustment



Illumination Enhancement



Texture Sharpening etc....

# Data Collection for Compound Image and Video Quality Mapping

# Data Collection I: Paired Image Retouching (expensive expert effort)



Original

Retouched by expert-A

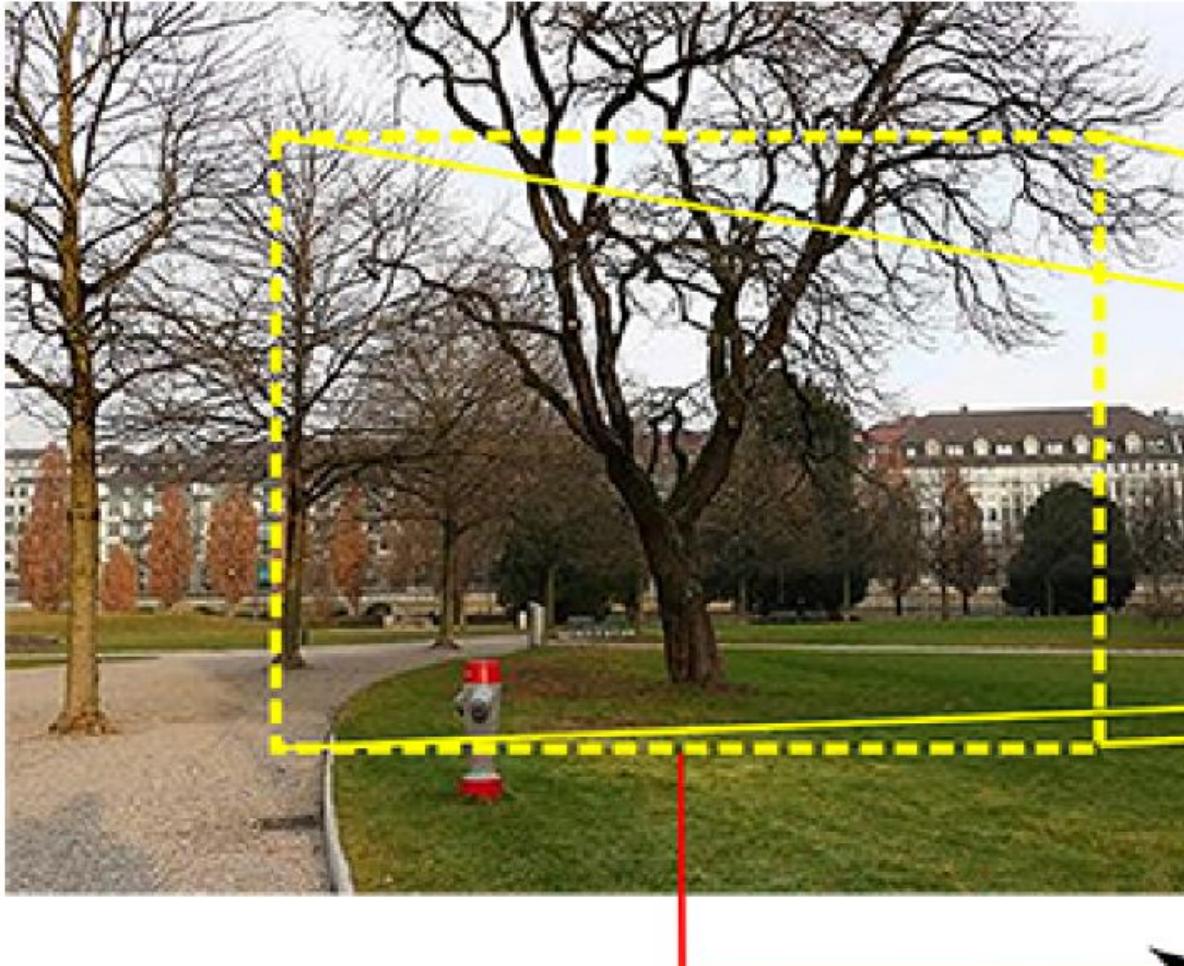
Retouched by expert-B

## Data Collection II: Weakly-paired Collection (expensive alignment)



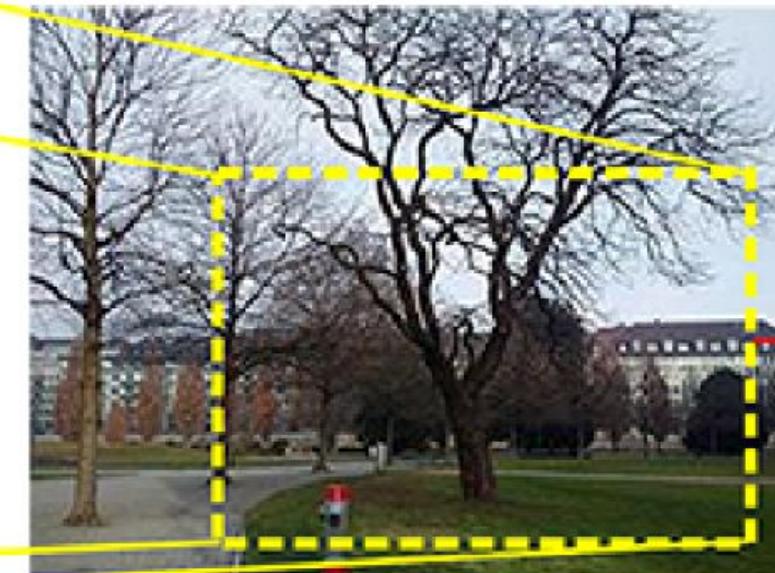
Camera	Sensor	Image size	Photo quality
<i>iPhone 3GS</i>	3 MP	2048 × 1536	Poor
<i>BlackBerry Passport</i>	13 MP	4160 × 3120	Mediocre
<i>Sony Xperia Z</i>	13 MP	2592 × 1944	Average
<i>Canon 70D DSLR</i>	20 MP	3648 × 2432	Excellent

*DSLR-Quality Photos on Mobile Devices with Deep Convolutional Networks, Ignatov et al., ICCV 2017.*



non-linear transform and a crop resulting in two images of the same resolution representing the same scene

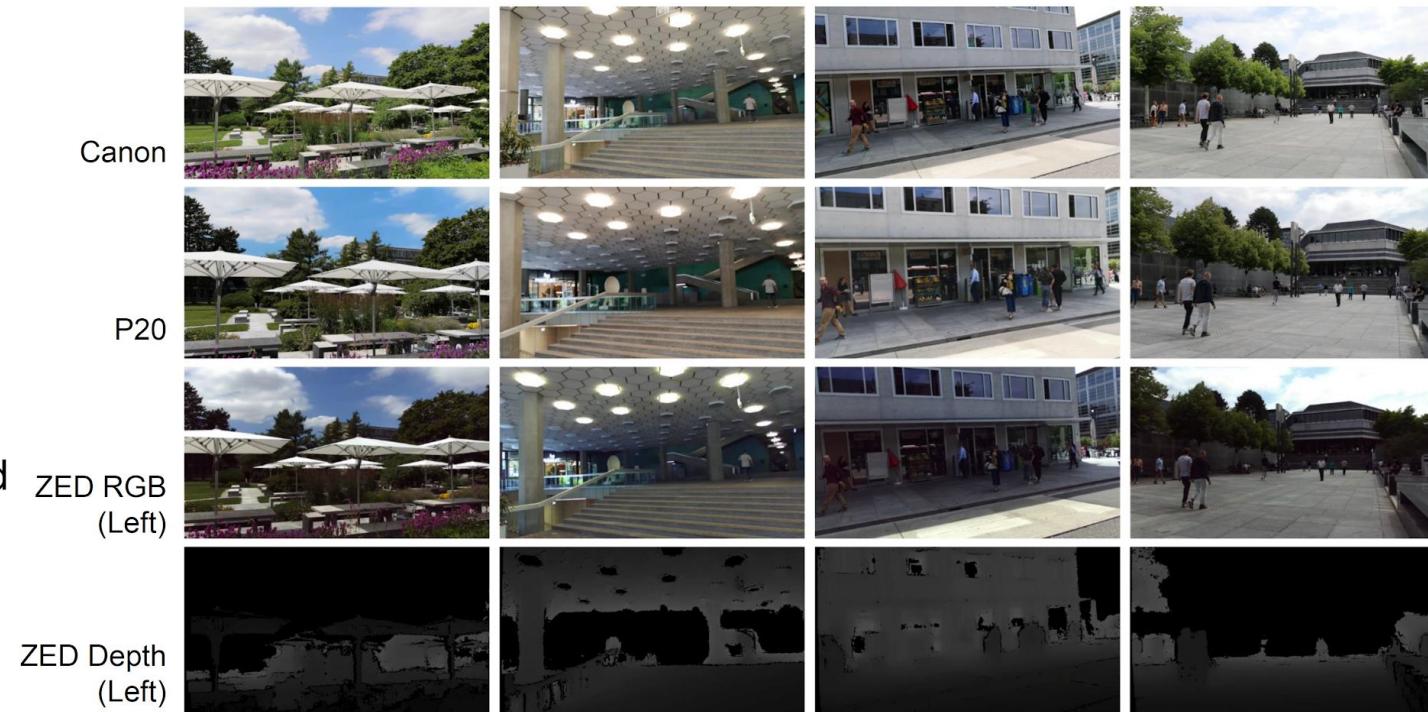
an overlapping region is determined by SIFT descriptor matching



## Data Collection II: Weakly-paired Video Collection? (more expensive)

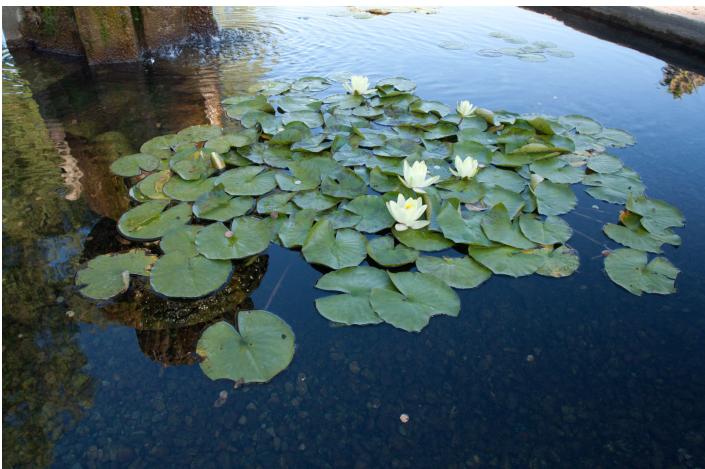
### Vid3oC Dataset

- Three cameras
- 82 recordings
- 328 videos
- Includes stereo depth
- **Canon 5D Mark IV** high quality DSLR
- **Huawei P20** high-end smartphone
- **ZED** stereo camera
- **AIM 2019** Video SR

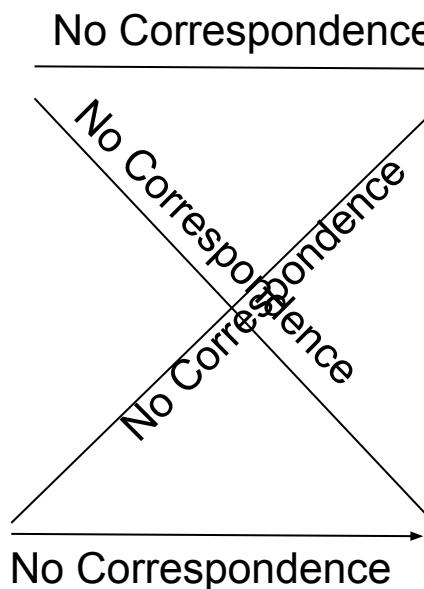


*The Vid3oC and IntVID Datasets for Video Super Resolution and Quality Mapping. Sohyeong Kim, Guanju Li, Dario Fuoli, Martin Danelljan, Zhiwu Huang, Shuhang Gu and Radu Timofte. ICCV 2019 Workshops.*

## Data Collection II: Unpaired Collection (Cheaper)



Poor-quality Image Dataset

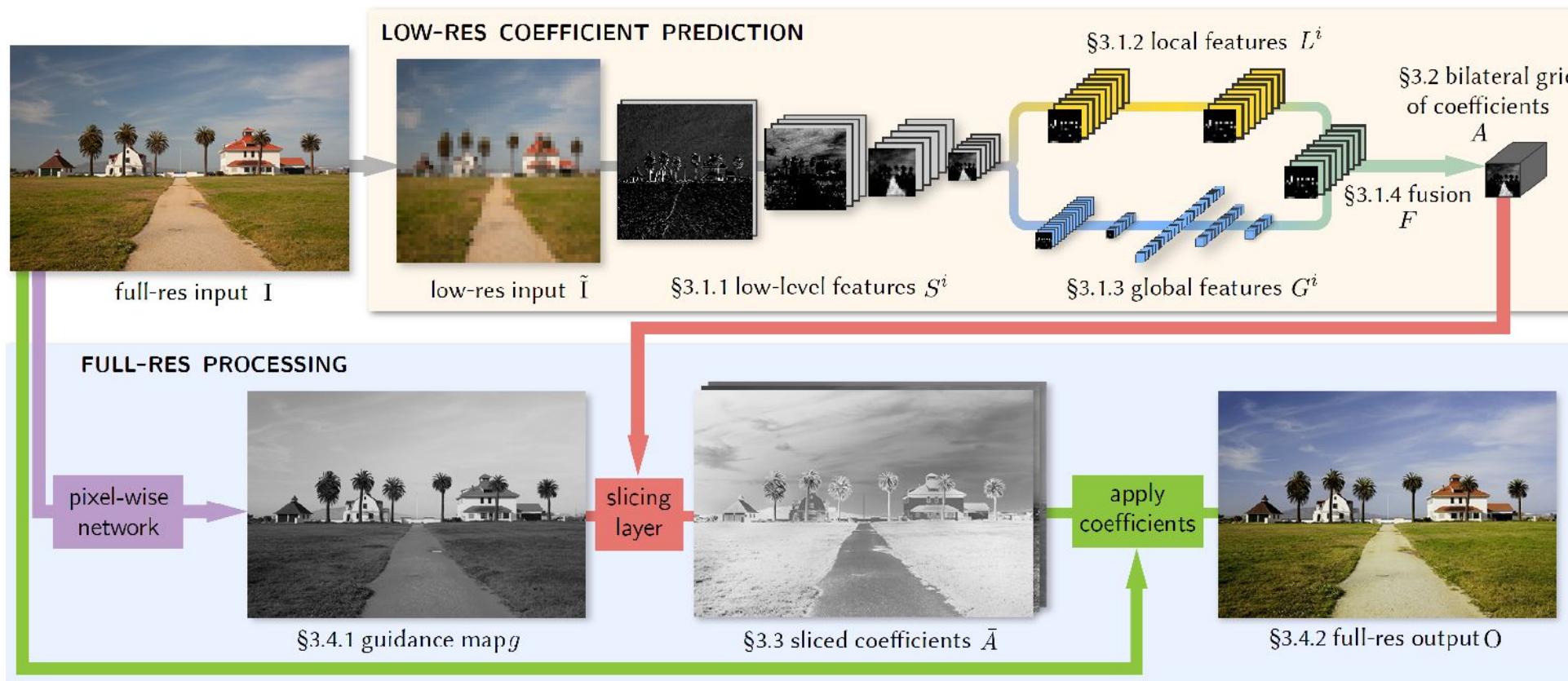


High-quality Image Dataset

# Supervised Deep Learning Methods for Compound Image Quality Mapping

# Deep Bilateral Learning for Real-Time Image Enhancement (HDRNet)

**Idea:** consumes a low-resolution version of the input image, followed by an edge-preserving upsampling to the full-resolution image in a bilateral filtering fashion



*Deep Bilateral Learning for Real-Time Image Enhancement, GHARBI et al., TOG 2017.*

# Deep Bilateral Learning for Real-Time Image Enhancement (HDRNet)



12 megapixel 16-bit linear input  
(tone-mapped for visualization)



tone-mapped with HDR+  
**400 – 600 ms**

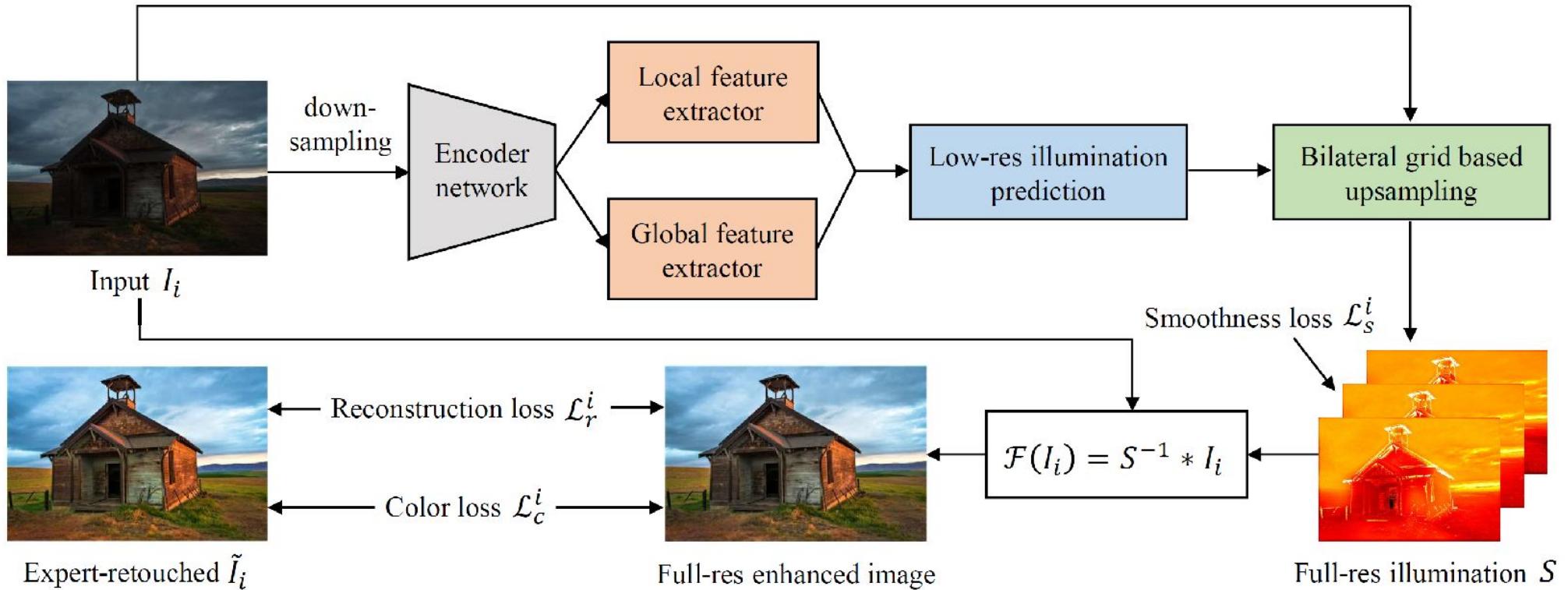


processed with our algorithm  
**61 ms, PSNR = 28.4 dB**

*Deep Bilateral Learning for Real-Time Image Enhancement, GHARBI et al., TOG 2017.*

# Underexposed Photo Enhancement (UPE)

**Idea:** learn an image-to-illumination (instead of image-to-image) mapping



*Underexposed Photo Enhancement using Deep Illumination Estimation, Wang et al., CVPR 2019.*

# Underexposed Photo Enhancement (UPE)

**Loss function:**

$$\mathcal{L} = \sum_{i=1}^N \omega_r \mathcal{L}_r^i + \omega_s \mathcal{L}_s^i + \omega_c \mathcal{L}_c^i$$

**Reconstruction**  $\mathcal{L}_r^i = \|I_i - S * \tilde{I}_i\|^2,$   
 $s.t. \quad (I_i)_c \leq (S)_c \leq 1, \quad \forall \text{ pixel channel } c$

**Smoothness**  $\mathcal{L}_s^i = \sum_p \sum_c \omega_{x,c}^p (\partial_x S_p)_c^2 + \omega_{y,c}^p (\partial_y S_p)_c^2$

**Color**  $\mathcal{L}_c^i = \sum_p \angle((\mathcal{F}(I_i))_p, (\tilde{I}_i)_p)$

*Underexposed Photo Enhancement using Deep Illumination Estimation, Wang et al., CVPR 2019.*

# Underexposed Photo Enhancement (UPE)



(a) Input



(b) JieP [4]



(c) HDRNet [13]



(d) DPE [9]



(e) White-box [15]



(f) Distort-and-Recover [22]



(g) Our result  
Visual Comparison on MIT-Adobe FiveK



(h) Expert-retouched

*Underexposed Photo Enhancement using Deep Illumination Estimation, Wang et al., CVPR 2019.*

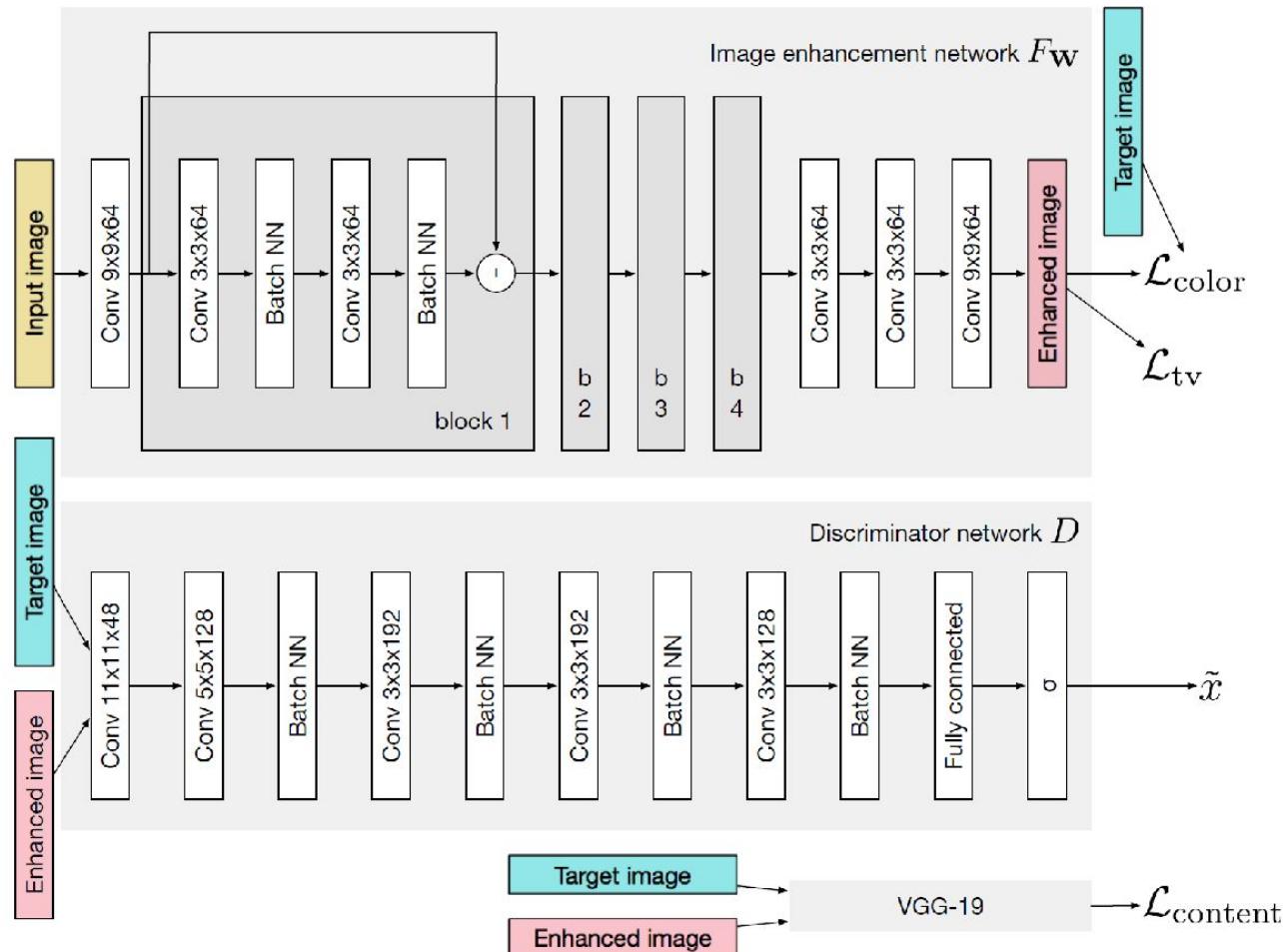
## Underexposed Photo Enhancement (UPE)

Method	PSNR	SSIM
HDRNet [2]	21.96	0.866
DPE [1]	22.15	0.850
White-Box [3]	18.57	0.701
Distort-and-Recover [4]	20.97	0.841
Ours w/o $\mathcal{L}_r$ , w/o $\mathcal{L}_s$ , w/o $\mathcal{L}_c$	21.97	0.867
Ours with $\mathcal{L}_r$ , w/o $\mathcal{L}_s$ , w/o $\mathcal{L}_c$	22.31	0.871
Ours with $\mathcal{L}_r$ , with $\mathcal{L}_s$ , w/o $\mathcal{L}_c$	22.89	0.884
Ours	<b>23.04</b>	<b>0.893</b>

Quantitative Comparison on MIT-Adobe FiveK

*Underexposed Photo Enhancement using Deep Illumination Estimation, Wang et al., CVPR 2019.*

# DSLR Photo Enhancement (DSLR-PE)



**Idea:** learn the translation function using a residual convolutional neural network with a composite perceptual error function that combines content, color and adversarial texture losses

$$\begin{aligned} \mathcal{L}_{total} = & \mathcal{L}_{content} + 0.4 \cdot \mathcal{L}_{texture} + 0.1 \cdot \mathcal{L}_{color} \\ & + 400 \cdot \mathcal{L}_{tv} \end{aligned}$$

$$\mathcal{L}_{color}(X, Y) = \|X_b - Y_b\|_2^2$$

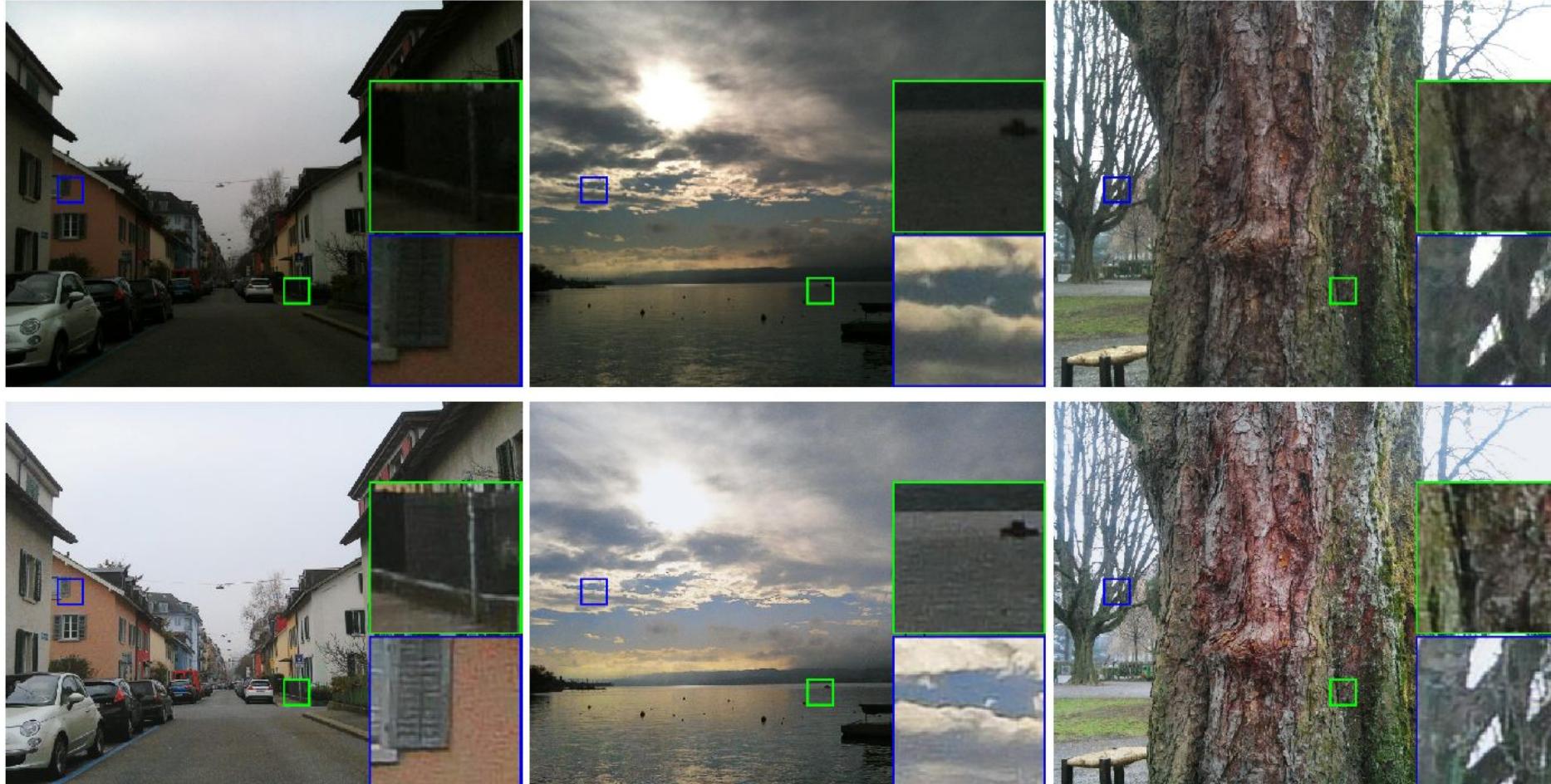
$$\mathcal{L}_{content} = \frac{1}{C_j H_j W_j} \|\psi_j(F_W(I_s)) - \psi_j(I_t)\|$$

$$\mathcal{L}_{tv} = \frac{1}{CHW} \|\nabla_x F_W(I_s) + \nabla_y F_W(I_s)\|$$

$$\mathcal{L}_{texture} = - \sum_i \log D(F_W(I_s), I_t)$$

DSLR-Quality Photos on Mobile Devices with Deep Convolutional Networks, Ignatov et al., ICCV 2017.

## DSLR Photo Enhancement (DSLR-PE)



Typical artifacts generated by our method (bottom) compared with original iPhone images (top)

DSLR-Quality Photos on Mobile Devices with Deep Convolutional Networks, Ignatov et al., ICCV 2017.

# Weakly-Supervised Deep Learning Methods for Compound Image Quality Mapping

# Weakly Supervised Photo Enhancer (WESPE)

-Content Consistency Loss

$$\mathcal{L}_{\text{content}} = \frac{1}{C_j H_j W_j} \|\psi_j(x) - \psi_j(\tilde{x})\|,$$

-Adversarial Color Loss

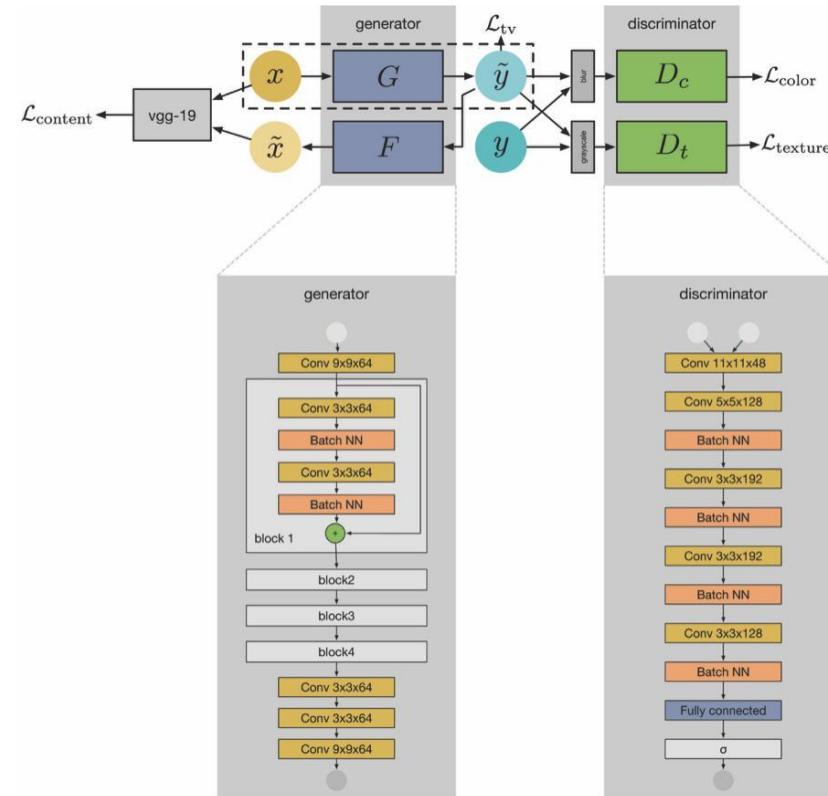
$$\mathcal{L}_{\text{color}} = - \sum_i \log D_c(G(x)_b).$$

-Adversarial Texture Loss

$$\mathcal{L}_{\text{texture}} = - \sum_i \log D_t(G(x)_g).$$

-Total Variation Loss

$$\mathcal{L}_{\text{tv}} = \frac{1}{CHW} \|\nabla_x G(x) + \nabla_y G(x)\|,$$



*WESPE: weakly supervised photo enhancer for digital cameras, Ignatov et al., CVPRW 2018.*

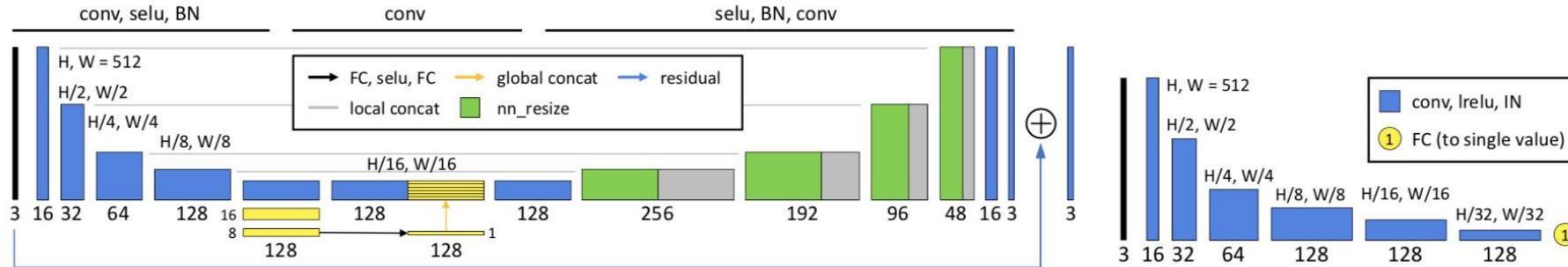
# Weakly Supervised Photo Enhancer (WESPE)

CVL

DPED images	APE		Weakly Supervised				Fully Supervised	
	PSNR	SSIM	WESPE [DIV2K]	WESPE [DPED]	PSNR	SSIM	[13]	PSNR
iPhone	17.28	0.86	17.76	0.88	18.11	0.90	<b>21.35</b>	<b>0.92</b>
BlackBerry	18.91	0.89	16.71	0.91	16.78	0.91	<b>20.66</b>	<b>0.93</b>
Sony	19.45	0.92	20.05	0.89	20.29	0.93	<b>22.01</b>	<b>0.94</b>

*WESPE: weakly supervised photo enhancer for digital cameras, Ignatov et al., CVPRW 2018.*

# Deep Photo Enhancer (DPE)



-Identity Mapping Loss:

$$I = \mathbb{E}_{x,y'} [MSE(x, y')] + \mathbb{E}_{y,x'} [MSE(y, x')].$$

-Cycle-Consistency Loss:

$$C = \mathbb{E}_{x,x''} [MSE(x, x'')] + \mathbb{E}_{y,y''} [MSE(y, y'')].$$

-Adversarial Loss:

$$\begin{aligned} A_D &= \mathbb{E}_x [D_X(x)] - \mathbb{E}_{x'} [D_X(x')] + \\ &\quad \mathbb{E}_y [D_Y(y)] - \mathbb{E}_{y'} [D_Y(y')], \\ A_G &= \mathbb{E}_{x'} [D_X(x')] + \mathbb{E}_{y'} [D_Y(y')]. \end{aligned}$$

# Deep Photo Enhancer (DPE)



Input

Enhanced by DPE

	CycleGAN	DPED	NPEA	CLHE	ours	total
CycleGAN	-	32	27	23	11	93
DPED	368	-	141	119	29	657
NPEA	373	259	-	142	50	824
CLHE	377	281	258	-	77	993
ours	389	371	350	323	-	1433

Preference matrix from AMT user study

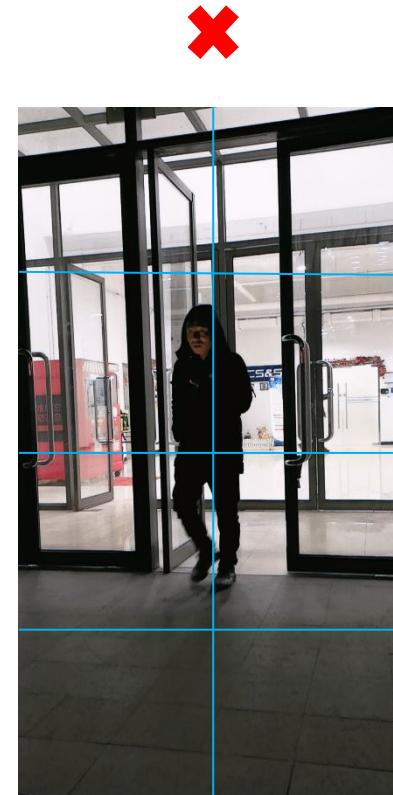
# Limitation for Compound Quality Mapping and High-Resolution Image Treatment



Input



Downscaling  
(low-res, noisy,  
blurry)



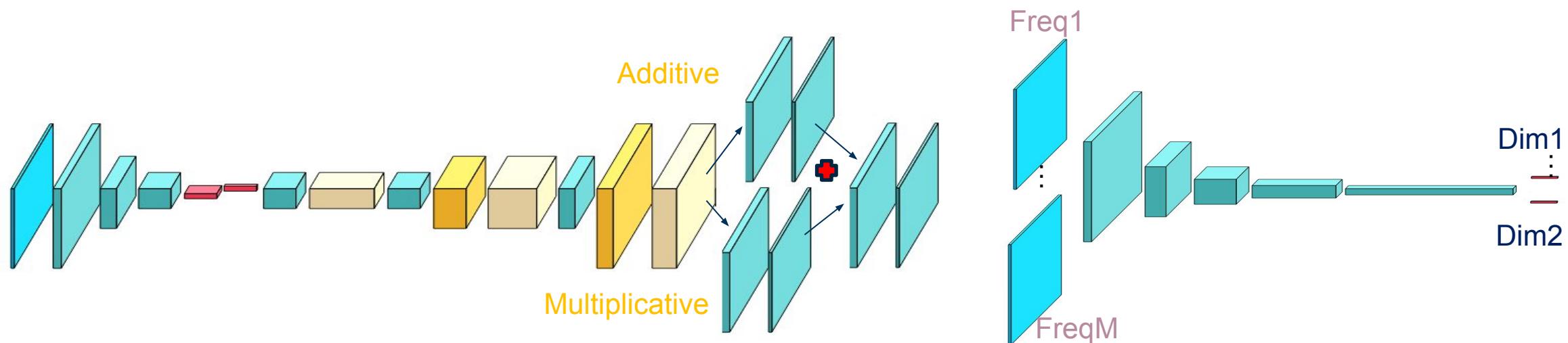
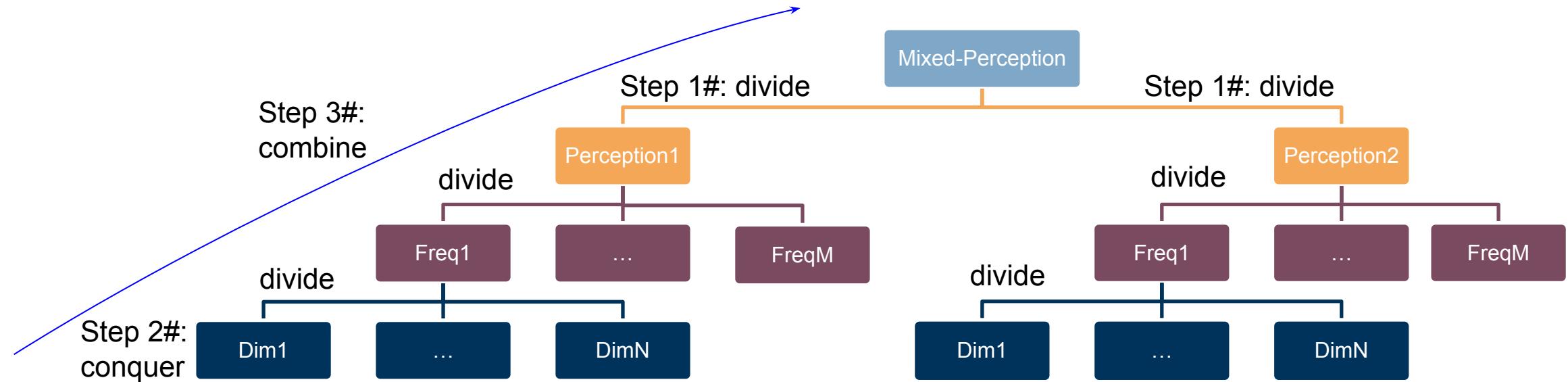
Patch-wise Enhancement  
(spatial inconsistency)

Model	Limitation (Compound Quality)	Limitation (High Resolution)
WESPE	Color and Texture(Not sufficient)	Patch-wise Enhancement
DPE	No consideration on mixed-percept ual improvement	Down-scaling

# Divide-and-Conquer Adversarial Learning for High-resolution Image and Video Enhancement

*Zhiwu Huang, Danda Pani Paudel, Guanju Li, Jiqing Wu, Radu Timofte, Luc Van Gool, arXiv preprint arXiv:1910.10455.*

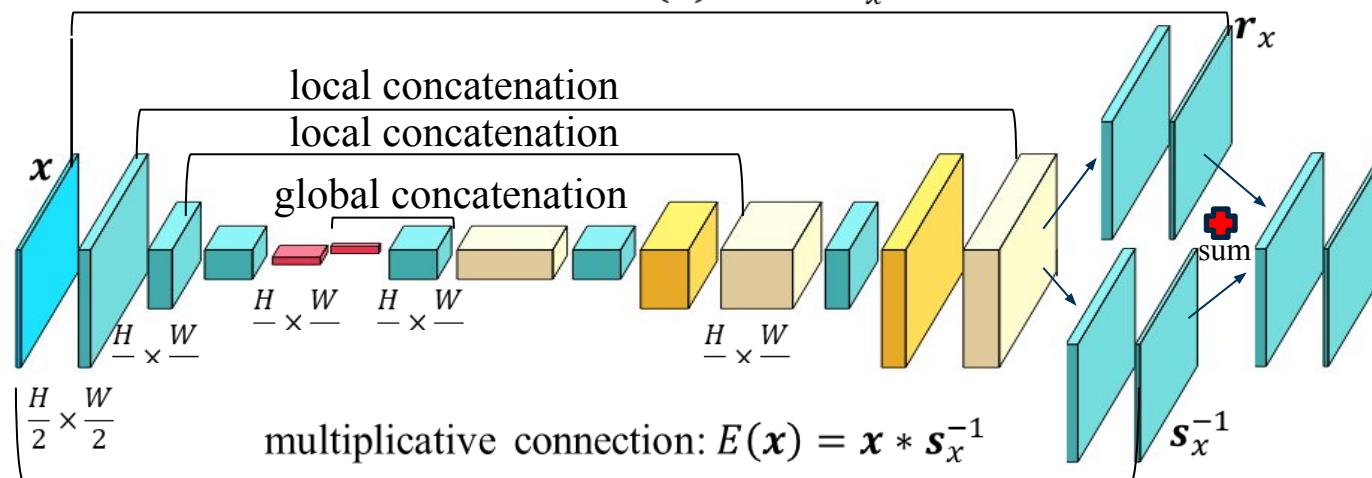
# Divide-and-Conquer Inspired Method



# Network Design

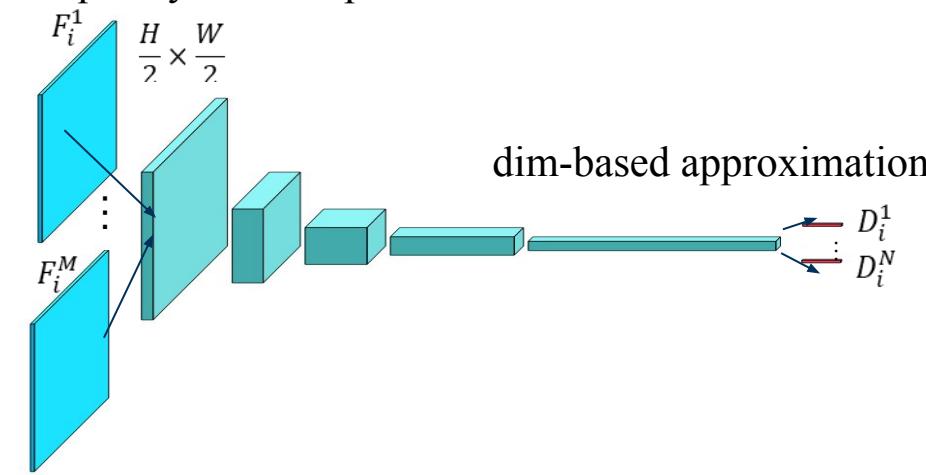
█ Input/output  
 █ Conv+ReLU+BN  
 █ Pooling / OrthoProj  
 █ Resize  
 █ Concat

additive connection:  $E(\mathbf{x}) = \mathbf{x} + \mathbf{r}_x$



(a) Enhancer for Perception-based Division

frequency-based input



(b) Discriminator for Freq- and Dim-based Division

# Perception-based Division



(a) Input



(b) Additive component



(c) Multiplicative component



(d) Additive map

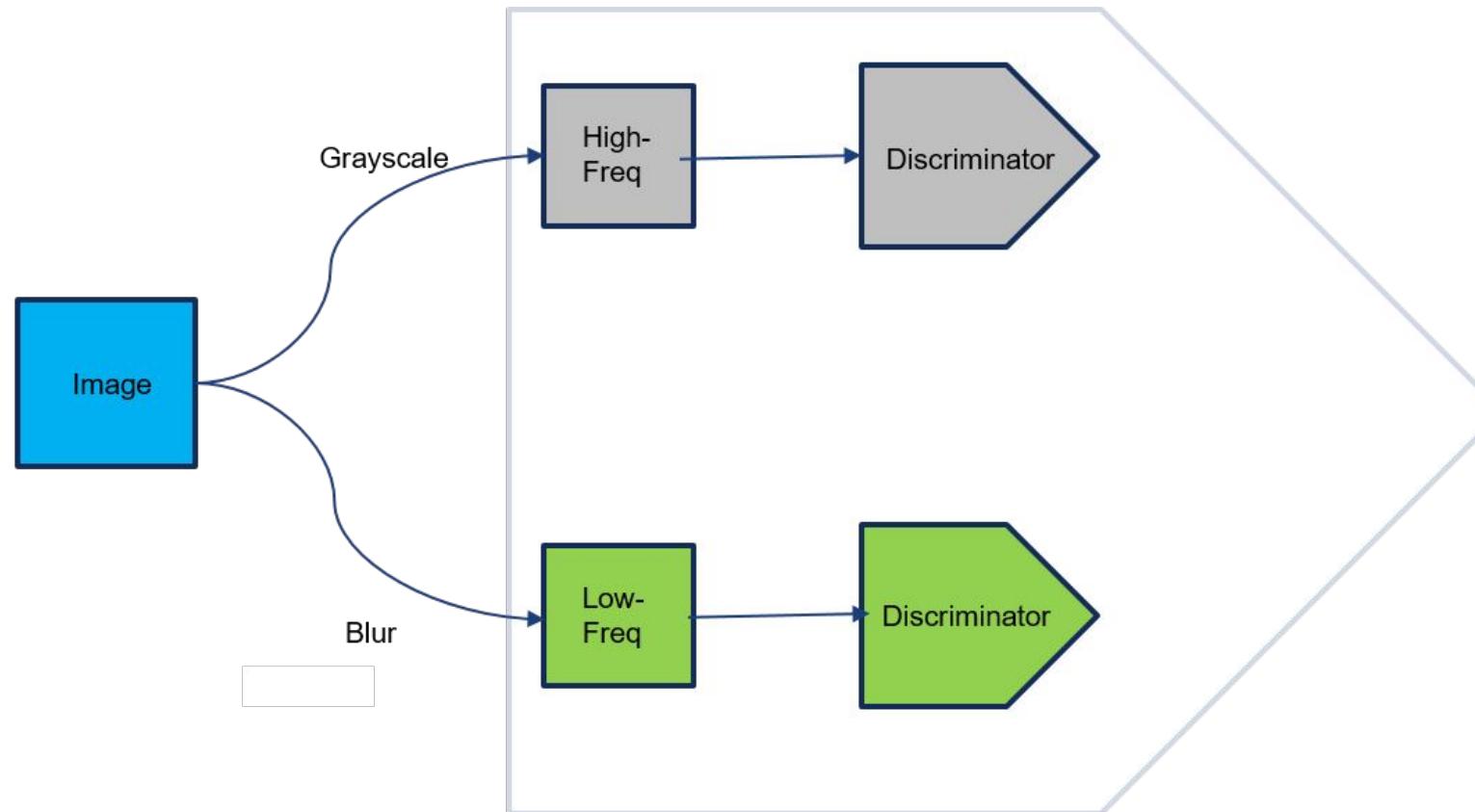


(e) Multiplicative map



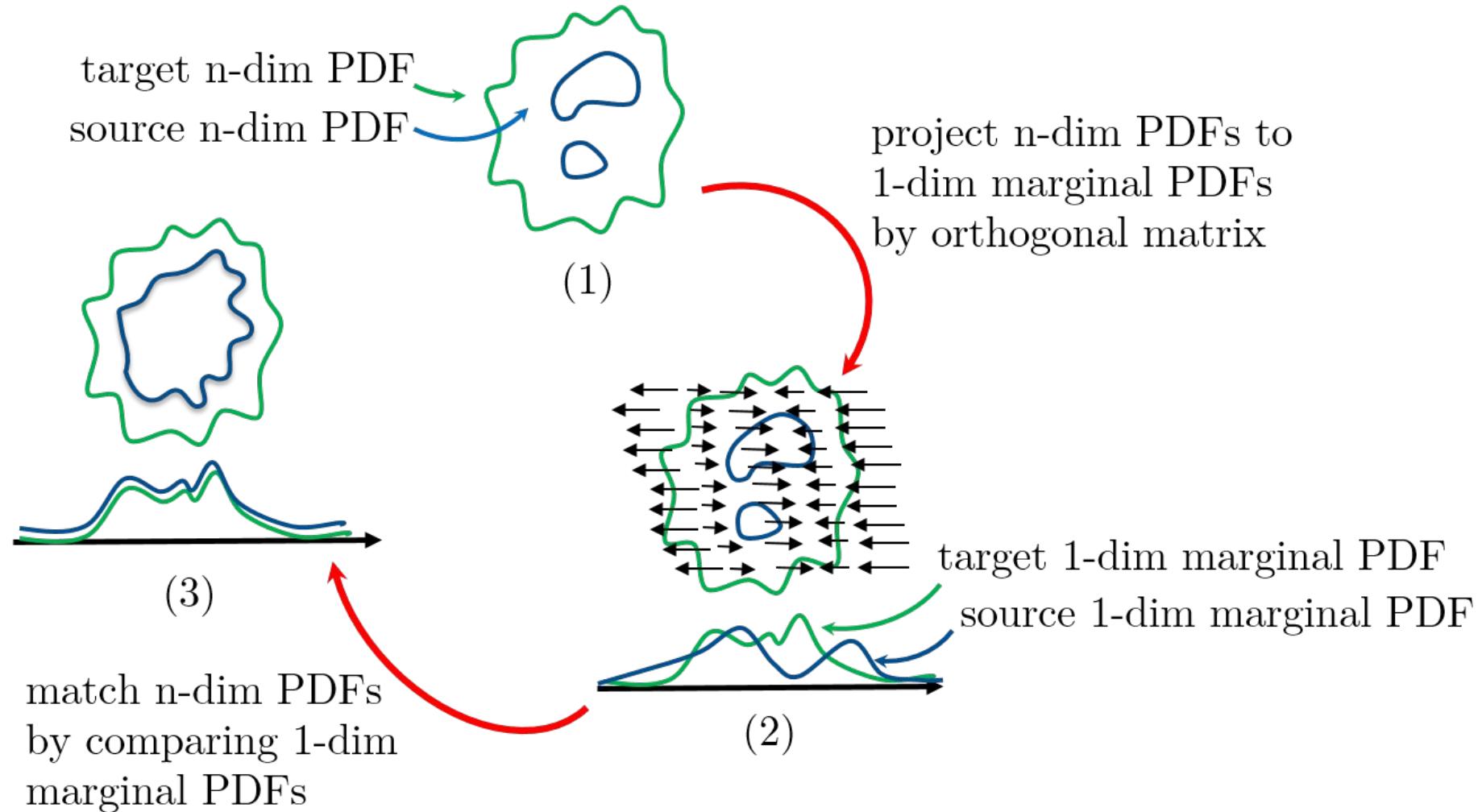
(f) Fused Map

# Frequency-based Division



- [1] Andrey Ignatov, Nikolay Kobyshev, Radu Timofte, Kenneth Vanhoey, and Luc Van Gool. Dslr-quality photos on mobile devices with deep convolutional networks.  
[2] Anoosheh, Asha, et al. "Night-to-day image translation for retrieval-based localization."

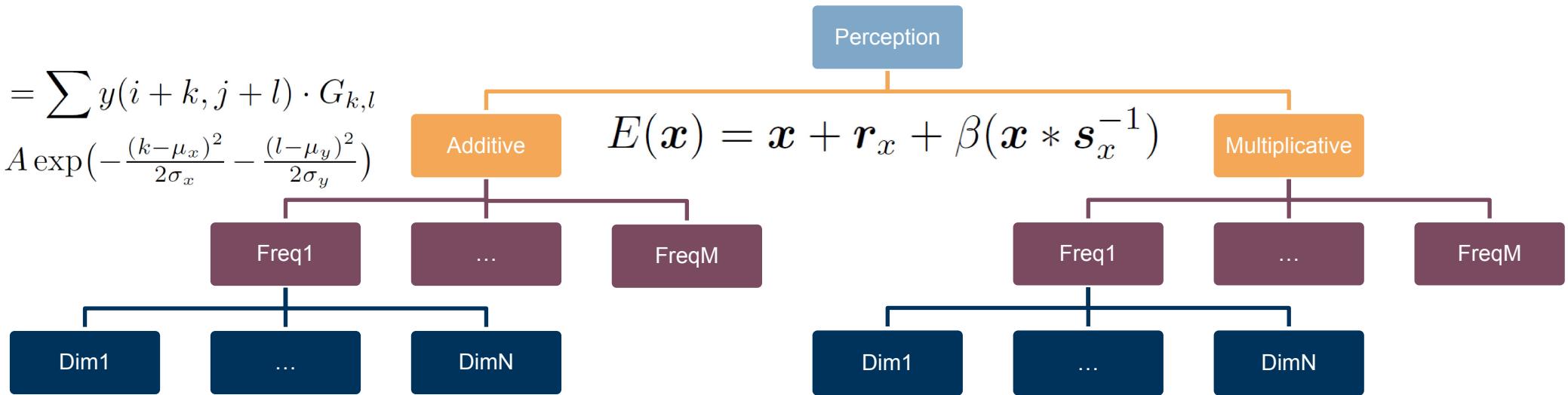
## Dimension-based Division and Optimization



# Loss Design

$$y_b(i, j) = \sum y(i + k, j + l) \cdot G_{k,l}$$

$$G_{k,l} = A \exp\left(-\frac{(k-\mu_x)^2}{2\sigma_x^2} - \frac{(l-\mu_y)^2}{2\sigma_y^2}\right)$$



## Adaptive SWGAN loss

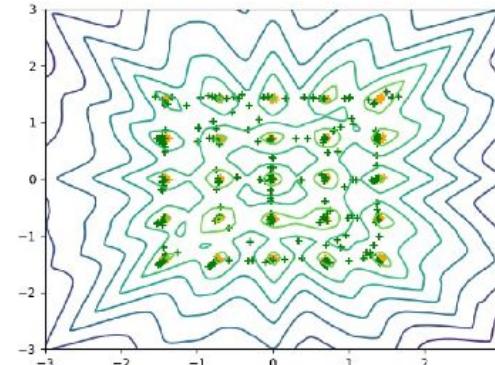
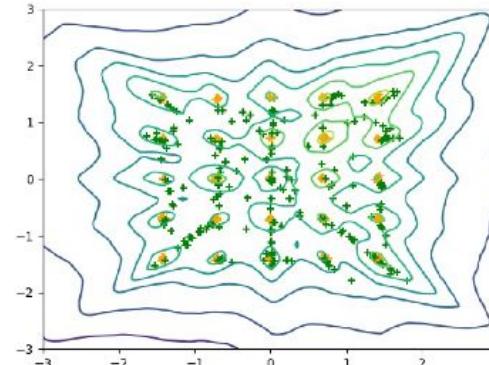
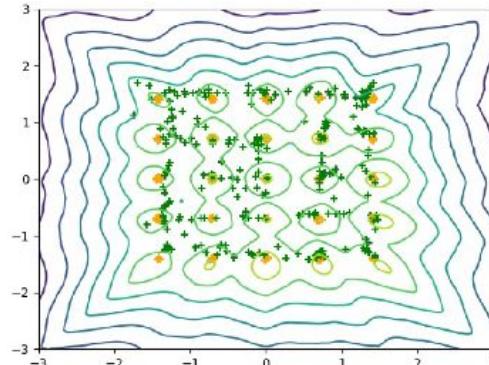
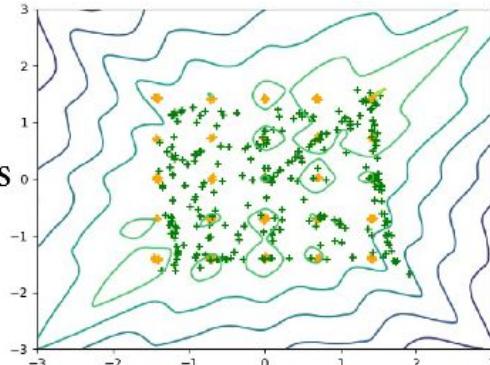
$$\min_E \max_C \int_{\theta \in \mathbb{S}^{n-1}} \left( \mathbb{E}_{\mathbf{y} \sim P_y} [C(\mathbf{y})] - \mathbb{E}_{\hat{\mathbf{y}} \sim P_E} [C(E(\mathbf{x}))] \right) + \lambda \mathbb{E}_{\hat{\mathbf{y}} \sim P_{\hat{\mathbf{y}}}} [\max(0, \|\nabla_{\hat{\mathbf{y}}} C(\hat{\mathbf{y}})\|_2 - 1)]$$

$$\overline{\nabla_{\hat{\mathbf{y}}} C(\hat{\mathbf{y}})} = \eta \overline{\nabla_{\hat{\mathbf{y}}} C(\hat{\mathbf{y}})} + (1 - \eta) \frac{\nabla_{\hat{\mathbf{y}}} C(\hat{\mathbf{y}})}{\lambda}$$

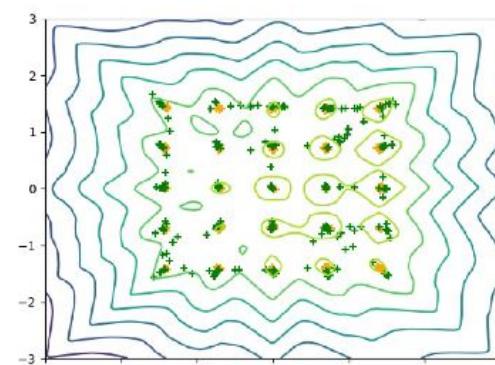
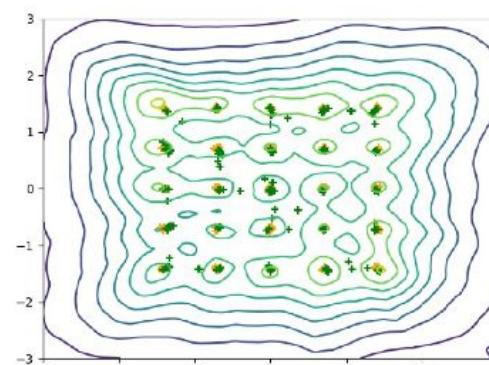
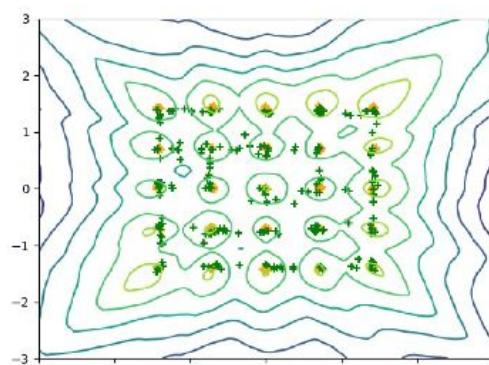
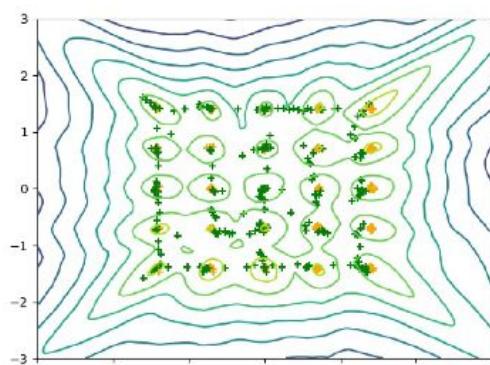
Adaptive Penalty

# Evaluation on Toy Data

2.5k iterations



5k iterations



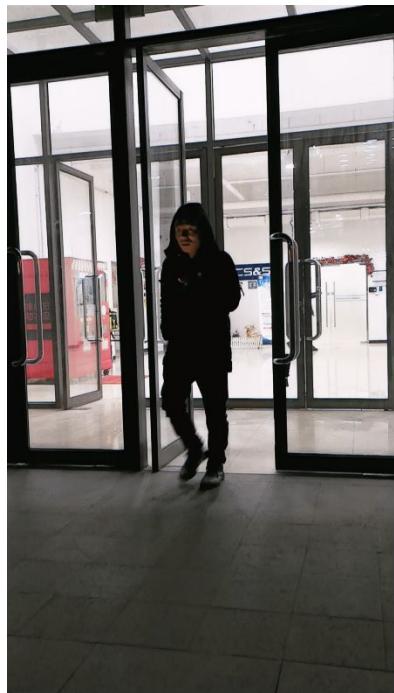
(a) WGAN

(b) AdaWGAN

(c) SWGAN

(d) Proposed AdaSWGАН

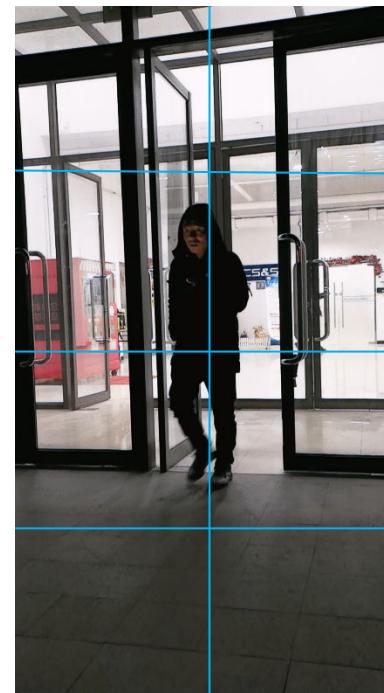
# High-resolution Issue for Image Enhancement



Input



Downscaling  
(low-res, noisy,  
blurry)  
*Deep Photo  
Enhancer (DPE)*  
*[Chen et al in  
CVPR '18]*



Patch-wise Enhancement  
(spatial inconsistency)  
*Weakly Supervised Photo  
Enhancer (WESPE) [our work  
in CVPR '18 workshop]*



Multi-scale Photo  
Enhancement (MUSPE)  
*Our current work*

coarse  
fine

# Multi-scale Extension of DACAL for Image Enhancement

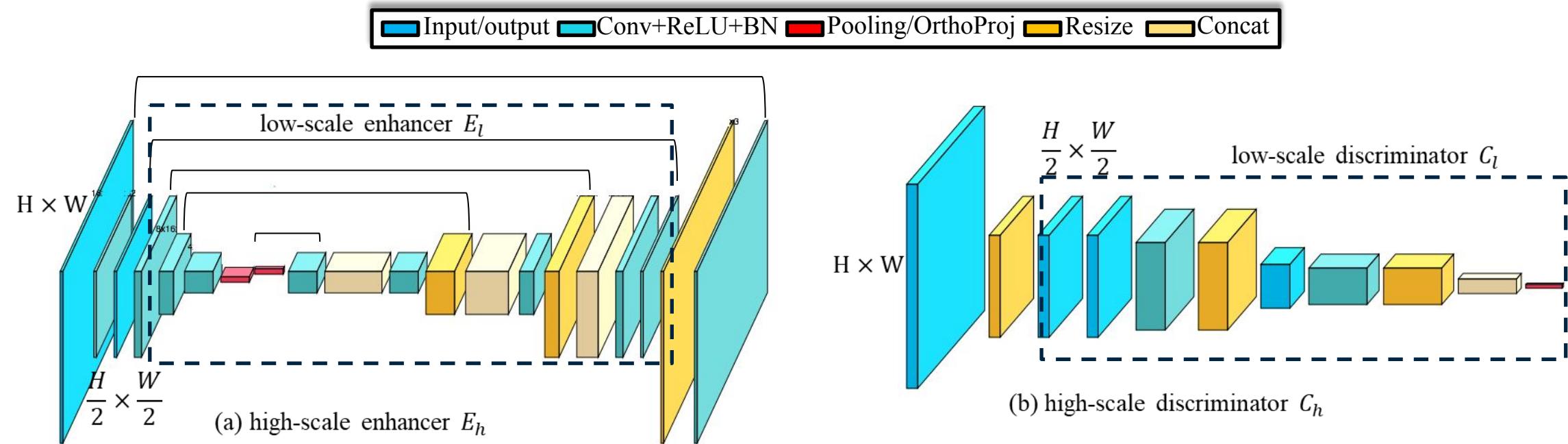


Table 1: PSNR and SSIM results for the MIT-Adobe FiveK [42] test images. Here, WB and DR indicate the White-Box and Distort-and-Recover methods, respectively.  $MUSPE_{l_1}$ ,  $MUSPE_{l_2}$ ,  $MUSPE_{l_3}$  and  $MUSPE_l$  represent the use of individual additive, individual multiplicative, multiplicative cascaded by additive, and our suggested parallel fusion (two-stream strategy), respectively.  $MUSPE_h$  is our higher-scale version.  $PSNR_d/SSIM_d$  and  $PSNR_f/SSIM_f$  indicate the results on downsampled images and full-resolution images, respectively.

	WB	DR	DPED	DPE	$MUSPE_{l_1}$	$MUSPE_{l_2}$	$MUSPE_{l_3}$	$MUSPE_l$	$MUSPE_h$
$PSNR_d$	18.86	21.64	21.05	22.10	22.73	22.99	23.01	23.52	<b>24.15</b>
$PSNR_f$	19.09	21.52	20.86	21.65	22.43	22.69	23.02	23.56	<b>24.07</b>
$SSIM_d$	0.928	0.936	0.922	0.947	0.958	0.942	0.949	0.959	<b>0.962</b>
$SSIM_f$	0.920	0.922	0.916	0.894	0.948	0.942	0.940	0.954	<b>0.956</b>

Table 2: PSNR and SSIM results for the DPED [14] test  $100 \times 100$  image patches. Here,  $l, f, d$  for MUSPE represent the use of our proposed sliced-perception, sliced-frequency and sliced-dimension learning respectively.  $MUSPE_h$  is our higher-scale version.

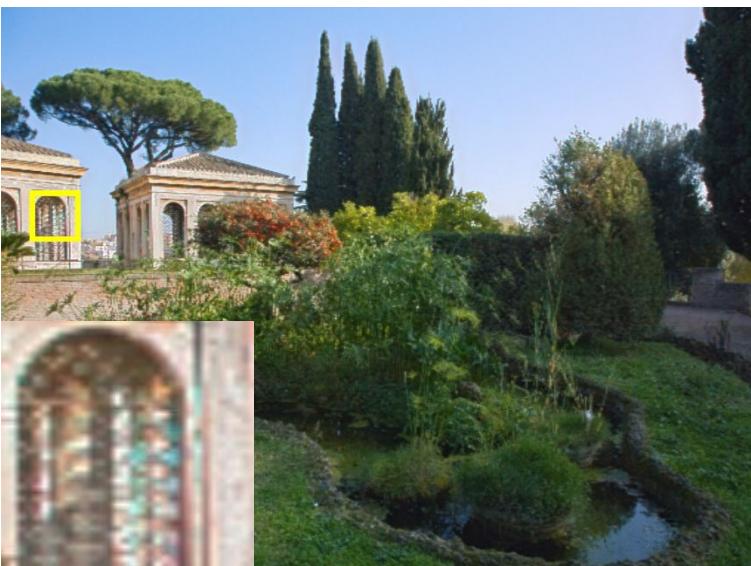
	WESPE	DPE	$MUSPE_l$	$MUSPE_{l+f}$	$MUSPE_{l+f+d}$	$MUSPE_h$
$PSNR_{100}$	17.45	18.53	19.62	20.01	20.43	<b>20.90</b>
$SSIM_{100}$	0.854	0.861	0.868	0.869	0.872	<b>0.874</b>



Input



WESPE [Ignatov, CVPRW'18]



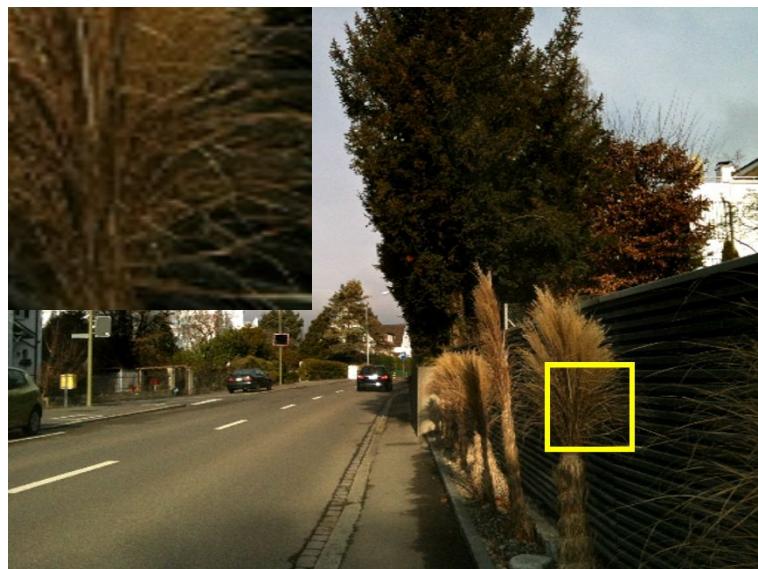
DPE [Chen, CVPR'18]



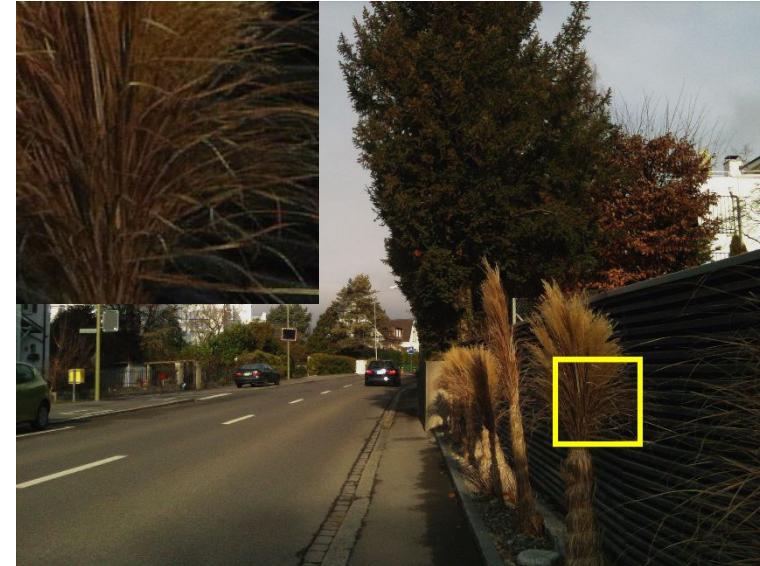
Proposed MUSPE [ICLR'20 submission]



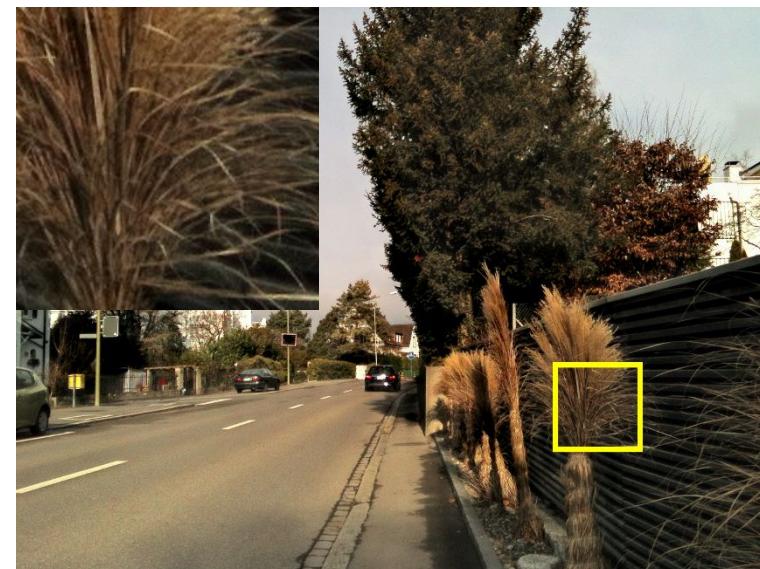
Input



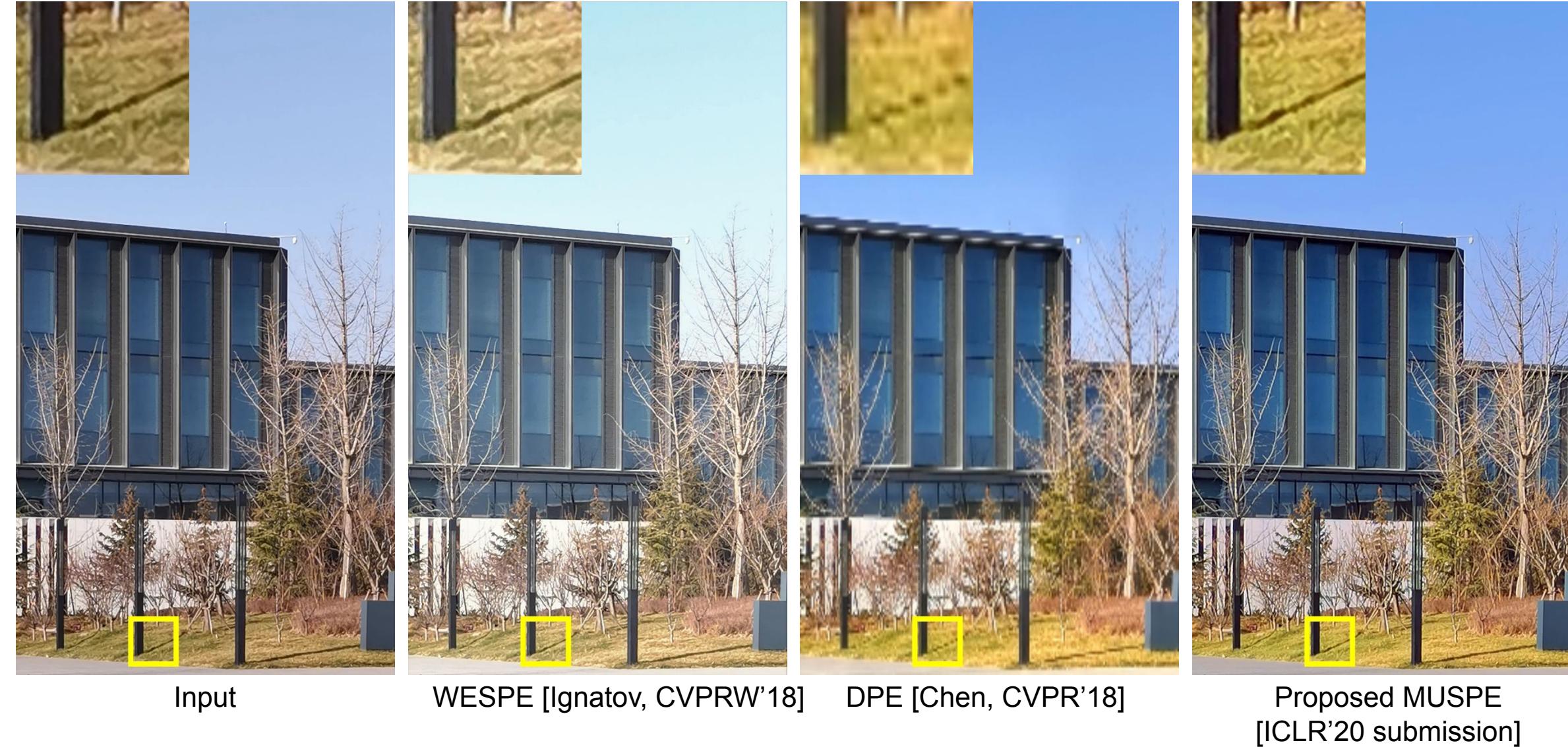
DPE [Chen, CVPR'18]



WESPE [Ignatov, CVPRW'18]

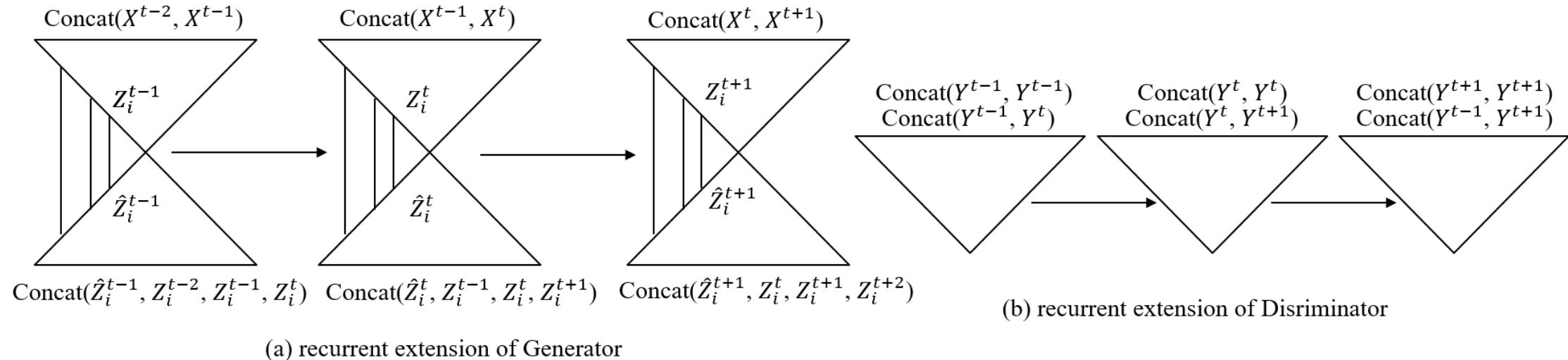


Proposed MUSPE [ICLR'20 submission]



**Video Quality Mapping = Image Quality Mapping + Temporal Smoothing**

# Recurrent Extension of DACAL for Video Enhancement





Perframe-DACAL



Recurrent-DACAL



Perframe-DACAL



Recurrent-DACAL (fine-tuned on Retouched&amp;DSLR images)



Perframe-DACAL



Recurrent-DACAL



Perframe-DACAL



Recurrent-DACAL (fine-tuned on Retouched&amp;DSLR images)

# Conclusion

- Supervision
  - Weak-supervision is cheaper
- Compound quality mapping
  - Divide-and-conquer inspired algorithm is promising
- High-resolution image treatment
  - Multiscaled training is helpful
- Video enhancement
  - Recurrent model works well