# Improve Mobile Phone's Image Quality by Deep Learning

Introduction & Review ·
Augmented Dataset ·
MPSE(Multi-Phase Semantic Enhancement) ·
MDIE(Multi-Domain Intrinsic Enhancement) ·
Future Work·

 $\cdot Q\&A \cdot$ 

HUANG Cen | LIU Yicun





# Problem Introduction

#### Mobile Phone vs. DSLR

- Small sensors, low-aperture optics
- Low-resolution, lack of details and contrast
- Poor dynamic range, high noise level
- Fall behind their DSLR counterparts





# Problem 1: Misalignment in PR Images

- Real-world Paired Data: Distortion and Misalignment
- End-to-End Training: Strait Registration is Essential
- Method Covered: SIFT-Feature Detection, Chessboard Registration, ECC Alignment, Optical Flow, Patch-Match and NNF







### Problem 2: Lack of Semantic Supervision in PR Training

#### Pairwise Enhancement Scheme:

• Random sampling, fixed-size cropping of small patches

#### Semantic Information Loss:

• Same structural details are enhanced differently

#### Major Existing Problems:

• Block artifacts, structure distortion, unbalanced enhancement





# Problem 3: Ambiguity in UPR Training

- Unpaired Methods: Ambiguity in training directions
- Lack of high-level structural supervision: Mismatching texture
- Divergence in weak supervision: weak constraints need to be modified





### Sub-Domain Problems

- Image Super-Resolution: SRCNN, FSRCNN, ESPCN, SR-GAN, ResNet, EDSR
- Image Denoising and Restoration: MLP, Auto Encoder, RedNet, BM3D, DND Benchmark Dataset
- Contrast Enhancement and Deblurring: HE, Dark Channel Prior, LIME, HDRNet, Gamma Correction
- Intrinsic Image Decomposition: Intrinsic in the Wild, Self-Supervised Intrinsic Decomposition
- Image-to-Image Transformation: Dual-Domain (Cycle-GAN, Conditional-GAN), Multi-Domain(Star-GAN)
- Photorealistic Style Transfer: Deep Image Analog, Deep Semantic Style, and MPSE
- Data-Driven Mixed Enhancement: DPED, DPED+, WESPE







## Synthesized Data vs. Real-world Data

- Image Super-Resolution: Deterministic Blur Kernel (Bicubic, Bilinear and etc.)
- Image Denoising and Restoration: Gaussian White Noise, Poisson Noise
- Model trained on synthesized data performs poorly in real-world blind denoising challenge
- Denoise Demo:
  - (a) RedNet trained on synthesized data, generated Gaussian white noise
  - (b) RedNet trained on real-world data, obtained by multi-frame denoise
- DND Denoising Benchmark: Approximate the Zero Noise ground-truth



(b)

# Sampling Paired Data

- Pro Camera
  - ≻ SONY A7,
  - Full-Frame Sensor with HDR
  - ➢ (6 stops of dynamic range)
  - $\blacktriangleright$  55 mm lens, 2x zooming relative to iPhone
  - ➤ (29 mm equivalent)
  - ➢ 24 mega pixels, nearly no noise
- Mobile Phone
  - ➢ iPhone 6s, normal capture mode, no HDR
  - > 29 mm lens, limit to its sensor and lens
  - ▶ 12 mega pixel, noisy, dark, lack of detail
- Registration
  - Adopt SIFT in patch level
  - Imperfect Alignment



## Our Dataset

- DIV 2K Dataset (with Night View Enrichment), 800 Images of 2K Resolution
- **PR Dataset**, DPED(7G) + Night View Enrichment(3G)
- UPR Dataset, 1500 HQ Pro Images + 300 LQ iPhone Images
- DND Denoising Benchmark Dataset, 50+ Real-Life Images with Ground-truth of Nearly Zero Noise
- MPI Sintel Dataset, Albedo and Original Image
- Adobe-5K Dataset, 5000 Raw Images with Human Retouch of 6 Experts (over 30 categories)

	004			<b>N</b>	607	<b>1</b>	D	0000 10000 251915	2	3					3503002220 (71	11314013045 (1	in solice		-		/ <b>0</b>	7 0
	(C)			(C)	Ne		-		N.C	"	3	4	- <u>1</u>	•	661e2ca_b	9f84941a_b	DSC07254	DSC07255	DSC07256	DSC07257	frame_0013	frame_0014
	ON				BR	DIS	11	15	18	17	20083878656_8d	20141042773_09	20171813788_a7	20345782054_a6	20372556141_ec	20386079443_7f	DSC07262	DSC07263	DSC07264	DSC07265	frame_0019	frame_0020
117 1111	014				650	104	89,1 m	80 <u>1</u> 100	CCCC EXCENSION S2:10					and the second							frame 0025	frame 0026
						9002 11-11-1	22_1 (R)	20,1 (0)	22,1 (13)	201010	21169447826_3f 9e4de399_b	21187742973_3f d4c999fd_b	21216561193_7e 537d5300_b	21235785831_3c f701b423_b	21277136892_dd 629f7c2b_b	21345746701_bb 40ca5f57_b	DSC07270	DSC07271	DSC07272	DSC07273		
HOD						000	80.08	20100	20_1 (17)	20,108	21752367029_6d 576abef2_b	21763805803_87 83dc093e_b	21805306620_2b cd1b85db_b	21810537134_b4 cb728cee_b	21921509076_26 5372f8f5_b	21924439119_4e 53161f43_b	DSC07278	DSC07279	DSC07280	DSC07281		frame_0032
,500 TH	Jucity				JICON JICON			1	2	3		Ì,			- Maria		DSC07286	DSC07287	DSC07288	DSC07289	frame_0013	frame_0014
, persua	лон	JCD4	627VB	, cazano	Jacina	Joches		4 1		12	22212605963_2a ff30a234_b	22222252864_24 7f671725_b	22366401225_57 3bc58aff_b	22366551336_00 b0438870_b	22428502998_71 442e9fbe_b	22472465975_37 ed2eecd8_b					frame_0019	frame_0020
_CECH65	JICEN	Jacoba III	KERN JOCINE	_cocine	лан	.goc.me		15	15	12	22702128072 (2	22712741777 0h	22710525425 26		22721125127.4	22752268602 42	DSC07294	DSC07295	DSC07296	DSC07297	frame 0025	frame 0026
,oran	JULINO		SCIANA JOCKET	, sense	,585383	9000	20_(1)	20, 07 20, 09	20, (1)	20, (11)	a9ce584d_b	4325e583_b	6dd0254f_b	24b19673_b	a4dc362d_b	20e940c9_b	DSC07304	DSC07305	DSC07306	DSC07307		
No.	1 C			a and	A.	Lenger L		i i i i			22818537423_9b	22822658252_20	22822952888_69	22826643732_ea	22828600726_9e	22842487410_f7					frame 0031	frame 0032



## Phase 1 : Paired Supervision (DPED+)

- Baseline Method: DPED
- Enhanced Implementation: Original feature maps scanned by different Gaussian Kernels



# Phase 1 : Paired Supervision (DPED+)

- Baseline Method: DPED
- Enhanced Implementation: Revised MobileNet Structure for better convergence





(c)  $1 \times 1$  Convolutional Filters called Pointwise Convolution in the context of Depthwise Separable Convolution

	load_dataset.py ×	train_model.py	test_model.py		models.py					
1 <b>impo</b> 2	rt tensorflow as tf									
3 def	lrelu(x, alpha):									
4	alpha = 0.005									
5 6	return tf.nn.relu(x) - alph	a * tf.nn.relu(-x)								
7 <b>def</b> 8	<pre>resnet(input_image):</pre>									
9 <sup>.</sup> 0	with tf.variable_scope("gen	erator"):								
1 2 2	W1 = weight_variable([9, 9, 3, 64], name="W1"); b1 = bias_variable([64], name="b1"); c1 = lrelu(conv2d(input_image, W1) + b1)									
.4										
.6 7 8	W2 = weight_variable([3 c2 = lrelu(_instance_no	, 3, 64, 64], name="W2"); b2 rm(conv2d(c1, W2) + b2))	<pre>= bias_variable([64],</pre>	, name="b2");						
9 0 1	W3 = weight_variable([3 c3 = lrelu(_instance_no	, 3, 64, 64], name="W3"); b3 rm(conv2d(c2, W3) + b3)) + c2	= bias_variable([64], 1	<b>, name="b</b> 3");						
2										
4 .5	W4 = weight_variable([3 c4 = lrelu(_instance_no	, 3, 64, 64], name="W4"); b4 rm(conv2d(c3, W4) + b4))	= bias_variable([64],	, name="b4");						

- Experimental Observations:
- ▶ 0.4dB loss (27.4132 vs 27.8025) of PSNR
- Accelerated Training: 2855 iter/hr (>1964 iter/hr)
- Accelerated Inference: 2.7 × speed



- Input: Generated Image I from DPED+, Reference Image R
- Weakly supervised manner, night-scene style tranfer network



• Passing both I & R into VGG-19 network and apply Gram Matrix to calculate the style loss

$$\mathcal{L}_s(O,\Lambda) = \sum_{r,l} \frac{\Lambda_r^l}{2N^{l^2}} \|G_r^l[\mathcal{O}] - G_r^l[\mathcal{R}]\|_2^2$$





• Semantically Adaptive Enhancement factor: special weights calculated by enhancement importance:  $e^{\alpha_r^l H_r^l[\mathcal{O},\mathcal{R}]}$ 

$$\Lambda_r^l = \frac{\mathrm{e}^{-rH_r}}{\sum_r \mathrm{e}^{\alpha_r^l H_r^l}[\mathcal{O},\mathcal{R}]}$$

• Photorealistic style transfer: A weighted combination of texture loss and content loss

$$\mathcal{L}_{c}(O) = \sum_{l} \alpha^{l} \frac{1}{2N^{l}D^{l}} \|F^{l}[\mathcal{O}] - F^{l}[\mathcal{I}]\|_{2}^{2}$$
$$\mathcal{L}_{m}(O) = \sum_{c \in r, g, b} V_{c}[\mathcal{O}]^{T} \mathcal{M}_{\mathcal{I}} V_{c}[\mathcal{O}]$$

313	# Content Loss	
	loss_content = content_loss(content_layer_const, content_layer_var, float(args.content_weight))	-
	tf.summary.scalar('content loss', loss_content)	5
	layer_class_style_weights = [[float(args.r1), 1.0, 5.0, 1.0, 1.0, 1.0, 1.0, float(args.s1), 1.0], [float(args.r2), 1.0, float(args.b2), 1.0, 1.0, 1	• 1
	[float(args.r3), 1.0, float(args.b3), 1.0, 1.0, 1.0, 1.0, float(args.s3), 1.0], [float(args.r4), 1.0, float(args.b4),	1
	[float(args.r5), 1.0, float(args.b5), 1.0, 1.0, 1.0, 1.0, float(args.s5), 1.0]]	

Final Objective Function and Optimization Our final objective function is a weighted combination of semantic adaptive style loss  $\mathcal{L}_s$ , content loss  $\mathcal{L}_c$  and photorealism regularization loss  $\mathcal{L}_m$ :

$$\mathcal{L}(\mathcal{O},\Lambda) = \mathcal{L}_c(\mathcal{O}) + \beta \mathcal{L}_s(\mathcal{O},\Lambda) + \gamma \mathcal{L}_m(\mathcal{O})$$

where  $\beta$  is the weight of the total style loss,  $\gamma$  is the weight to balance the photorealism regularization loss  $L_m$ .

#### # Affine Loss

if Matting:

loss\_affine = affine\_loss(input\_image\_plus, M, args.affine\_weight)
else:

loss\_affine = tf.constant(0.00001) # junk value

#### # Total Variational Loss

loss\_tv = total\_variation\_loss(input\_image, float(args.tv\_weight))

• Semantic Enhancement Network Architecture:



Figure 7. The network architecture of jointly enhancement factor  $\Lambda$  and the output image. The right side illustrates the network structure for style loss computing. We adaptively adjust our enhancement factor in regional basis and take into account the proportional information iteratively. The left side is the Gram matrix distance calculation from each region and each layer.

## Phase 3 : Optimization & Post-processing

• Recommendation Network:

FC-8 layer from the pre-trained VGG-16 network

$$D_{sem} = \|f_{fc8}(I_1 - f_{fc8}(I_1))\|_2$$



• Anchor x : minimize positive reference distances, maximize negative reference distances

$$\mathcal{L}_{ref} = \arg\min_{\mathcal{L}} \sum_{j=1}^{10} \|f_j(x_i) - f_j(x_i^{p_j})\|_2 - \sum_{k=1}^{50} \|f_j(x_i) - f_j(x_i^{n_k})\|_2 + \alpha]_{\dagger}$$

## Phase 3 : Optimization & Post-processing

• Post-Processing:

Histogram Equalization: maximize the information entropy by unify image histogram

$$T(k) = floor(L-1) \sum_{n=0}^{k} p_n$$
$$\frac{d}{dy} \left( \int_0^y p_Y(z) dz \right) = p_x(T^{-1}(y)) \frac{d}{dy}(T^{-1}(y))$$
$$\frac{d}{dx} \Big|_{x=T^{-1}(y)} \frac{d}{dy}(T^{-1}(y)) = 1$$

• Gamma Correction:

Human eyes-to-illuminance logarithm

Curve

1/2.2 gamma correction



Figure 9. Comparison between responses to different light intensity of human eye (Gamma 1/2.2) and commonly-used cameras (Linear).

## Experiment : Statistics about Training Details

• Experimental Details

Quantitative Measurement $(2 \times 10^4 \text{ and } 7 \times 10^4 \text{ iters})$									
Revised Models	PSNR	SSIM	PSNR	SSIM					
Baseline	26.3455	0.9112	27.8025	0.9338					
VGG Loss(VL)	27.4288	0.9493	28.8042	0.9703					
l1 + l2	26.2453	0.8954	26.8251	0.9122					
MobileNet	25.2410	0.9321	27.4123	0.9416					
LReLU(LR)	26.5448	0.9107	27.7239	0.9293					
Contrast(HE)	27.0145	0.9220	28.0112	0.9445					
VL+LR+HE	27.4338	0.9420	29.2002	0.9754					



#### self.input = x = tf.placeholder(tf.float32,[None,img\_size,img\_size,output\_channels]) self.target = y = tf.placeholder(tf.float32,[None,img\_size\*scale,img\_size\*scale,output\_channels]) Preprocessing as mentioned in the paper, by subtracting the mean However, the subtract the mean of the entire dataset they use. As of now, I am subtracting the mean of each batch mean\_x = tf.reduce\_mean(self.input) image\_input =x- mean\_x mean\_y = tf.reduce\_mean(self.target) image\_target =y- mean\_y x = slim.conv2d(image\_input,feature\_size,[3,3]) scaling\_factor = 0.1 for i in range(num layers): x = utils.resBlock(x,feature\_size,scale=scaling\_factor) x = slim.conv2d(x,feature\_size,[3,3]) $x += conv_1$ self.out = tf.clip\_by\_value(output+mean\_x,0.0,255.0)

self.loss = loss = tf.reduce\_mean(tf.losses.absolute\_difference(image\_target,output))

## Entire Pipeline: 3 Phases of MPSE



- Phase 1: PR Learning DPED+ [detail & texture augmentation]
- Phase 2: UPR Learning [semantic style transfer, structural enhancement]
- Phase 3: Post-processing [global contrast, illuminance map adjustment]

## MPSE Entire Pipeline: Result





### LIME

• The task of low-light enhancement can be viewed as 'finding a better illuminance map'

 $\mathbf{L} = \mathbf{I} \circ \mathbf{T},$ 

- L is the original input and I is the desired recovery, T is the illuminance map
- Given the same original image, different estimated illuminance map results diverting enhancement results
- The wanted illuminance map, smooth, local-consistent and texture-free
- The wanted illuminance map should exploit the maximum potential of the original image



### Intrinsic Image Decomposition

• Decompose the Image into Albedo (reflectance) and Shading

$$\mathcal{I} = \mathcal{A} \cdot \mathcal{S}$$

- Shading gives important cues on the object geometry and material
- Albedo(reflectance) is a lighting invariant quantity
- Albedo is similar to the previous illuminance map







Albedo

Shading

• Lighting-texture-structure separation can be possible (supervision)



### Observation on Real-World Intrinsic

- In daylight condition, the ambient light is relatively abundant: DSLR-Mobile is mainly **Shading-different**
- In night-scene, with poor lighting condition:
   DSLR-Mobile is mainly Illuminance-different (minor Shading-different)
- Structure information mainly exists in **Reflectance**
- Texture information mainly exists in **Shading**
- Noise pattern mainly exists in **Shading**



(b)

(c)

### Self-Supervised Intrinsic Decomposition

#### Problem and Motivation:

- Real-world intrinsic dataset is rare, manually labelling is extremely expensive
- CG-based dataset (e.g. MPI Sintel) performs disappointingly on real-world task

#### Solutions:

- Symmetric network structure, two encoder-decoders
- Reconstruction loss, similar to cycle loss

$$\|\mathcal{I} - \mathcal{A} \cdot \mathcal{S}\|_{2}^{2}$$
  $\mathcal{A} = \mathcal{I}, \mathcal{S} = 11^{T}$   $\mathcal{L}_{rec} = \|\mathcal{I} - \mathcal{F}'(\mathcal{A}, \mathcal{S})\|^{2}$ 

- Pretrain the first encoder-decoder on synthesized dataset
- Train the entire network with reconstruction loss on real-world unlabeled dataset



### Multi-Domains Translation: Data-Driven Initiatives

#### Absence of Overall Enhancement Dataset:

- Alignment Issue: Overall DSLR-Mobile Dataset is still absent
- Partial enhancement datasets largely exist (denoise, color, contrast, texture, resolution)
- Denoise (DND), Color and Contrast (Adobe 5K), Resolution (DIV 2K)

#### Unclear Training Directions:

- The actual contributing factors behind 'high-quality' is hidden
- In overall enhancement training, noise and texture is sometimes indistinguishable



#### Solution:

- Three datasets dedicated for denoise, color & contrast, texture
- Multi-domain image translation GAN, even if the overall enhancement dataset is lacked

### Star-GAN: Unified Solution for Multi-Domains

#### Previous Image-to-Image Translation Network:

- General solutions require training N(N-1)/2 dual-domain image translation network
- Previous solution cannot translate to a domain with no paired data

#### Star-GAN

- Only one universal generator is needed
- Attribute: enhancement factor inherent in an image such as color, texture, or noise level.
- Attribute Value: a particular level of the attribute, such as low-quality or high-quality.
- Domain: a set of images sharing the same attribute value. E.g., images of low-quality in color perspective.

#### Our Goal:

• Train G to translate an input image x into an output image y conditioned on the target domain label c

$$G(x,c) \to y.$$

• Train D to produce probability distributions over both sources and domain labels

$$D: x \to \{D_{src}(x), D_{cls}(x)\}.$$



### Network Design

#### Masked Vector:

- We use an n-dimensional **one-hot vector** to denotes attribute value vector **m**
- If here we input the i-th dataset, then only  $c_i$  is assigned to 1 and others is assigned to 0
- We have three datasets dedicated to color, texture, and noise. Then the n is set to be 3.
- In each dataset, we only have two attribute value: low-quality(0) and high quality(1).

$$\widetilde{c} = [c_1, ., c_i, ., c_n, m]$$

#### Domain Classification Loss:

• An auxiliary classifier is added on top of the discriminator D to impose domain classification

$$\mathcal{L}_{cls}^{r} = \mathbb{E}_{x,c'}[-log D_{cls}(c'|x)]$$

- $D_{cls}(c'|x)$  denotes a probability distribution over domain labels calculated by D.
- D learns to classify a real image x to its corresponding original domain c'



### Network Design

#### Cycle Construction Loss:

• G takes the translated image G(x, c) and the domain label of the original input c' as input and attempts to reconstruct the original image x.

$$\mathcal{L}_{cyc} = \mathbb{E}_{x,c,c'} [\|x - G(G(x,c),c')\|]_1$$

#### **Total Objective Function**

•  $\lambda_{cls}$  and  $\lambda_{cyc}$  are hyper-parameters that balance the relative weight of domain classification loss and cycle re- construction loss.

$$\mathcal{L}_D = -\mathcal{L}_{adv} + \lambda_{cls} \mathcal{L}_{cls}^r$$

$$\mathcal{L}_G = \mathcal{L}_{adv} + \lambda_{cls} \mathcal{L}_{cls}^f + \lambda_{cyc} \mathcal{L}_{cyc}$$



## Training Strategy

#### Our Previous Observations:

- Color and contrast information exist both in **Reflectance** and **Shading**
- Texture information mainly exists in **Shading**
- Noise pattern mainly exists in Shading

#### Three Partially Enhanced Datasets:

- DIV2K: resolution and texture enhancement (train on **Shading**)
- Adobe-5K: color and contrast enhancement (train on **Reflectance** and **Shading**)
- DND Denoise Benchmark Dataset: noise reduction (train on Shading)

#### Inference Setting:

• We Set all c<sub>i</sub> to 1 with the m value of 1, indicating the translation from one LQ domain with all 'bad factors' to the HQ domain with all 'good' factors.

$$\widetilde{c} = [c_1, ., c_i, ., c_n, m]$$

Total Objective:

$$\mathcal{L}_D = -\mathcal{L}_{adv} + \lambda_{cls} \mathcal{L}_{cls}^r$$
$$\mathcal{L}_G = \mathcal{L}_{adv} + \lambda_{cls} \mathcal{L}_{cls}^f + \lambda_{cyc} \mathcal{L}_{cyc}$$
$$\mathcal{L}_{rec} = \alpha \|\mathcal{I}_{LQ} - \mathcal{F}'(\mathcal{A}_{LQ}, \mathcal{S}_{LQ})\|^2 + \beta \|\mathcal{I}_{HQ} - \mathcal{F}'(\mathcal{A}_{HQ}, \mathcal{S}_{HQ})\|^2$$

## MDIE: Decompose-Enhance-Reconstruct



Multi-Domain Translation

## MDIE Experiment





### Future Works - MPSE

• Augmented Semantic Regions: current three main regions - building, river, sky



• Semantic Enhancement Factor: Weight Selection

Enhancement Importance: Softmax, Texture Calculation, Total Variation





## Future Works - MDIE

#### Texture Enrichment:

• Not limited to resolution, but using texture dataset for training

Joint Training of the Intrinsic and Translation Network:

- Pretrain the first encoder-decoder using MPI Sintel
- Input three partially enhanced datasets, joint training and balancing the two networks
- The intrinsic network should learn real-life decomposition
- The translation network should learn the partial enhancement for a overall enhancement

$$\mathcal{L}_{D} = -\mathcal{L}_{adv} + \lambda_{cls} \mathcal{L}_{cls}^{r}$$
$$\mathcal{L}_{G} = \mathcal{L}_{adv} + \lambda_{cls} \mathcal{L}_{cls}^{f} + \lambda_{cyc} \mathcal{L}_{cyc}$$
$$\mathcal{L}_{rec} = \alpha \|\mathcal{I}_{LQ} - \mathcal{F}'(\mathcal{A}_{LQ}, \mathcal{S}_{LQ})\|^{2} + \beta \|\mathcal{I}_{HQ} - \mathcal{F}'(\mathcal{A}_{HQ}, \mathcal{S}_{HQ})\|^{2}$$

#### Minimizing the Separation Loss:

- Using more sophisticated network for intrinsic decomposition
- Purify the separation of albedo and shading in real world data

# Improve Mobile Phone's Image Quality by Deep Learning

Introduction & Review ·
Augmented Dataset ·
MPSE(Multi-Phase Semantic Enhancement) ·
MDIE(Multi-Domain Intrinsic Enhancement) ·
Future Work·

 $\cdot Q\&A \cdot$ 

HUANG Cen | LIU Yicun