

Exposure-slot: Exposure-centric representations learning with Slot-in-Slot Attention for Region-aware Exposure Correction

-Supplementary Material-

Anonymous CVPR submission

Paper ID 10942

001 A. Encoder and Decoders

002 We use an Image Encoder (**Enc**), an Image Decoder (**Dec**_{enhance}) , and a Slot Decoder (**Dec**_{slot}) for slot
003 reconstruction as the backbone network of Exposure-slot.
004

005 Table S1 and S2 presents the detailed architecture of the
006 **Enc**, **Dec**_{enhance} and **Dec**_{slot}. The Conv-block includes
007 a convolution operation with a stride of 1 and padding of 1,
008 followed by the GeLU [4] activation function.

Stage	Operations	Outputs
Enc-1	Conv-block, 3 × 3	$h \times w \times 32$
	Conv-block, 3 × 3	$h \times w \times 32$
	batchnorm2d(32)	$h \times w \times 32$
Enc-2	Conv-block, 3 × 3	$h \times w \times 16$
	PixelShuffle(2)	$h/2 \times w/2 \times 64$
	Conv-block, 3 × 3	$h/2 \times w/2 \times 64$
	batchnorm2d(64)	$h/2 \times w/2 \times 64$
Enc-3	Conv-block, 3 × 3	$h/2 \times w/2 \times 32$
	PixelShuffle(2)	$h/4 \times w/4 \times 128$
	Conv-block, 3 × 3	$h/4 \times w/4 \times 128$
	batchnorm2d(128)	$h/4 \times w/4 \times 128$
Dec _{enhance} -1	Conv-block, 3 × 3	$h/4 \times w/4 \times 128$
	PixelUnshuffle(2)	$h/2 \times w/2 \times 64$
-	Skip-connection with Enc-2	$h/2 \times w/2 \times 128$
Dec _{enhance} -2	Conv-block, 3 × 3	$h/2 \times w/2 \times 128$
	PixelUnshuffle(2)	$h \times w \times 32$
	Conv-block, 3 × 3	$h \times w \times 32$
	Conv-block, 3 × 3	$h \times w \times 32$
-	Skip-connection with Enc-1	$h \times w \times 64$
Dec _{enhance} -3	Conv-block, 3 × 3	$h \times w \times 32$
	Conv-block, 3 × 3	$h \times w \times 32$
	Conv-block, 1 × 1	$h \times w \times 3$

Table S1. Specification of Image Encoder (**Enc**) and Decoder (**Dec**_{enhance}) architecture.

Stage	Operations	Outputs
Dec _{slot} -1	PixelUnshuffle(2)	$h/2 \times w/2 \times 64$
	Conv-block, 3 × 3	$h/2 \times w/2 \times 64$
	Conv-block, 3 × 3	$h/2 \times w/2 \times 64$
Dec _{slot} -2	PixelUnshuffle(2)	$h \times w \times 32$
	Conv-block, 3 × 3	$h \times w \times 32$
	Conv-block, 3 × 3	$h \times w \times 16$
Dec _{slot} -3	Conv-block, 3 × 3	$h \times w \times 3$

Table S2. Specification of Slot Decoder (**Dec**_{slot}) architecture.

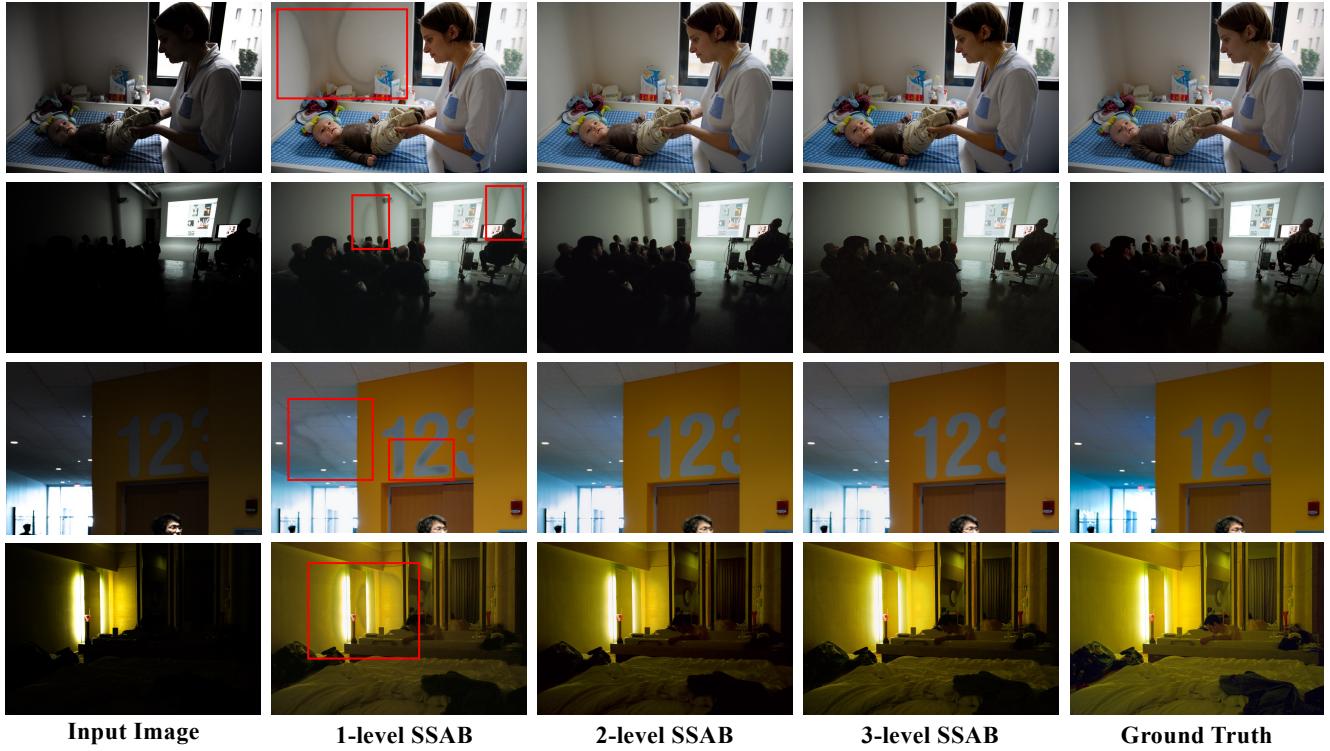
009 B. More Results of Each Structural Levels

010 In the main manuscript, we set the default configuration of
011 the SSAB block to 2-level. Additionally, in Sec.4.3 of the
012 main manuscript, Table 5, we validated the effectiveness of
013 different structural levels on the SICE [3] dataset. Furthermore,
014 in Table S3, we provide additional results for 1 and
015 3-level SSAB on the MSEC [1] and LCDPNet [8] datasets.

Dataset	Model	K ^{main}	K ^{sub-1}	K ^{sub-2}	PSNR↑	SSIM↑
SICE [3]	1-level	3	-	-	22.02	0.7131
	2-level	3	7	-	<u>22.81</u>	<u>0.7236</u>
	3-level	3	7	10	<u>23.06</u>	<u>0.7306</u>
MSEC [1]	1-level	3	-	-	23.04	0.8668
	2-level	3	7	-	<u>23.18</u>	<u>0.8697</u>
	3-level	3	7	10	<u>23.25</u>	<u>0.8700</u>
LCDP [8]	1-level	3	-	-	23.81	<u>0.8596</u>
	2-level	3	7	-	<u>24.03</u>	0.8592
	3-level	3	7	10	<u>24.13</u>	<u>0.8629</u>

Table S3. Extended version of Table 5 in the main manuscript. Ablation studies on not only SICE [3] dataset but also MSEC [1] and LCDP [8] dataset is additionally provided.

016 As the level increases, the performance metric values
017 for PSNR and SSIM improve across all benchmark

Figure S1. Visual comparisons across different n -level SSAB structures. ($n = 1, 2, 3$)

018 datasets [1, 3, 8]. Notably, the 3-level SSAB achieves a 0.07
 019 performance gain in SSIM on the SICE dataset compared to
 020 the 2-level SSAB. Depending on computer resources, users
 021 can set SSAB levels beyond 2-levels to achieve better
 022 results.

Model	Params (M)	FLOPs (G)	Time (S)
1-level	1.079	14.064	0.0644
2-level	1.229	14.175	0.0676
3-level	1.465	14.618	0.0938

Table S4. Computational cost across different levels of SSAB.

023 In Table S4, we present the computational cost of differ-
 024 ent levels of SSAB in terms of parameters, FLOPs (Floating
 025 point operations per second), and execution time. The Flops
 026 is calculated on $3 \times 256 \times 256$ input, and execution time
 027 measurements taken on a single 844×1500 RGB image us-
 028 ing an NVIDIA RTX 4090 GPU. The 1-level and the 2-level
 029 SSAB, which we used as the default in the main manuscript,
 030 show minimal differences in computational cost. However,
 031 the 3-level configuration requires slightly more runtime,
 032 with an execution time increase of approximately 0.1 sec-
 033 onds.

034 Additionally, Fig. S1 presents the visual output results

035 corresponding to each level. The 1-level SSAB tends to
 036 produce artifacts such as blotches caused by lighting (red
 037 boxes), whereas such issues are absent with configura-
 038 tions of 2 and 3-level SSAB. To aid understanding, Fig. S2 visu-
 039 alizes the slot attention maps for each level. In the 1-level
 040 SSAB, the attention maps show stark partitioning based on
 041 light sources, whereas this issue is resolved in configura-
 042 tions of 2-level SSAB or higher.

C. More Quantitative results on Perceptual Quality

043 In Table S5, we also provide a perceptual comparison of
 044 the results with other methods. The evaluation is con-
 045 ducted on SICE [3] dataset. To measure the perceptual
 046 quality, we adopt Learned Perceptual Image Patch Simi-
 047 larity (LPIPS) [9] and Perception Index (PI) [2].

Model	#Params	Time (S)	LPIPS \downarrow	PI \downarrow
ECLNet [6]	0.018	0.1328	0.268	3.346
FECNet [5]	0.150	<u>0.0746</u>	0.298	3.679
CSEC [7]	1.364	2.3633	<u>0.208</u>	<u>2.993</u>
Exposure-slot	1.229	0.0676	0.161	2.949

Table S5. Results of Perceptual Quality

049

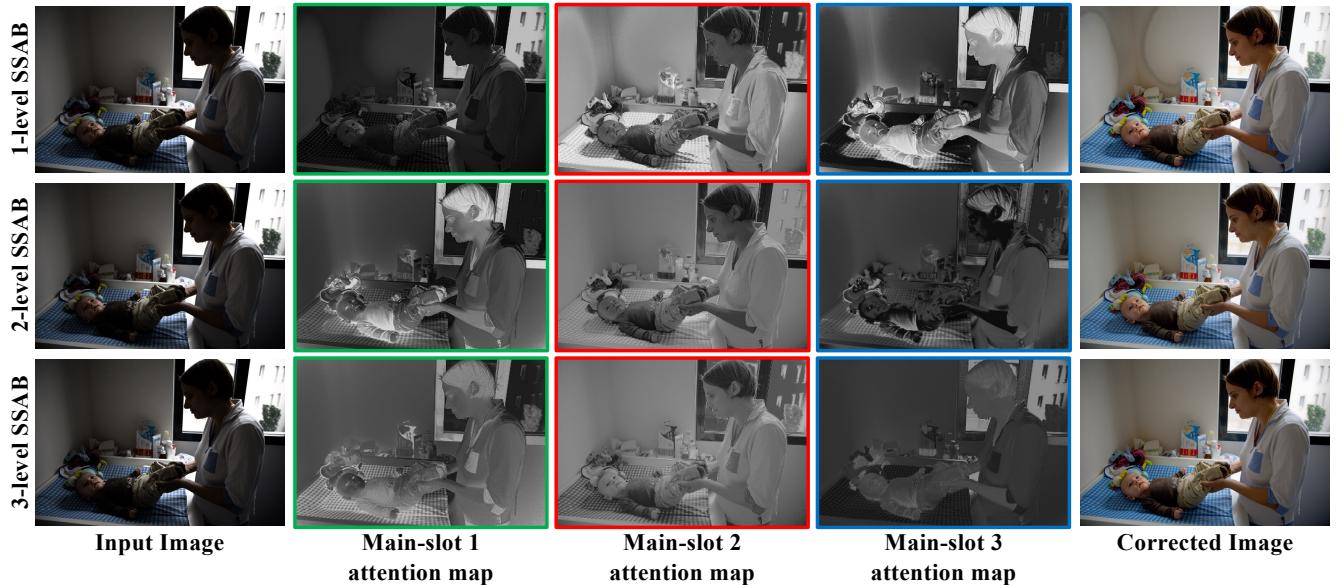


Figure S2. Visual comparisons of slot attention maps across different n -level SSAB structures ($n = 1, 2, 3$).

050 The execution time measurements taken on a single
 051 844×1500 RGB image using an NVIDIA RTX 4090 GPU.
 052 Exposure-slot delivers faster processing speeds compared
 053 to existing approaches while maintaining high perceptual
 054 quality.

055 D. More Qualitative results.

056 Fig. S3, S4 and S5 present additional visual results on the
 057 LCDP [8] dataset. For comparison, we include LCDP-
 058 Net [8] and CSEC [7] as baseline methods. LCDPNet
 059 shows more robust results against light source diffusion
 060 compared to CSEC but struggles with accurate color cor-
 061 rection, whereas CSEC excels at color correction but gen-
 062 erates artifacts due to light source diffusion. In contrast,
 063 our proposed method, Exposure-slot achieves robust cor-
 064 rection results, effectively addressing both challenges.

065 E. Visualization of Slot Attention Maps

066 We visualize the progressive refinement of each slot atten-
 067 tion map. In Fig. S6 – S18, we provide attention maps of
 068 main- and sub-slots at each iteration. As shown in figures,
 069 attention maps progressively evolve into features optimized
 070 for exposure correction, reflecting gradual feature improve-
 071 ment and smoothing with each iteration.

072 F. Visualization of t-SNE.

073 This section provides details regarding the t-SNE plots
 074 shown in Fig. 5 of the main manuscript. In Fig. 5, t-SNE
 075 is visualized based on the feature vectors from \mathbf{V} in Eq. 2.

076 The left side of the figure illustrates the t-SNE plot of \mathbf{V}
 077 before applying prompts, while the right side displays the
 078 t-SNE plot of $\mathbf{V} \cdot \mathbf{P}^{final}$, representing the prompt applied
 079 version.

080 For further information, in this supplementary material,
 081 we provide t-SNE plots for the 1, 2, and 3-level SSAB struc-
 082 ture in Fig. S19, S20, and S21. For each figure, the plots
 083 from left to right represent features before prompts, features
 084 after prompts, prompts, and slots. Features before and af-
 085 ter prompts are visualized using the method employed in
 086 Fig. 5, while \mathbf{S}^{final} and \mathbf{P}^{final} from Eq. 9 are used for slots
 087 and prompts, respectively. Additionally, from top to bot-
 088 tom, each row corresponds to plots clustered at the main-
 089 slot, sub-slot, and sub2-slot levels. We observe that not only
 090 the prompts and slots exhibit distinct distributions and rela-
 091 tionships, but the features after prompts also display unique
 092 patterns for each level of the SSAB structure. This indicates
 093 that as the level increases, the SSAB structure becomes pro-
 094 gressively more adept at performing sophisticated exposure
 095 correction.

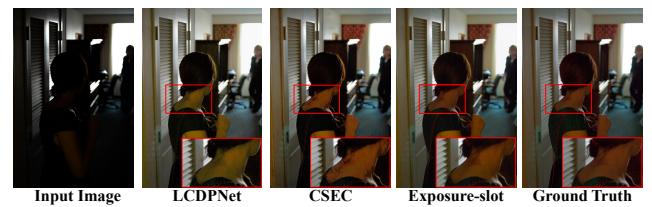


Figure S3. Visual results from LCDP [8] dataset. We utilize LCDPNet [8] and CSEC [7] as comparison methods.

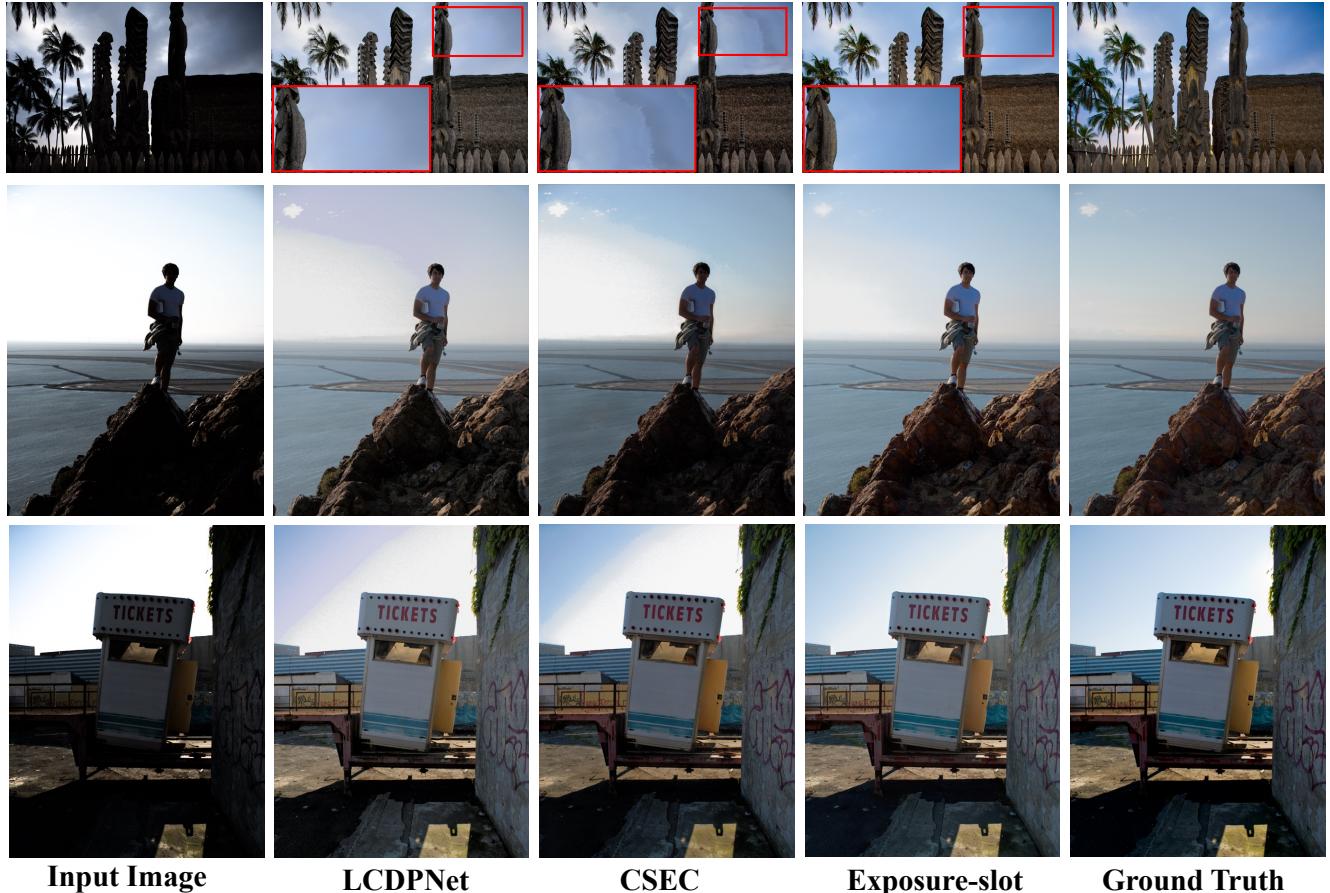


Figure S4. Visual results from LCDP [8] dataset. We utilize LCDPNet [8] and CSEC [7] as comparison methods.



Figure S5. Visual results from LCDP [8] dataset. We utilize LCDPNet [8] and CSEC [7] as comparison methods.

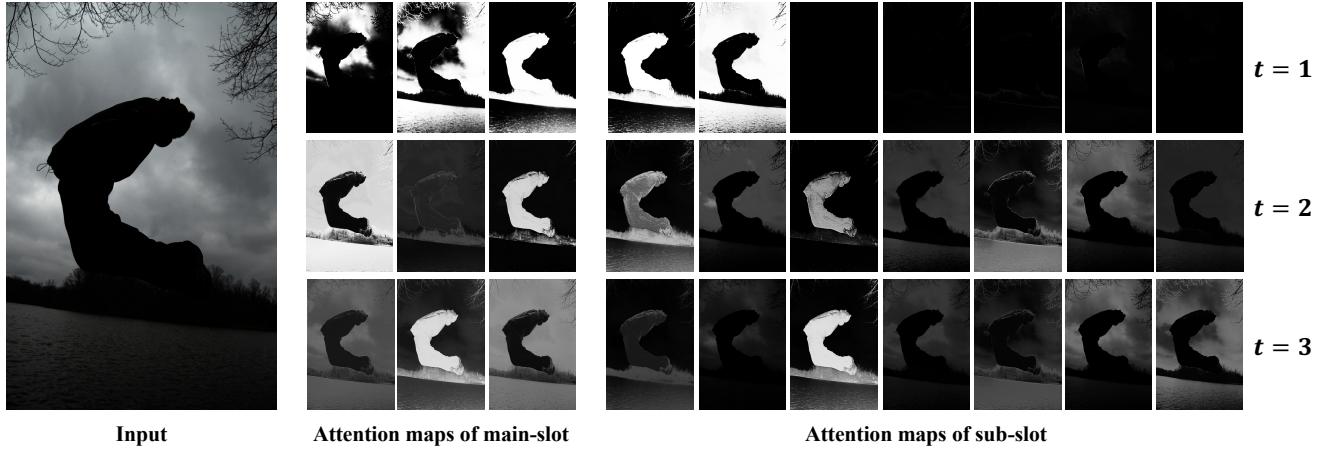


Figure S6. Visual result of attention maps from LCDP [8] dataset. From top to bottom, we present the attention maps of the main-slot and sub-slot at each iteration ($t = 1, 2, 3$).

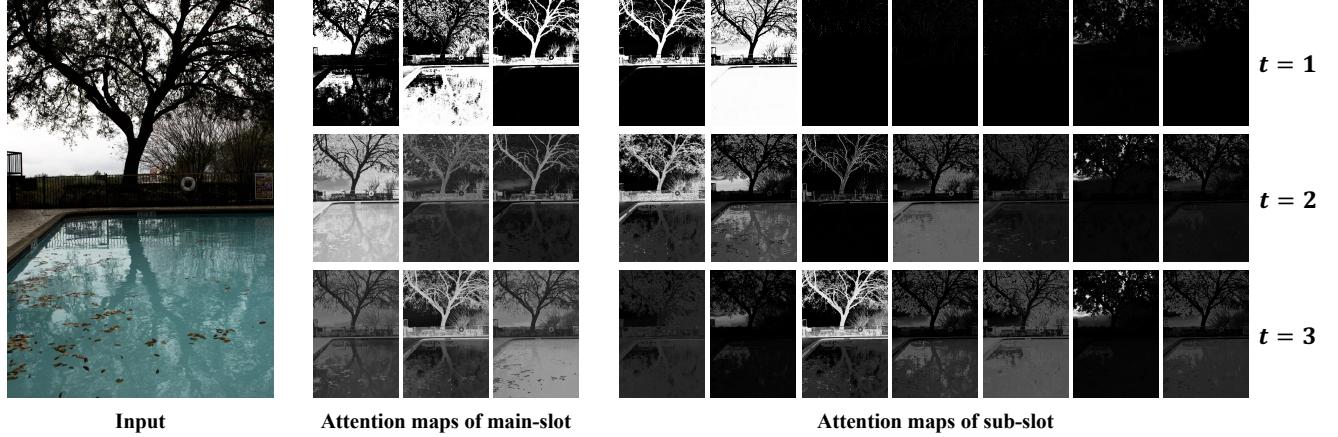


Figure S7. Visual result of attention maps from LCDP [8] dataset. From top to bottom, we present the attention maps of the main-slot and sub-slot at each iteration ($t = 1, 2, 3$).

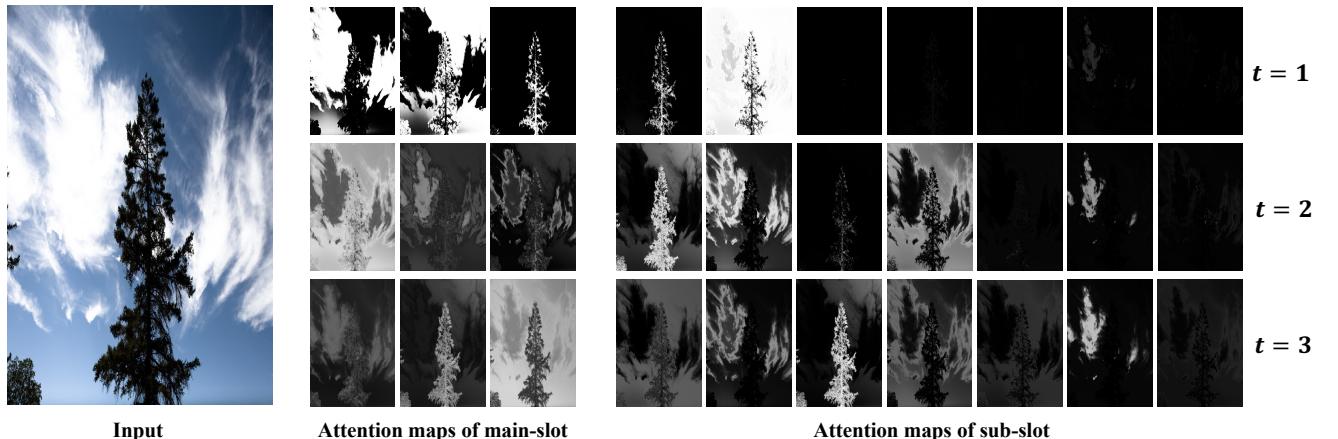


Figure S8. Visual result of attention maps from LCDP [8] dataset. From top to bottom, we present the attention maps of the main-slot and sub-slot at each iteration ($t = 1, 2, 3$).

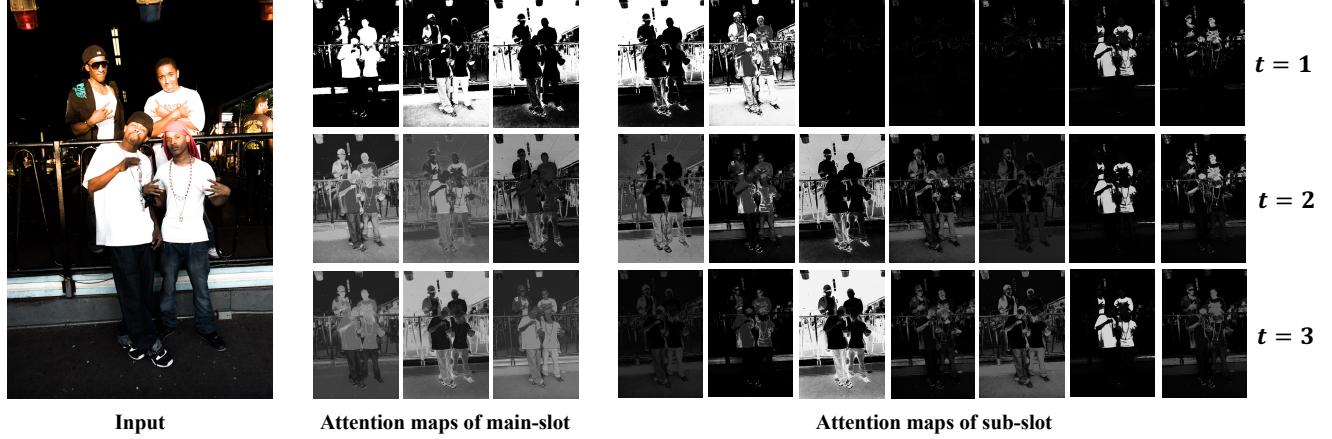


Figure S9. Visual result of attention maps from LCDP [8] dataset. From top to bottom, we present the attention maps of the main-slot and sub-slot at each iteration ($t = 1, 2, 3$).

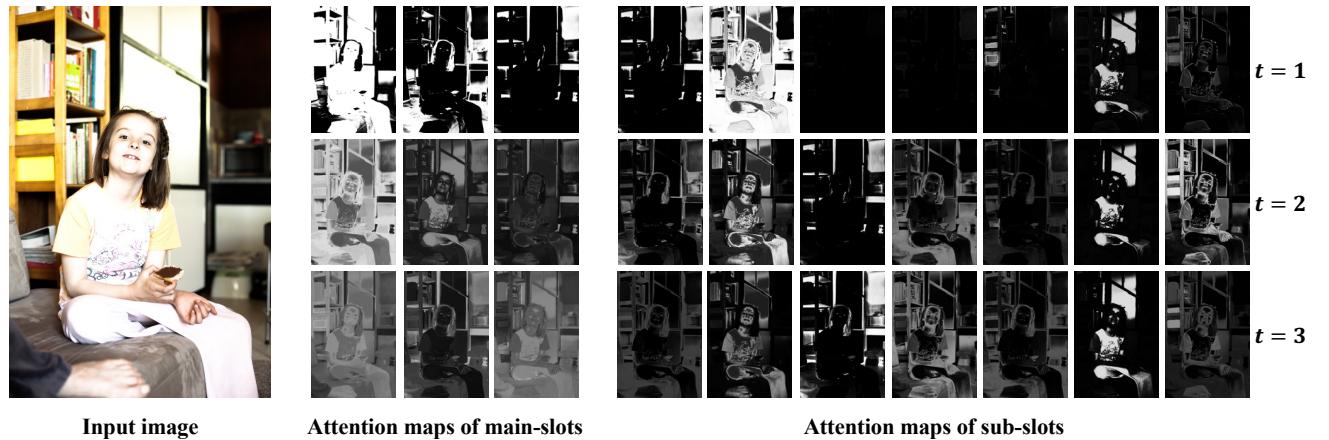


Figure S10. Visual result of attention maps from LCDP [8] dataset. From top to bottom, we present the attention maps of the main-slot and sub-slot at each iteration ($t = 1, 2, 3$).

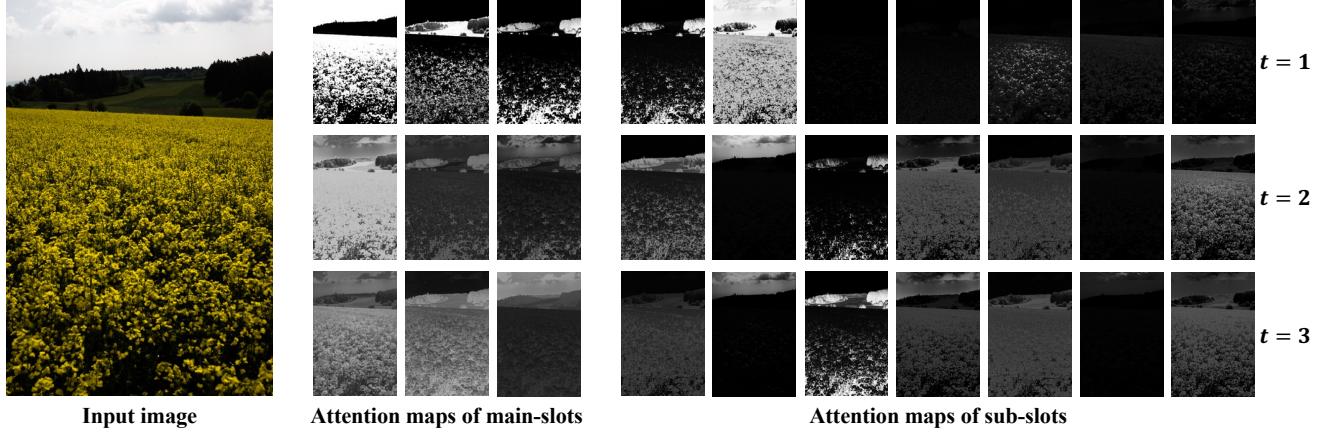


Figure S11. Visual result of attention maps from LCDP [8] dataset. From top to bottom, we present the attention maps of the main-slot and sub-slot at each iteration ($t = 1, 2, 3$).

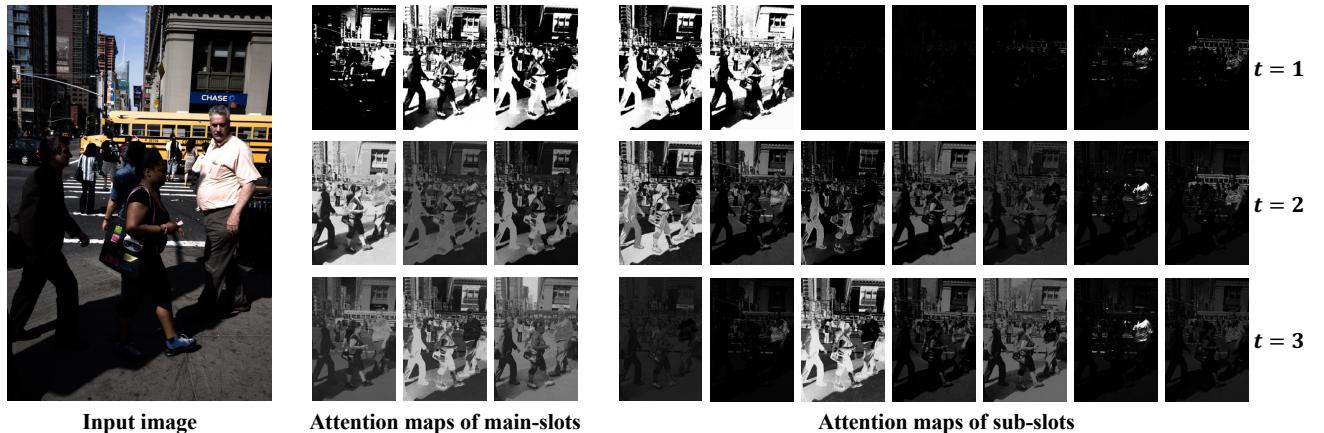


Figure S12. Visual result of attention maps from LCDP [8] dataset. From top to bottom, we present the attention maps of the main-slot and sub-slot at each iteration ($t = 1, 2, 3$).



Figure S13. Visual result of attention maps from LCDP [8] dataset. From top to bottom, we present the attention maps of the main-slot and sub-slot at each iteration ($t = 1, 2, 3$).

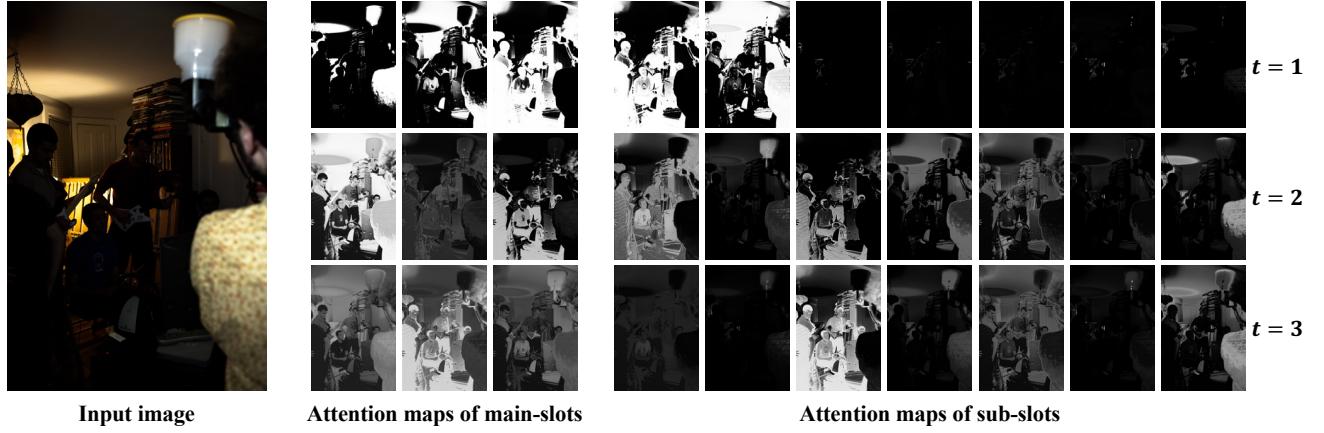


Figure S14. Visual result of attention maps from LCDP [8] dataset. From top to bottom, we present the attention maps of the main-slot and sub-slot at each iteration ($t = 1, 2, 3$).



Figure S15. Visual result of attention maps from LCDP [8] dataset. From top to bottom, we present the attention maps of the main-slot and sub-slot at each iteration ($t = 1, 2, 3$).

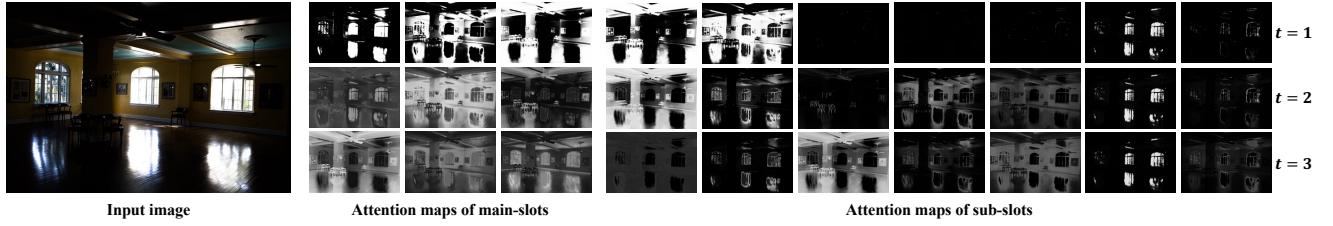


Figure S16. Visual result of attention maps from LCDP [8] dataset. From top to bottom, we present the attention maps of the main-slot and sub-slot at each iteration ($t = 1, 2, 3$).



Figure S17. Visual result of attention maps from LCDP [8] dataset. From top to bottom, we present the attention maps of the main-slot and sub-slot at each iteration ($t = 1, 2, 3$).



Figure S18. Visual result of attention maps from LCDP [8] dataset. From top to bottom, we present the attention maps of the main-slot and sub-slot at each iteration ($t = 1, 2, 3$).

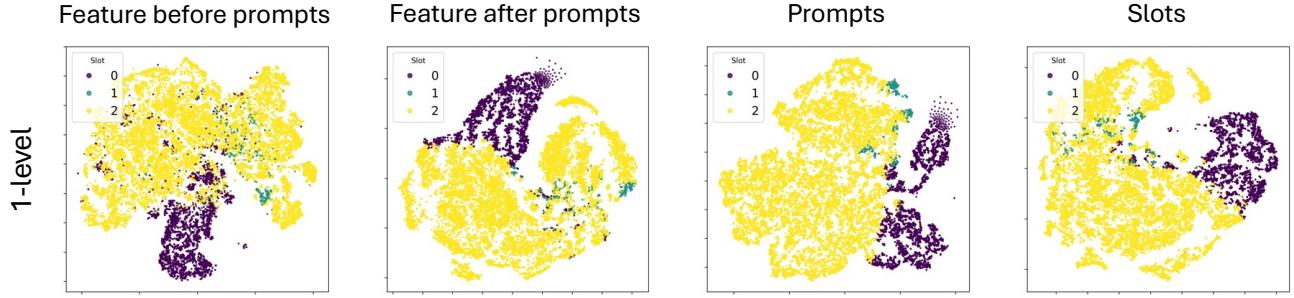


Figure S19. The t-SNE plots for the 1-level SSAB structure are presented. From left to right, the plots represent features before prompts (\mathbf{V}), features after prompts ($\mathbf{V} \cdot \mathbf{P}^{\text{final}}$) as defined in Eq. 2, prompts ($\mathbf{P}^{\text{final}}$), and slots ($\mathbf{S}^{\text{final}}$) as defined in Eq. 9, respectively.

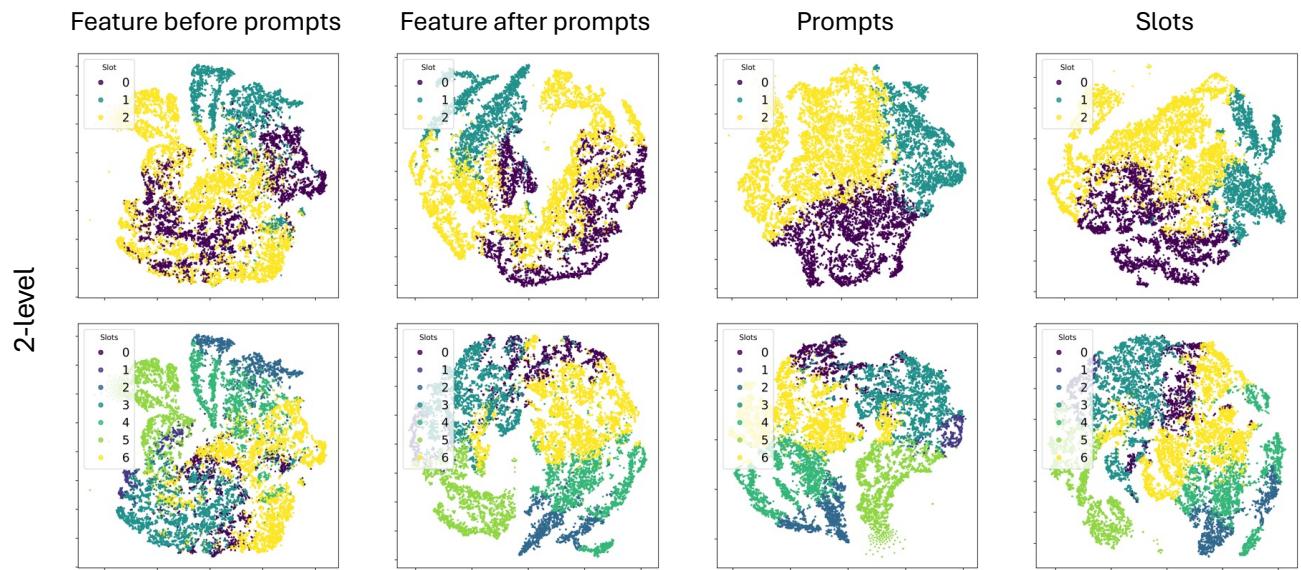


Figure S20. The t-SNE plots for the 2-level SSAB structure are presented. From left to right, the plots represent features before prompts (\mathbf{V}), features after prompts ($\mathbf{V} \cdot \mathbf{P}^{\text{final}}$) as defined in Eq. 2, prompts ($\mathbf{P}^{\text{final}}$), and slots ($\mathbf{S}^{\text{final}}$) as defined in Eq. 9, respectively. Additionally, from top to bottom, each row corresponds to plots clustered at the main-slot and sub-slot, respectively.

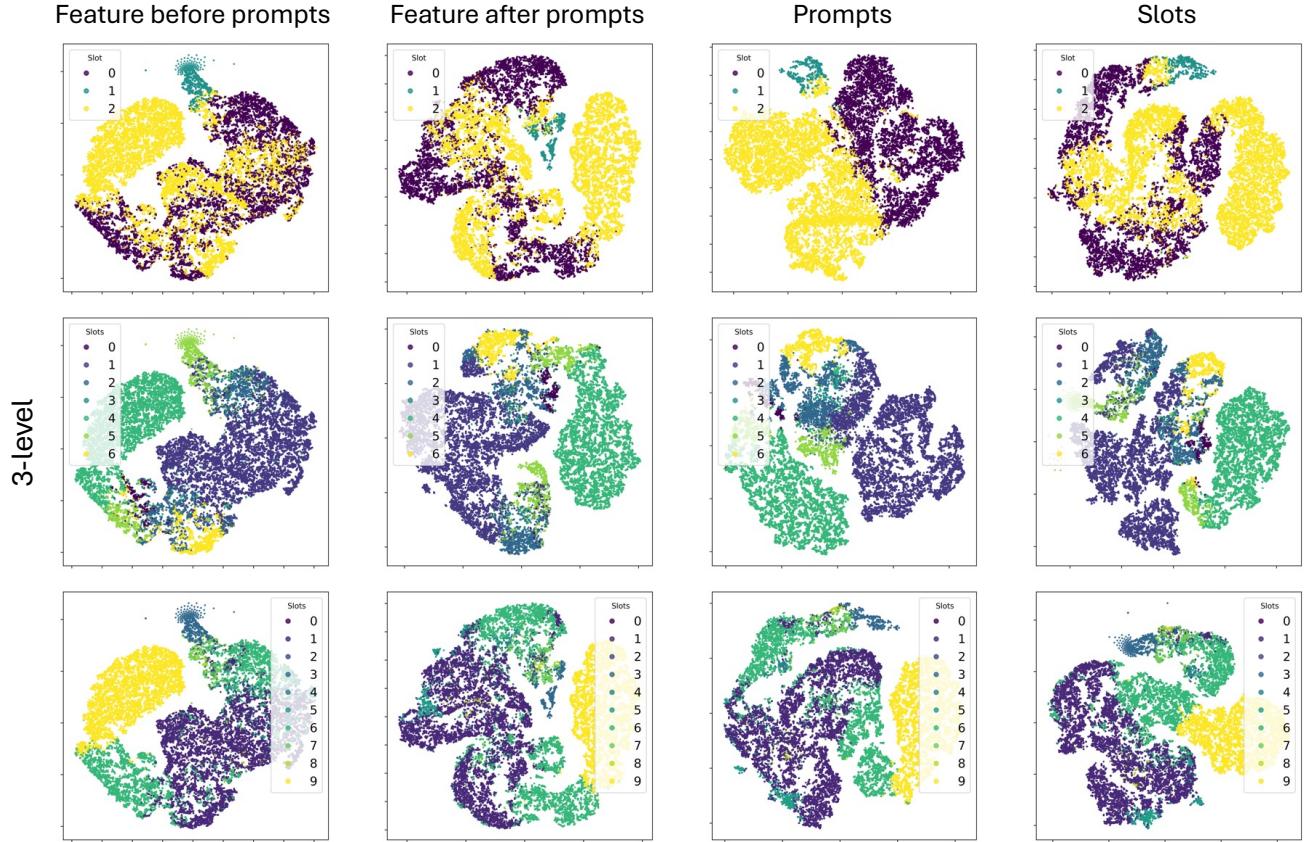


Figure S21. The t-SNE plots for the 2-level SSAB structure are presented. From left to right, the plots represent features before prompts (\mathbf{V}), features after prompts ($\mathbf{V} \cdot \mathbf{P}^{\text{final}}$) as defined in Eq. 2, prompts ($\mathbf{P}^{\text{final}}$), and slots ($\mathbf{S}^{\text{final}}$) as defined in Eq. 9, respectively. Additionally, from top to bottom, each row corresponds to plots clustered at the main-slot, sub-slot, and sub2-slot, respectively.

096 **References**

- 097 [1] Mahmoud Afifi, Konstantinos G Derpanis, Bjorn Ommer, and
098 Michael S Brown. Learning multi-scale photo exposure cor-
099 rection. In *CVPR*, 2021. [1](#), [2](#)
- 100 [2] Yochai Blau, Roey Mechrez, Radu Timofte, Tomer Michaeli,
101 and Lihi Zelnik-Manor. The 2018 pirm challenge on percep-
102 tual image super-resolution. In *ECCVW*, 2018. [2](#)
- 103 [3] Jianrui Cai, Shuhang Gu, and Lei Zhang. Learning a deep
104 single image contrast enhancer from multi-exposure images.
105 *IEEE Transactions on Image Processing*, 27, 2018. [1](#), [2](#)
- 106 [4] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units
107 (gelus). *arXiv preprint arXiv:1606.08415*, 2016. [1](#)
- 108 [5] Jie Huang, Yajing Liu, Feng Zhao, Keyu Yan, Jinghao Zhang,
109 Yukun Huang, Man Zhou, and Zhiwei Xiong. Deep fourier-
110 based exposure correction network with spatial-frequency in-
111 teraction. In *ECCV*, 2022. [2](#)
- 112 [6] Jie Huang, Man Zhou, Yajing Liu, Mingde Yao, Feng Zhao,
113 and Zhiwei Xiong. Exposure-consistency representation
114 learning for exposure correction. In *ACMMM*, 2022. [2](#)
- 115 [7] Yiyu Li, Ke Xu, Gerhard Petrus Hancke, and Rynson WH
116 Lau. Color shift estimation-and-correction for image enhance-
117 ment. In *CVPR*, 2024. [2](#), [3](#), [4](#), [5](#)
- 118 [8] Haoyuan Wang, Ke Xu, and Rynson WH Lau. Local color dis-
119 tributions prior for image enhancement. In *ECCV*. Springer,
120 2022. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [9](#), [10](#)
- 121 [9] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman,
122 and Oliver Wang. The unreasonable effectiveness of deep
123 features as a perceptual metric. In *Proceedings of the IEEE*
124 *conference on computer vision and pattern recognition*, pages
125 586–595, 2018. [2](#)