#### **Iterative Prompt Learning for Unsupervised Backlit Image Enhancement**

Zhexin Liang Chongyi Li Shangchen Zhou Ruicheng Feng Chen Change Loy S-Lab, Nanyang Technological University

{zliang008, chongyi.li, s200094, ruicheng002, ccloy}@ntu.edu.sg https://zhexinliang.github.io/CLIP\_LIT\_page/

ICCV 2023

# Background

#### Low-light images:



(a) Insufficient lighting



(b) back lighting



(c) Uneven lighting

Backlit image enhancement is more challenging as it requires **preserving well-lit areas** while **enhancing underexposed areas**.

### Background

#### Conventional methods.

Difficult to process real-world backlit images.

#### Supervised methods. Retinex-Net、KID、 URetinex-Net

Over-enhancement in well-lit areas or under-enhancement in low-light areas.

#### > Unsupervised methods. Zero-DCE, EnlightenGAN

Rely on ideal assumptions such as mean brightness and gray world models or directly learn the distribution of normal light images through adversarial training.

### Introduction



0.44 0.91

0.73

0.95

A photo of a low light scene. A photo of a backlit scene. Low light. Backlit.

A photo of a low light scene.	0.64
A photo of a backlit scene.	0.27
Low light.	0.47
Backlit.	0.33

 $s = \frac{\boldsymbol{x} \odot \boldsymbol{t}}{||\boldsymbol{x}|| \cdot ||\boldsymbol{t}||},$ 

The optimal prompts could vary on a case-by-case basis due to the complex illuminations in the scene.

It is unlikely to achieve optimal performance with **fixed prompts**.

Replace fixed prompts with **learnable prompts** and continuous **prompt refinement** helps achieve image enhancement





(b) Prompt Refinement and Enhancement Model Fine-tuning

The first stage involves the **initialization of negative and positive (learnable) prompts** to roughly characterize backlit and well-lit images, as well as the **training of the initial enhancement network**.



Antonym prompt pairing





#### The CLIP-Enhance loss *L<sub>clip</sub>*

$$\mathcal{L}_{clip} = \frac{e^{\cos(\Phi_{image}(I_t), \Phi_{text}(T_n))}}{\sum_{i \in \{n, p\}} e^{\cos(\Phi_{image}(I_t), \Phi_{text}(T_i))}}$$

#### The identity loss *L<sub>identity</sub>*

$$\mathcal{L}_{identity} = \sum_{l=0}^{4} \alpha_l \cdot ||\Phi_{image}^l(I_b) - \Phi_{image}^l(I_t)||_2,$$

The enhancement loss *L*<sub>enhance</sub>

$$\mathcal{L}_{enhance} = \mathcal{L}_{clip} + w \cdot \mathcal{L}_{identity},$$

In the second stage, we iteratively perform **prompt refinement** and **enhancement network tuning**. The prompt refinement and the tuning of the enhancement network are conducted in an **alternating manner**.



The negative similarity score between the prompt pair and an image:

$$S(I) = \frac{e^{\cos(\Phi_{image}(I), \Phi_{text}(T_n))}}{\sum_{i \in \{n, p\}} e^{\cos(\Phi_{image}(I), \Phi_{text}(T_i))}},$$

Margin ranking loss:

$$\mathcal{L}_{prompt1} = \max(0, S(I_w) - S(I_b) + m_0) + \max(0, S(I_t) - S(I_b) + m_0) + \max(0, S(I_w) - S(I_b) + m_0) + \max(0, S(I_w) - S(I_t) + m_1), + \max(0, S(I_w) - S(I_t) + m_1) + \max(0, S(I_t) - S(I_{t-1}) + m_2),$$



# **Experiments**



Input



Afifi et al.



URetinex-Net



SNR-Aware-LOLv1



SNR-Aware-LOLv2Real



SNR-Aware-LOLv2syn



**RUAS-retrained** 



SCI-retrained







**CLIP-LIT (Ours)** 

Zhao et al.-LOL

#### **Experiments**

Table 1: Quantitative comparison on the BAID test dataset. The best and second Table 2: Quantitative comparison on the performance are marked in red and blue.

Backlit300 test dataset.

Type	Methods	<b>PSNR</b> ↑	SSIM↑	LPIPS↓	MUSIQ↑	Methods	<b>MUSIQ</b> ↑
Input		16.641	0.768	0.197	52.115	Input	51.900
	Afifi et al. [1]	15.904	0.745	0.227	52.863	Afifi et al. [1]	51.930
Supervised	Zhao et alMIT5K [34]	18.228	0.774	0.189	51.457	Zhao et alMIT5K [34]	50.354
	Zhao et alLOL [34]	17.947	0.822	0.272	49.334	Zhao et alLOL [34]	48.334
	URetinex-Net [27]	18.925	0.865	0.211	54.402	URetinex-Net [27]	51.551
	SNR-Aware-LOLv1 [28]	15.472	0.747	0.408	26.425	SNR-Aware-LOLv1 [28]	29.915
	SNR-Aware-LOLv2real [28]	17.307	0.754	0.398	26.438	SNR-Aware-LOLv2real [28]	30.903
	SNR-Aware-LOLv2synthetic [28]	17.364	0.752	0.403	23.960	SNR-Aware-LOLv2synthetic [28]	29.149
	Zero-DCE [9]	19.740	0.871	0.183	51.804	Zero-DCE [9]	51.250
Unsupervised	Zero-DCE++ [15]	19.658	0.883	0.182	48.573	Zero-DCE++ [15]	48.216
	RUAS-LOL [18]	9.920	0.656	0.523	37.207	RUAS-LOL [18]	40.329
	RUAS-MIT5K [18]	13.312	0.758	0.347	45.008	RUAS-MIT5K [18]	44.523
	RUAS-DarkFace [18]	9.696	0.642	0.517	39.655	RUAS-DarkFace [18]	48.216
	SCI-easy [20]	17.819	0.840	0.210	51.984	SCI-easy [20]	50.642
	SCI-medium [20]	12.766	0.762	0.347	44.176	SCI-medium [20]	48.216
	SCI-diffucult [20]	16.993	0.837	0.232	52.369	SCI-diffucult [20]	49.428
	EnlightenGAN [10]	17.550	0.864	0.196	48.417	EnlightenGAN [10]	48.308
	ExCNet [31]	19.437	0.865	0.168	52.576	ExCNet [31]	50.278
Unsupervised (retrained)	Zero-DCE [9]	18.553	0.863	0.194	49.436	Zero-DCE [9]	48.491
	Zero-DCE++ [15]	16.018	0.832	0.240	47.253	Zero-DCE++ [15]	46.000
	RUAS [18]	12.922	0.743	0.362	45.056	RUAS [18]	45.251
	SCI [20]	16.639	0.768	0.197	52.265	SCI [20]	51.960
	EnlightenGAN [10]	17.957	0.849	0.182	53.871	EnlightenGAN [10]	48.261
	CLIP-LIT (Ours)	21.579	0.883	0.159	55.682	CLIP-LIT (Ours)	52.921

# **Ablation Studies**

#### Necessity of Learned Prompt



Method	<b>PSNR</b> ↑	SSIM↑
fixed prompts (backlit/well-lit)	14.748	0.823
w/o ranking losses (w/o Eqs. (7) and (8))	20.884	0.865
w/o $t - 1$ outputs (w/o Eq. (8))	20.146	0.866
Ours	21.579	0.883

#### Effectiveness of Iterative Learning



# **Ablation Studies**

#### Impact of Training Data

Table 4: Comparison of training data impact. The quantitative comparisons are conducted on BAID test dataset.

Reference images	<b>PSNR</b> ↑	SSIM↑	LPIPS↓	<b>MUSIQ</b> ↑
MIT5K [4]+DIV2K [2]	21.413	0.881	0.162	56.494
DIV2K [2]	21.579	0.883	0.159	55.682



Input Trained with DIV2K Trained with DIV2K+MIT5K Figure 13: Visual comparisons of our method trained using different reference images.

➢ Advantage of CLIP-Enhance Loss over the Adversarial Loss.

Table 5: Comparison between CLIP-Enhance loss and adversarial loss on the BAID test dataset.

loss	<b>PSNR</b> ↑	SSIM↑	LPIPS↓	<b>MUSIQ</b> ↑
Adversarial loss	17.407	0.785	0.194	52.416
<b>CLIP-Enhance</b> loss	21.579	0.883	0.159	55.682

## **Ablation Studies**

Impact of Different Prompt Initialization



Settings	<b>PSNR</b> ↑	<b>SSIM</b> ↑	LPIPS↓	<b>MUSIQ</b> ↑
Pure random initialization	21.237	0.884	0.158	55.959
Random initialization+backlit/well-lit	21.527	0.882	0.159	55.946
Random initialization+low light/normal light	21.579	0.883	0.159	55.682

Figure 15: Prompt initialization learning investigation.

# Thanks