

# Alignment-free HDR Deghosting with Semantics Consistent Transformer

## - Supplementary Material -

### Abstract

In the supplementary material, we provide in Section 1 the ablation study on the hyperparameters, in Section 2 a detailed study on the existing benchmark, in Section 3 the generalization capability with qualitative comparisons, and in Section 4 the full quantitative comparisons on our dataset.

### 1. Ablation Studies on the Number of Layers

In this section, we conduct studies to analyze the influence of hyperparameters, *i.e.*, numbers of attention blocks, on the HDR deghosting performance. We have three hyperparameters ( $N_L$ ;  $N_G$ ;  $N_S$ ), as shown in Figure 2 of the manuscript.  $N_L$  stands for the number of attention layers,  $N_G$  stands for the number of global spatial attention blocks, and  $N_S$  stands for the number of semantic-consistent attention blocks. The quantitative results under different settings can be found in Table 1. We can conclude that our method performs well with a large variation of hyperparameters.

### 2. Study on Kalantari *et al.* [2] Testing Samples

In this section, we first conduct a detailed analysis of the current benchmark Kalantari *et al.* In Figure 2 we plot the general performance by averaging the  $l$ -PSNR and  $\mu$ -PSNR scores of SOTA methods [4, 3], including ours, on each Kalantari’s testing sample. It can be seen that for samples 7 to 10, HDR deghosting networks produce low linear  $l$ -PSNR values (in blue color). However, after tone mapping, the  $\mu$ -PSNR (in orange color) becomes significantly higher, leading to a large but abnormal gap (in red color) between these two metrics. We think that this phenomenon may be majorly coupled with the presence of over-exposed/underexposed regions in the predicted images or ground truth.

To validate our initial guess, we evaluate the number of over-exposed pixels in the corresponding ground truth. In our study, we define a pixel as over-exposed if its value is greater than 95% of the maximum value that can be encoded in the HDR ground truth. The proportion of over-exposed

Table 1. Ablation study on the hyperparameters. The comparison is conducted on the Kalantari *et al.* dataset [2].

$N_L$	$N_G$	$N_S$	Mb	$\mu$ -PSNR	$l$ -PSNR	$\mu$ -SSIM	$l$ -SSIM
4	6	4	29	<b>44.49</b>	42.29	<b>0.9924</b>	<b>0.9887</b>
4	7	4	31	43.59	42.25	0.9912	0.9885
4	5	4	21	43.85	42.21	0.9912	0.9883
4	4	4	21	43.81	42.08	0.9913	0.9886
4	6	2	28	43.82	<b>42.30</b>	0.9915	0.9882
4	6	5	30	43.76	42.18	0.9913	0.9880
4	5	5	23	43.57	41.09	0.9910	0.9879
5	6	4	37	43.91	42.05	0.9911	0.9881
3	6	4	23	43.23	41.43	0.9910	0.9878
2	6	4	15	43.62	41.42	0.9909	0.9874

Table 2. Proportion of over-exposed pixels in Kalantari *et al.* [2] testing samples.

Sample ID	Num. over-exposed pixels	Proportion %
1	35	< 1
2	5	< 1
3	1060	< 1
4	10793	< 1
5	251	< 1
6	5685	< 1
7	187506	12.50
8	111788	7.45
9	213696	14.25
10	458368	30.56
11	14035	< 1
12	1060	< 1
13	1	< 1
14	251	< 1
15	5644	< 1

pixels in each sample can be found in Table 2. The visualization can be found in Figure 3. It can be seen that samples 7 to 10 contain more over-exposed pixels than others. The corresponding over-exposed mask can be found in Figure 1 for these samples. It can be seen that a large number of the background pixels are over-exposed. We think that the over-exposition may be linked to several factors: (1) a too-long medium exposure time making the reference image already over-exposed; (2) too-small differences in the

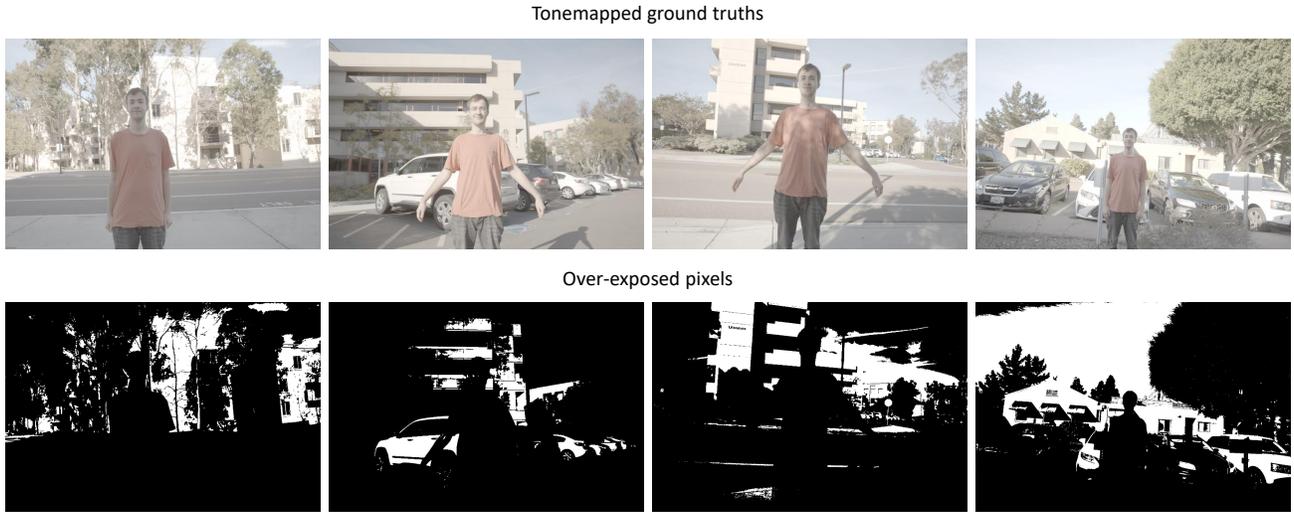


Figure 1. Exposition masks for Kalantari [2] testing samples where state-of-the-art dehazing methods [4, 3] yield unsatisfactory results in the linear domain. The exposition masks use white pixels to denote over-exposed regions and black pixels to represent reasonably-exposed regions. It can be observed that most of the backgrounds are over-exposed.

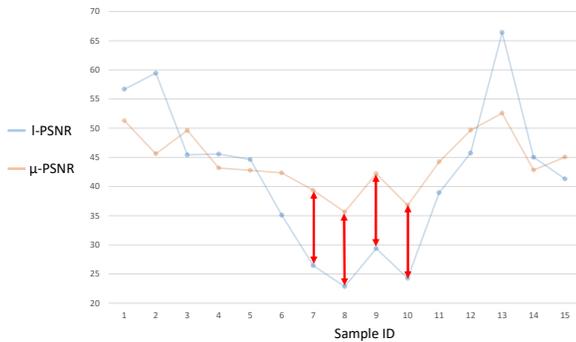


Figure 2. Average  $l$ -PSNR and  $\mu$ -PSNR scores from SOTA methods, including ours, on Kalantari [2] testing samples. For samples 7 to 10, the general performance is poor in the linear domain. Meanwhile, after the tonemapping function, the  $\mu$ -PSNR score becomes significantly higher, yielding an undesirable but important gap, in red color, between these metrics. This phenomenon may be caused due to the large number of over-exposed pixels in these samples. Please zoom in for more details.

exposure time during input LDR images, making it impossible to cover a larger dynamic range; (3) the conventional 3 input LDR images may not be sufficient.

Therefore, while collecting our dataset, we followed a very rigorous processing as described in our main manuscript. Following the same protocol, we show in Figure 4 that none of our samples contains more than 2% of over-exposed values. In addition, even though our dataset follows the conventional setting with 3 LDR inputs, we will release the whole bracket of 9 input exposures. In such a case, future works can benefit from a larger dynamic range to design better dehazing methods.

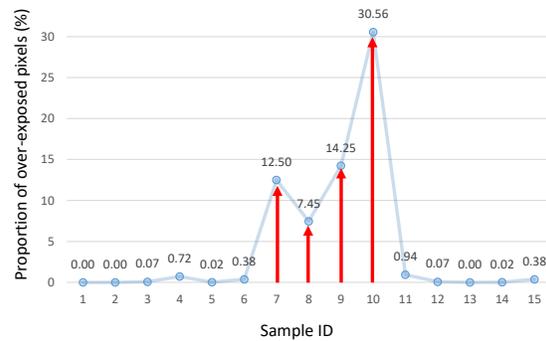


Figure 3. Percent of over-exposed pixels for each Kalantari [2] testing sample. We can find the anomalies in samples 7 to 10 where more than 5% of pixels are over-exposed.

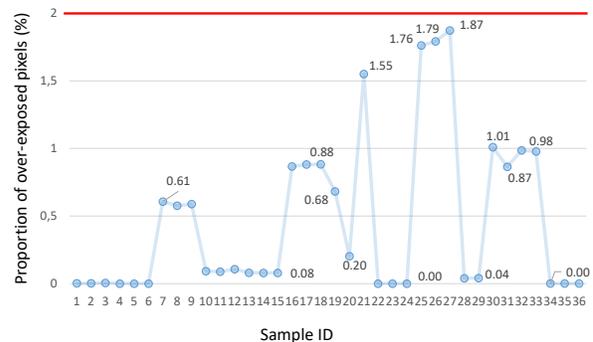


Figure 4. Compared to the existing benchmark, the proportion of over-exposed pixels is consistently lower than 2% in our dataset.

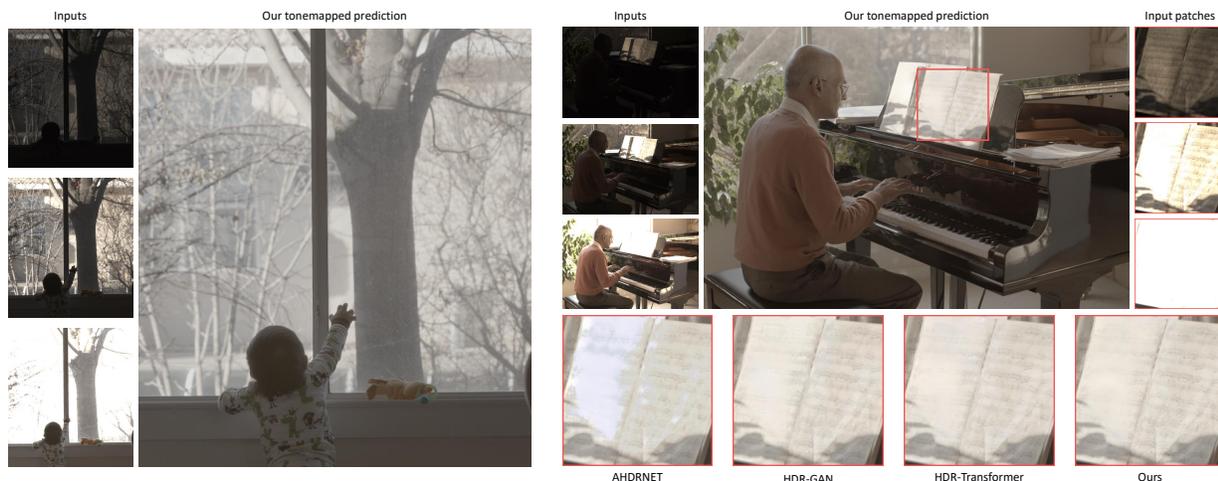


Figure 5. Evaluation of the generalizability of our solution using the Sen *et al.* [5] unsupervised dataset. All the compared networks are trained with our proposed dataset. It can be seen that our network reproduces a better texture of the piano scores book. Please zoom in for more details.

Table 3. Quantitative comparison with state-of-the-art methods on our proposed dataset.  $l$ -PSNR and  $l$ -SSIM are computed in the linear domain while  $\mu$ -PSNR and  $\mu$ -SSIM are computed after  $\mu$ -law tone mapping. PU-PSNR and PU-SSIM are calculated by applying the encoding function proposed in [1]. The compared methods are trained through their official implementation.

Method	$\mu$ -PSNR	PU-PSNR	$l$ -PSNR	$\mu$ -SSIM	PU-SSIM	$l$ -SSIM	HDR-VDP2
Sen <i>et al.</i> [5]	39.97	39.53	44.21	0.9792	0.9397	0.9932	67.2021
DHDRNet[2]	41.67	39.49	47.33	0.9838	0.9779	0.9961	72.2474
AHDRNet[6]	44.16	42.78	50.29	0.9896	0.9795	0.9971	78.1245
HDR-Transformer[3]	44.88	43.11	51.09	0.9904	0.9852	<b>0.9981</b>	78.8682
<b>Ours</b>	<b>44.93</b>	<b>44.78</b>	<b>51.73</b>	<b>0.9906</b>	<b>0.9950</b>	<b>0.9981</b>	<b>79.5267</b>

### 3. Qualitative Evaluation on Unsupervised Benchmarks

In order to verify the generalizability of our method, we conducted evaluations on the datasets proposed by Sen *et al.* [5]. All the compared networks are trained using our proposed dataset. As depicted in Figure 5, all other methods [6, 4, 3] exhibit distortion in over-exposed areas, while our network is able to reproduce the texture of the piano scores book most accurately.

### 4. Quantitative Performance on Our Dataset

The full performance can be found in Table 3.

### References

- [1] Maryam Azimi et al. Pu21: A novel perceptually uniform encoding for adapting existing quality metrics for hdr. In *2021 Picture Coding Symposium (PCS)*, pages 1–5. IEEE, 2021. 3
- [2] Nima Khademi Kalantari, Ravi Ramamoorthi, et al. Deep high dynamic range imaging of dynamic scenes. *ACM TOG*, 36(4):144–1, 2017. 1, 2, 3
- [3] Zhen Liu, Yinglong Wang, Bing Zeng, and Shuaicheng Liu. Ghost-free high dynamic range imaging with context-aware transformer. In *ECCV*, 2022. 1, 2, 3
- [4] Yuzhen Niu, Jianbin Wu, Wenxi Liu, Wenzhong Guo, and Rynson WH Lau. Hdr-gan: Hdr image reconstruction from multi-exposed ldr images with large motions. *IEEE TIP*, 30:3885–3896, 2021. 1, 2, 3
- [5] Pradeep Sen, Nima Khademi Kalantari, Maziar Yaesoubi, Soheil Darabi, Dan B Goldman, and Eli Shechtman. Robust Patch-Based HDR Reconstruction of Dynamic Scenes. *ACM TOG*, 31(6):203:1–203:11, 2012. 3
- [6] Qingsen Yan, Dong Gong, Qinfeng Shi, Anton van den Hengel, Chunhua Shen, Ian Reid, and Yanning Zhang. Attention-guided network for ghost-free high dynamic range imaging. *CVPR*, 2019. 3