

# How to cheat with metrics in single-image HDR reconstruction

## Supplementary material

Gabriel Eilertsen<sup>1</sup>, Saghi Hajisharif<sup>1</sup>, Param Hanji<sup>2</sup>, Apostolia Tsirikoglou<sup>1</sup>,  
Rafał K. Mantiuk<sup>2</sup>, Jonas Unger<sup>1</sup>

<sup>1</sup> Dept. of Science and Technology, Linköping University, Sweden

<sup>2</sup> Dept. of Computer Science and Technology, University of Cambridge, UK

### 1. Dataset and camera simulation

Figure 1 shows examples of HDR scenes used in the experiments. Figure 2 demonstrates the impact of the different camera simulations that were used.

### 2. Complementing results

- Figure 3 complements Figure 3 in the main paper, with the same results but for EV-10 instead of EV-5.
- Figure 4 and Figure 5 complement Figure 4 in the main paper, with statistical testing of rankings for all metrics and camera simulations.
- Figure 8 complements Figure 5 in the main paper, with plots for all three different metrics.
- Figure 9 and Figure 10 complement Figure 6 in the main paper, with similar examples for other scenes.

### 3. Supplementary results

Figure 6 shows the same evaluations as in Figure 3 in the main paper, but evaluated only on saturated pixels. That is, if the metric is  $d(\hat{H}, H)$ , these results have been computed using  $d(\alpha\hat{H}, \alpha H)$ , where  $\alpha = \max(0, L - 0.9)/0.1$  masks out only the saturated pixels. The rankings with testing of significant differences are given in Figure 7.

The comparison of only saturated regions can clearly separate P-rec as the best model. This is expected, since this has been composed using the ground truth information in the saturated pixels. It is also possible to see how the naive model is inferior, since this does not contain any information in the saturated pixel areas. However, for all other methods it is difficult to clearly separate between methods. Many methods also show differences in ranking depending on camera simulation.

P-lin does not contain any information in saturated regions, similar as with the naive model. Still, P-lin shows slightly higher mean. This is likely due to the blending

performed by  $\alpha$ , which incorporates also some information around the saturated regions.

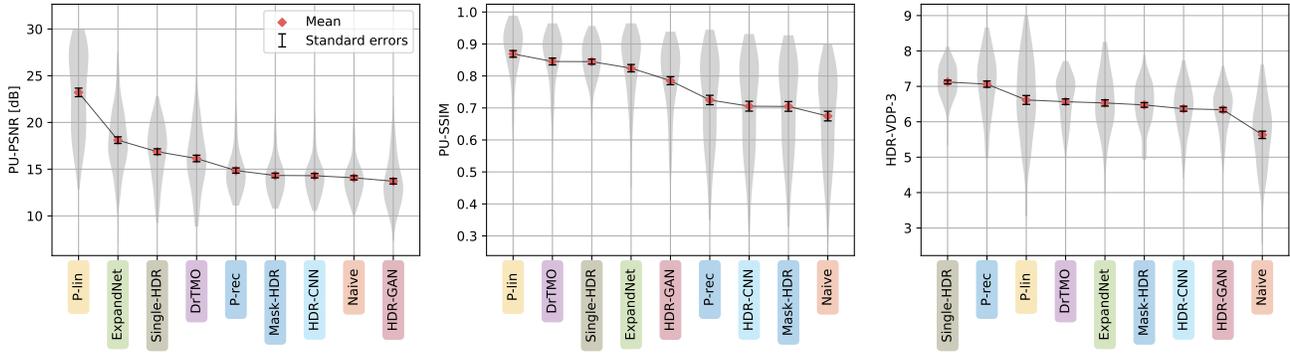
The high variance, and small differences in mean between good models and the naive one, points to how this type of comparison also is inadequate, where quality of linearization still has a prominent effect on the results. There is a need for better separation of the variations in the reconstruction problem.



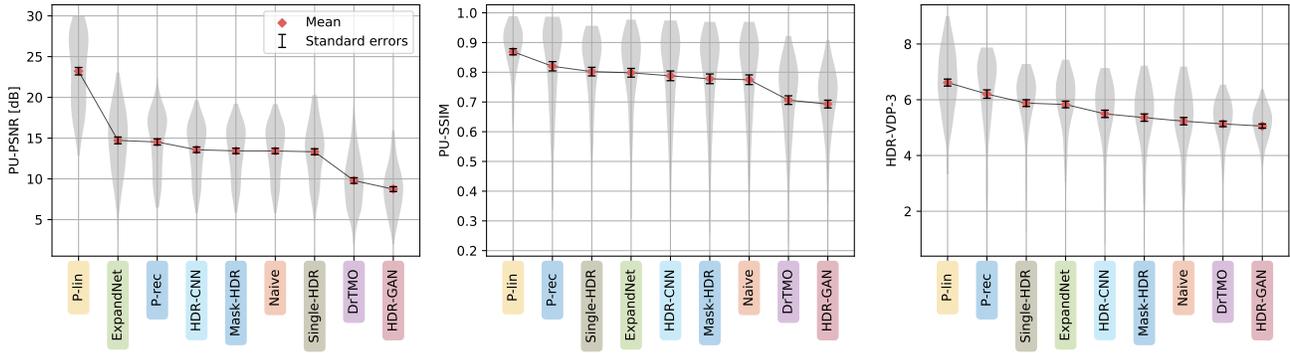
Figure 1: 15 example HDR scenes from the 96 scenes used in evaluation. Images have been gamma-encoded for display.



Figure 2: Examples of camera simulation, showcasing the results of different CRFs and exposures used for evaluation.



(a) Camera simulation: M-CRF, EV-10



(b) Camera simulation: CLAHE, EV-10

Figure 3: The distribution of metric values over the 96 tested scenes, where methods have been sorted by mean value to facilitate comparing differences in ranking. (a) uses camera simulation with M-CRF, while (b) is with CLAHE, and both have been simulated with EV-10. Left, middle, and right show results with PU-PSNR, PU-SSIM, and HDR-VDP-3, respectively.

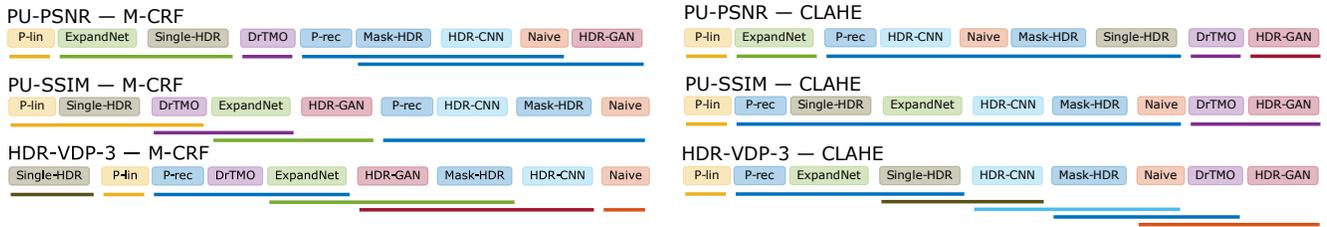


Figure 4: Rankings for different camera simulations with EV-5. The lines connect methods where the differences cannot be deemed statistically significant in a t-test, with a p-value threshold of 0.05.

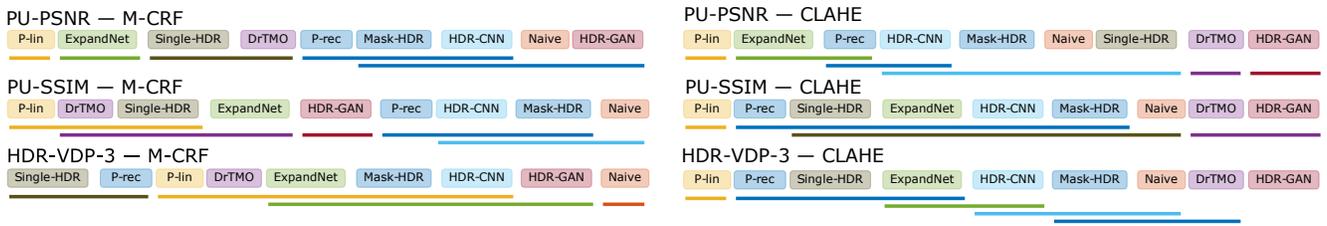
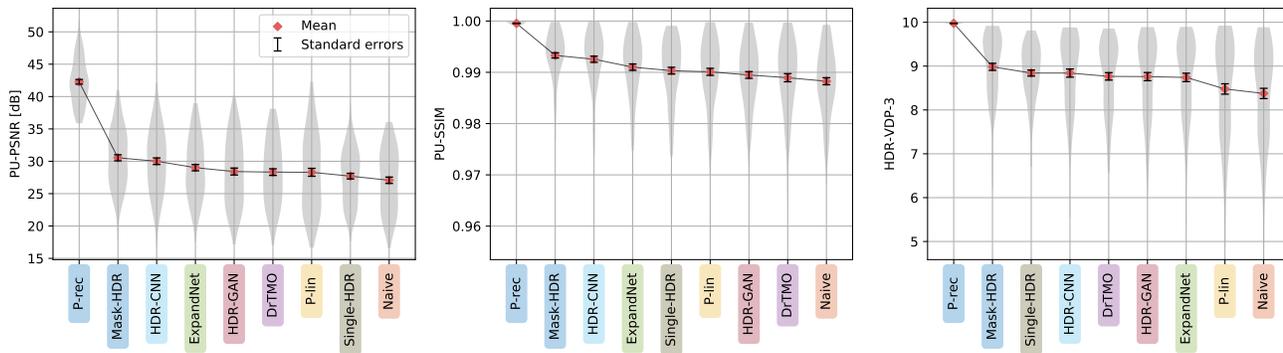
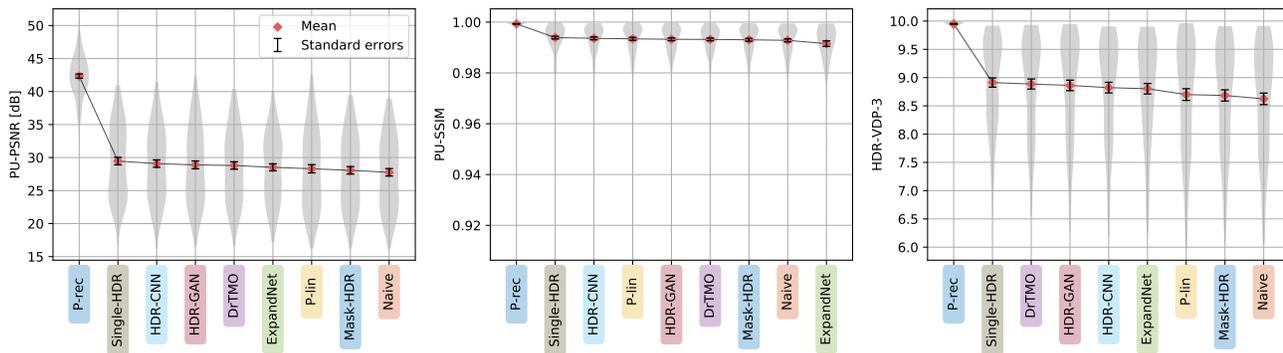


Figure 5: Rankings for different camera simulations with EV-10. The lines connect methods where the differences cannot be deemed statistically significant in a t-test, with a p-value threshold of 0.05.



(a) Camera simulation: M-CRF, EV-5



(b) Camera simulation: CLAHE, EV-5

Figure 6: Same as Figure 3, but for EV-5 and only evaluated in saturated regions of the images.

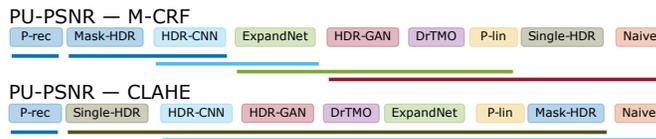


Figure 7: The rankings provided by PU-PSNR with M-CRF and CLAHE, when evaluated only in saturated image regions.

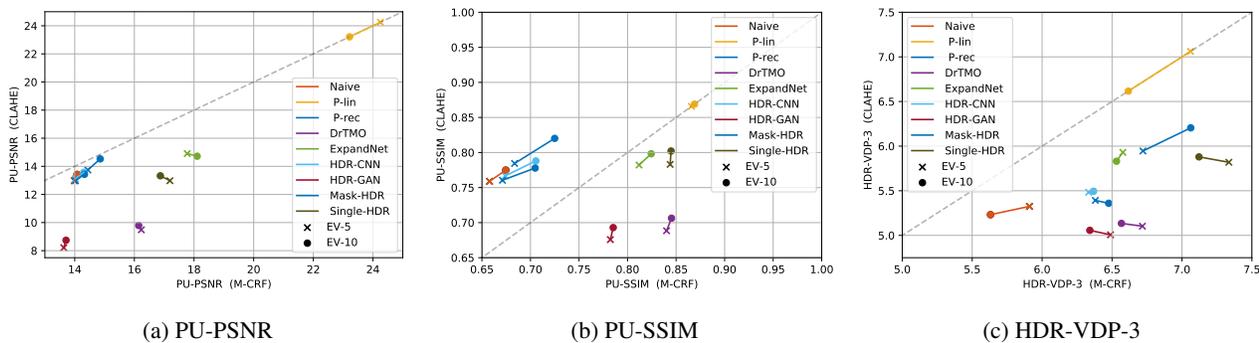


Figure 8: Differences in PU-PSNR (a), PU-SSIM (b) and HDR-VDP-3 (c) when using M-CRF and CLAHE. The two points for each method are with EV-5 and EV-10, demonstrating how many methods do not show an expected reduction in quality with increased camera exposure (more challenging reconstruction problem).

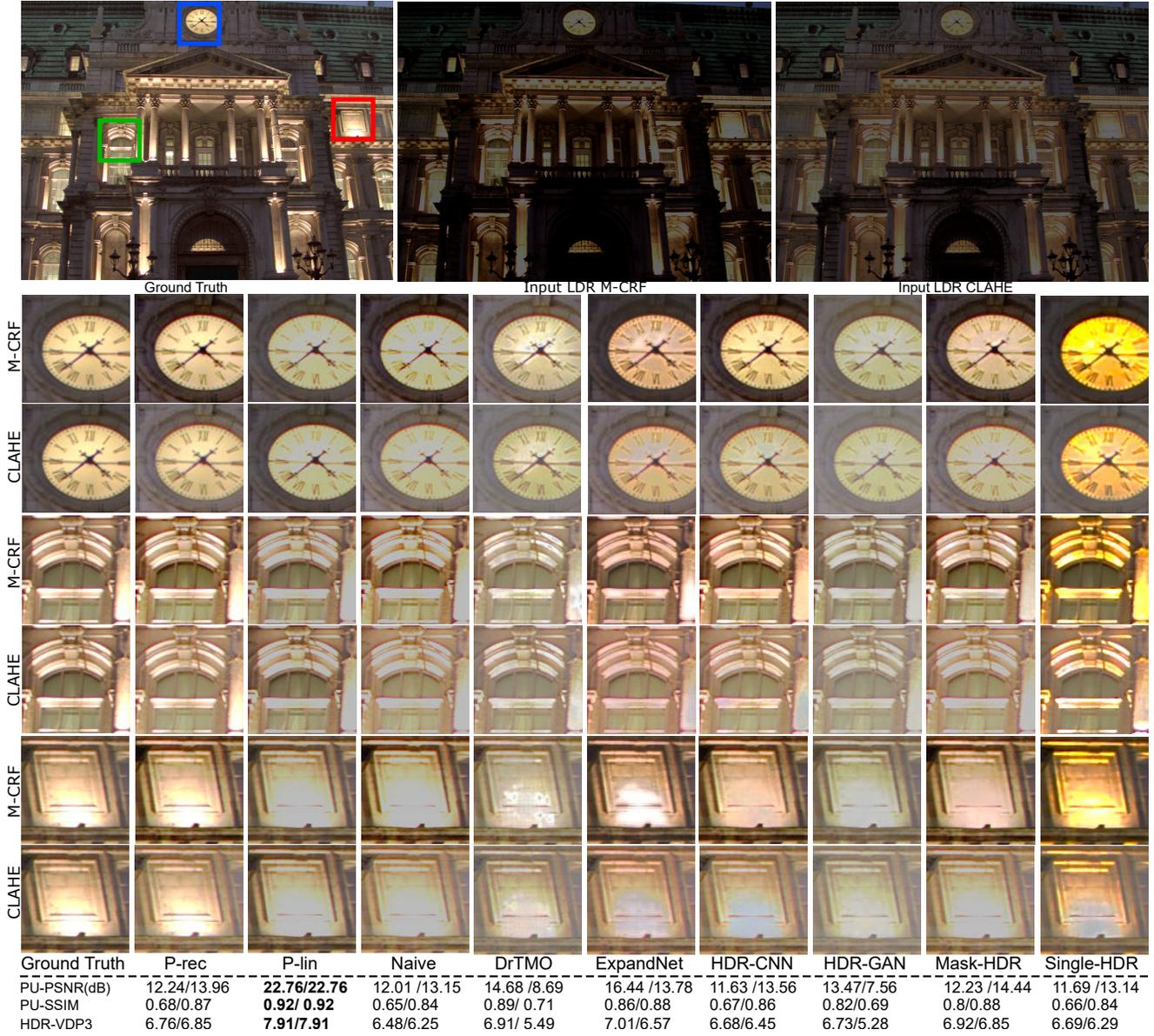


Figure 9: Selected scene areas for different reconstructions, with input LDR images simulated using M-CRF and CLAHE. The metrics in the bottom show the performance with M-CRF/CLAHE. The exposure time was set such that 5% of pixels are saturated (EV-5). Ground truth and reconstructed HDR images have been gamma-encoded for display.

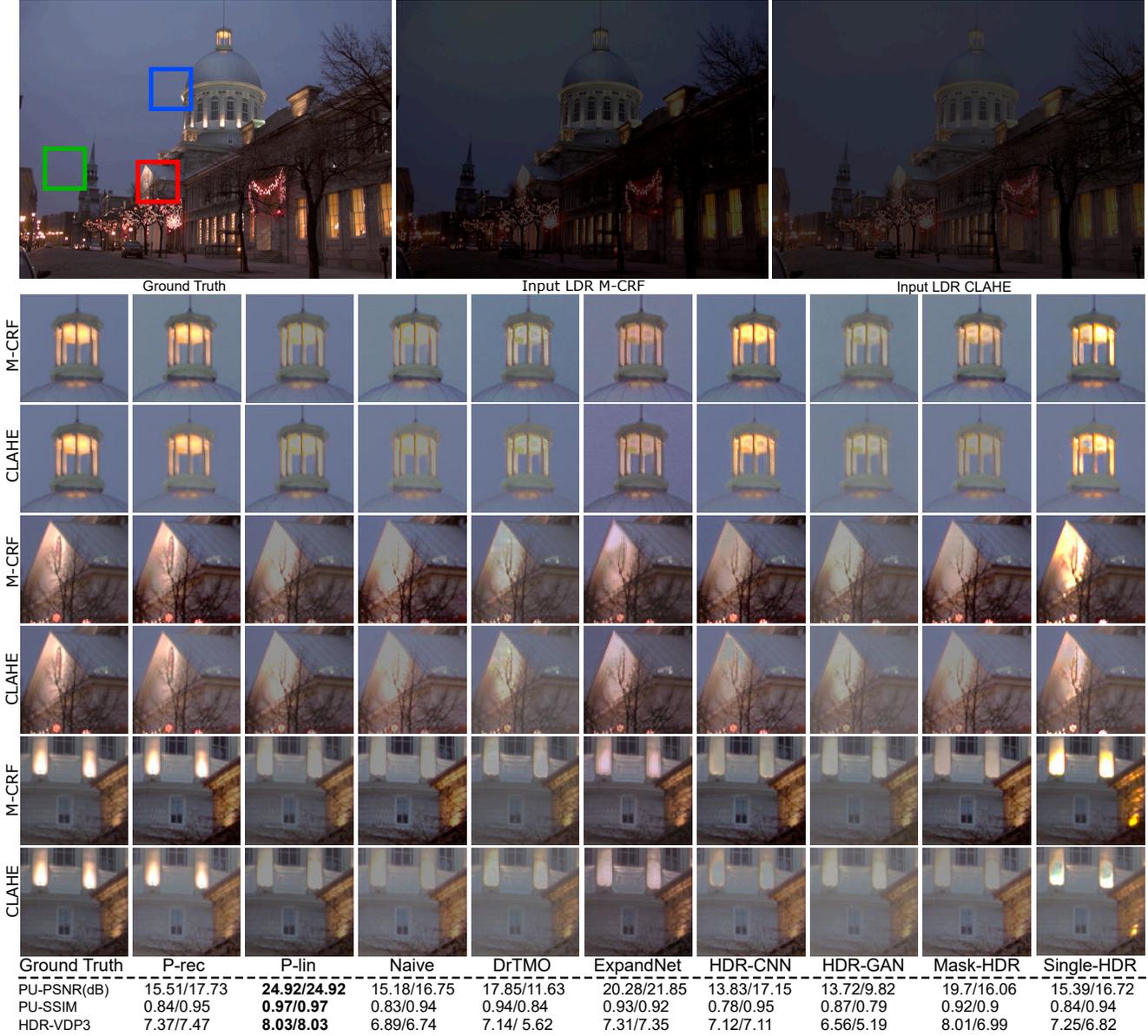


Figure 10: Selected scene areas for different reconstructions, with input LDR images simulated using M-CRF and CLAHE. The metrics in the bottom show the performance with M-CRF/CLAHE. The exposure time was set such that 5% of pixels are saturated (EV-5). Ground truth and reconstructed HDR images have been gamma-encoded for display.