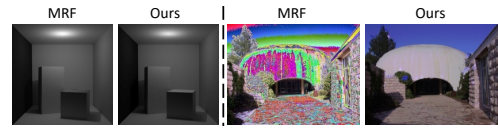


1 We sincerely thank all reviewers for their valuable comments and suggestions. We respectfully think it should be  
2 interesting to share with the NeurIPS community how to learn unwrapping a natural image, though the camera might be  
3 obscure. We will fix the typos and improve the figures as suggested in the final version. Below we respond to specific  
4 comments and concerns.

5 **R1&R2: Results of the MRF algorithm [55].** Those results are indeed generated with the original implementation (we have got the  
6 code from the authors). We noticed and confirmed with the authors that the parameters of the MRF algorithm (especially the energy function  
7 and the cliques to consider) need carefully fine-tuning on different input data to achieve optimal results, and the fine-tuning process is quite tricky. Moreover, for RGB data, it independently  
8 unwraps each channel with different parameters. According to our experiments, the images shown in [55] should be  
9 "cherry-picked", as mentioned by R1. Since our method uses the same set of parameters for all test cases, for a fair  
10 comparison, we fix the parameters of the MRF algorithm for evaluation as well. As shown in the right figure (examples  
11 are from Figure 4 and 9 of [55]), by applying the MRF algorithm using the same set of parameters, the left image is  
12 well recovered while the right one is not. In comparison, our learning-based method gives consistently good results.



16 **R1: Other non-learned techniques for processing modulo imagery.** Please kindly note that our method only takes a  
17 *single* modulo image as input. We did not find other suitable baselines. Papers mentioned in the review "Robust Multi-  
18 Image HDR Reconstruction for the Modulo Camera", "Reconstruction from Periodic Nonlinearities, with Applications  
19 to HDR Imaging", "Signal Reconstruction from Modulo Observations", all require multiple images as input (what's  
20 more, the model presented in the first paper is limited to two modulo periods). Multiple modulo images could be  
21 merged for HDR as done in [55], but such a multi-image approach is beyond the scope of our paper. We will clarify this  
22 assumption in the final version, and consider learning to unwrap multiple modulo images as our future work.

23 **R1: Evaluations on other conventional CNN models.** We respectfully point out that we have included such results in  
24 our ablation study. For example, we have tested a vanilla attention U-Net structure that predicts the number of rollovers  
25 **K** (referring to "Directly predicting **K**" in Table 2 of the paper). To show the effectiveness of modulo edge separator, we  
26 have also tried a single rollover mask predictor to directly output binary rollover mask (referring to "Removing modulo  
27 edge separator" in Table 2 of the paper). The final rollover mask predictor adopts a modified attention U-Net structure,  
28 and it outperforms a vanilla attention U-Net structure by 2.21 dB for PSNR according to our previous experiments  
29 (when we were designing our network). We will include the results of such a vanilla attention U-Net in the final version.

30 **R2: Other deep-learning based 2D phase unwrapping methods.** We have tried our best to search for the code of  
31 such methods, *e.g.*, [41, 44, 49], but none of them is available now. We are afraid that re-implementing them within  
32 the rebuttal period could not produce reliable results. We will try to include such a comparison in the final version.  
33 However, we suspect these methods designed for phase maps instead of modulo images could not compete with either  
34 [55] or our method, due to the essentially different natures of phase maps and wrapped natural images.

35 **R2: Details about the tone mapping method.** We use the OpenCV-Python's Reinhard tone mapping algorithm  
36 (`cv2.createTonemapReinhard(intensity=-1.0, light_adapt=0.8, color_adapt=0.0)`) for visualization.

37 **R2&R3: Implementation details of the MRF algorithm [55].** The code of [55] is written in MATLAB script and  
38 cannot run on GPU, so we run it on an Intel Core i7-8700K CPU (using a single core), while UnModNet runs on an  
39 NVIDIA 2080 Ti GPU (thanks to the PyTorch's CUDA acceleration). At each iteration, UnModNet takes around 200ms  
40 to process a  $512 \times 512$  modulo image, and [55] takes around 24s.

41 **R2: Evaluation metrics (SSIM, PSNR, MS-SSIM).** We have further conducted  
42 evaluations on the tone-mapped HDR images with metrics SSIM, PSNR, and MS-  
43 SSIM, and the advantage of our method did not change. For example, the PSNR  
44 (dB) results in Table 1 of the paper become 30.33 for our method, 7.73 for MRF  
45 [55], 19.99 for DrTMO [7], 18.79 for ExpandNet [30], and 18.68 for HDRCNN [6].  
46 We will include complete results in our final version.



47 **R2: High-resolution results.** Examples (scaled to the same height) are shown in the right figure.

48 **R2: Physical limitations and properties of the sensor.** The robustness might be improved if we include a proper  
49 camera noise model in our dataset generation pipeline, and we consider to add it in our future work. Our method could  
50 adapt to other bit depth (*e.g.*, 12-bit) if we change the bit depth of the modulo images in the training dataset.

51 **R3: Pros and cons of different single-image HDR reconstruction methods.** We will improve the the related work  
52 section by comprehensively discussing pros and cons of each method.

53 **R4: Details of our dataset.** Our training dataset consists of 21312 pairs of modulo and updated (referring to Equation  
54 (1) of the paper) image patches in  $256 \times 256$  resolution, which are randomly cropped and generated from 600 HDR  
55 images collected from various sources [12, 11, 13, 14, 37, 53, 10, 27]. The details of our dataset creation pipeline and  
56 the parameter settings could be found in Section 7 of the supplementary material.

57 **R4: Evidence of not overfitting and the generalization applicability.** The evaluation on two types of real captured  
58 data could prove that our method does not overfit, and has a good generalization ability.