

Joint High Dynamic Range Imaging and Super-Resolution from a Single Image

Jae Woong Soh¹ Jae Sung Park² Nam Ik Cho¹

Department of ECE, INMC, Seoul National University, Seoul, Korea¹
 Samsung Electronics Co. Ltd., Suwon, Korea²

soh90815@ispl.snu.ac.kr, jason79.park@samsung.com, nicho@snu.ac.kr

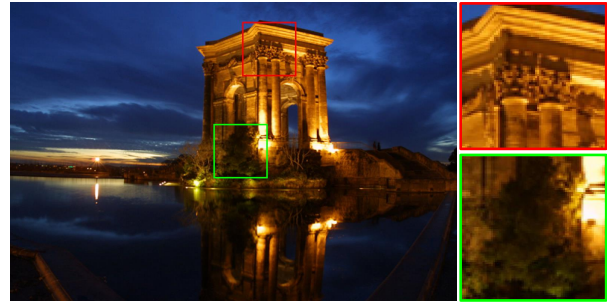
Abstract

This paper presents a new framework for jointly enhancing the resolution and the dynamic range of an image, i.e., simultaneous super-resolution (SR) and high dynamic range imaging (HDRI), based on a convolutional neural network (CNN). From the common trends of both tasks, we train a CNN for the joint HDRI and SR by focusing on the reconstruction of high-frequency details. Specifically, the high-frequency component in our work is the reflectance component according to the Retinex-based image decomposition, and only the reflectance component is manipulated by the CNN while another component (illumination) is processed in a conventional way. In training the CNN, we devise an appropriate loss function that contributes to the naturalness quality of resulting images. Experiments show that our algorithm outperforms the cascade implementation of CNN-based SR and HDRI.

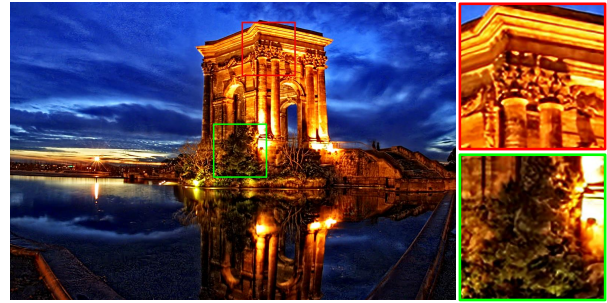
1. Introduction

With the advent of ultra high definition television (UHDTV) with HDR rendering [36], the techniques for capturing high resolution (HR) and high dynamic range (HDR) images have become important. In addition to developing the methods to capture the new UHD contents, it is also important to convert the vast amount of existing low-resolution (LR) and low dynamic range (LDR) images to the HDR-HR contents for rendering them on the UHDTV. For these purposes, there have been many methods to convert the LR to HR, which is called the single image SR (SISR). Also, there have been some single image HDRI algorithms to produce an HDR image from a single LDR input.

Since the SISR is an important problem that finds many useful applications, it has long been studied by many researchers [55, 48, 58, 16, 50, 5, 29, 30, 60]. Some earlier works exploited the statistical priors of images for the SISR [48, 46]. The learning based methods, specifically the ones based on the neighbor embedding [9] and sparse



(a) The input LDR-LR image.



(b) Our HDRI-SR-B result.

Figure 1: An example of LDR-LR input and the output of our method. It can be seen that both the dynamic range and resolution are enhanced.

coding methods [58, 57, 50] were also introduced for better SR. Recently, the state-of-the-art methods are based on the CNN [5, 22, 29, 28, 51, 30, 60, 17, 15], which show further enhanced results than the previous ones. Generative methods [29, 43, 54] based on generative adversarial networks (GAN) [10] have also been introduced for the better perceptual quality of SR images.

For the HDRI, multiple images or sensors with different exposures are used, or a single image is reversely tonemapped to generate an HDR output [38, 35, 2, 3, 26, 33, 6, 21, 7]. For the enhancement of undesirably illuminated images, some methods generated virtual exposure images

from a single input and applied the conventional HDRI process [49, 53, 52, 8, 14, 19, 24, 40]. Recently, deep CNN-based methods have also been proposed [6, 21, 7]. There are also some works that use both HDRI and SR for image enhancement: Schubert *et al.* [44] developed a method for the HDRI-SR which consists of several steps. Park *et al.* tried several cascade implementations of HDRI-SR in different color spaces [40].

Recently, deep CNNs have shown dramatic improvements in most of low to high-level vision tasks including super-resolution [5, 22, 29, 47, 28, 30, 60, 15], Gaussian denoising [59], and JPEG artifact reduction [4, 12]. For these problems, the deep CNN is shown to work as a proper mapping function from the degraded image to the original. In this respect, we adopt the CNN for the joint HDRI-SR problem.

The most straightforward method of using a CNN is to design an end-to-end network, *i.e.*, a deep network which takes LDR-LR image as the input and generates the corresponding HDR-HR. However, in the case of joint HDRI-SR task, we have a problem that there is no standardized labeled data. Precisely, the HDR images' luminance range and tonemapping function from the HDR to the LDR are not well specified in the current HDR datasets. Moreover, the dynamic ranges of target images are usually different from each other, and the tonemapping functions are also different and nonlinear due to the use of locally adaptive nonlinear mapping in most images. Thus, by directly training a discriminative CNN to map LDR images to HDR ones, the network usually fails to converge. Hence, we need to use a transformed image or find another domain that is less affected by the luminance range and tonemapping function.

Some of the previous works also showed that it is important to find an appropriate domain when applying a CNN for image enhancement. For example, it is shown in [12] that applying the dual domain representation, *i.e.*, using the image and DCT domain priors increases the performance of JPEG artifacts reduction compared to the methods without the DCT prior. Also, the SR with wavelet domain priors [13] enhanced the performance compared to the conventional image domain methods. Additionally, recent SR and denoising CNNs such as VDSR [22] and DnCNN [59] focus on the residual structure, because the SR task is to find the high-frequency details and the Gaussian denoising is also to estimate the noise which is the residual signal.

The HDRI also focuses on the reconstruction of image details rather than the low-frequency components. Specifically, recent single image HDRI methods process the illumination and reflectance components separately [8, 14, 19, 24, 40], where the illumination corresponds to the low-frequency component and the reflectance amounts to the image details. The illumination is simply scaled according to certain virtual exposure levels, while the image details are

locally and sophisticatedly manipulated to reveal the details in the saturated regions.

Considering that both HDRI and SR try to find the lost image details in the high-frequency region rather than the low-frequency, we design our joint HDRI-SR CNN to focus on the high-frequencies. Furthermore, it is possible to address the above-stated problem of HDRI training datasets by excluding the low-frequency illumination component in the training and inference, because the illuminations usually have very wide and different ranges from each other. Specifically, we propose a CNN architecture and its training schemes for enlarging the resolution and dynamic range of reflectance component. We also adopt the generative scheme [10], for generating better details and textures.

In summary, the main contributions of this paper are as follows.

- The proposed method performs HDRI and SR using a single CNN with high generalization performance, especially without well-organized labeled datasets for HDRI.
- The proposed method performs better than adopting separate CNNs for HDRI and SR in terms of perceptual quality and some no-reference metrics.
- The proposed single CNN requires much fewer parameters than the cascade of state-of-the-art SR and HDRI networks.
- Unlike the conventional single image HDRI that needs to generate virtual multi-exposure images, the proposed method directly produces the HDR image through the CNN which simultaneously performs the SR.

2. Related Work

2.1. Single Image Super-Resolution based on CNN

In recent years, the CNN-based SISR algorithms outperform most of the conventional non-CNN based methods. Since Dong *et al.* first proposed a CNN for the SISR [5], many other deep networks have been proposed. For some examples, the residual structure is used for better performance in [22], and the sub-pixel convolution layer is introduced in [45] to speed up. The SRGAN proposed in [29] generates the photorealistic SR images by exploiting the generative adversarial losses. Guo *et al.* [13] proposed a CNN for the SR in the wavelet domain, and have shown that Haar wavelet domain is an efficient one for the SISR. There are also many other structures and methods such as the recursive architecture [23, 47] and the Laplacian pyramids [28]. Very recently, enormously deep structures have been proposed and achieved the state-of-the-art performances [30, 60].

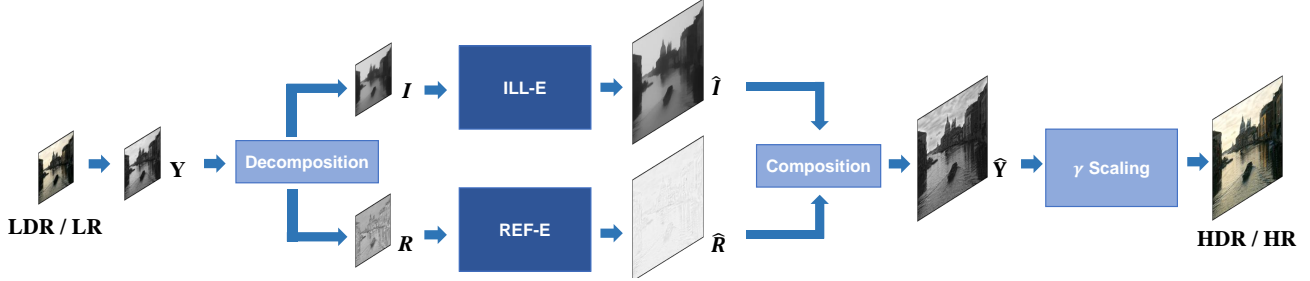


Figure 2: The overall process of the proposed scheme. We first decomposed the LDR-LR input into the illumination I and the reflectance R . ILL-E and REF-E denote ILLUMINATION Enhancement and REFLECTANCE Enhancement, respectively. The CNN is designed and trained only for the reflectance component R , and the illumination component I is simply up-scaled to increase the dynamic range.

2.2. High Dynamic Range Imaging from a Single Image

For generating an HDR image from a single input, most of the conventional methods generate the virtual multi-exposure images by applying brightness enhancement functions and then fuse them with the appropriate weight maps obtained from each of the virtual images [49, 53, 8, 14, 19, 24, 40]. Among these techniques, the Retinex based approaches [8, 14, 19, 24, 40] enhance the illumination and reflectance components separately. The undesirably illuminated regions are compensated by the estimated illumination, while the saturated details are enhanced by controlling the reflectance component. For generating a real (not the exposure fusion in low dynamic range) HDR image, the reverse tone mapping operators (rTMOs) with well-designed expanding maps are presented in [3, 18, 26, 33, 34]. Recently, a few CNN-based approaches [6, 21, 7] have been proposed. In [21], a CNN-based multi-exposure fusion scheme is proposed. In [7], a CNN-based reverse tone mapping method is proposed where virtual multi-exposure images are predicted by a CNN, and then they are fused for generating the HDR image. In [21, 7], the HDR image is generated as a weighted sum of multi-exposure images (real or virtually generated). In [6], the interest is on the saturated highlights and the CNN predicts these saturated regions to generate an HDR image.

3. Proposed HDRI-SR Scheme

The overview of our method is illustrated in Figure 2, which shows the process of only Y component because the $CbCr$ components are just bicubic interpolated (not shown in the figure). The figure shows that the Y is decomposed into the illumination (I) and reflectance (R), which undergo different processes and finally fused again to be an enhanced \hat{Y} .

For explaining the process formally in the rest of paper, we denote the HDR-HR image as X_{HH} , HDR-LR as X_{HL} ,

LDR-HR as X_{LH} , and LDR-LR as X_{LL} , where X is an RGB image. Also, $\{Y, Cb, Cr, I, R\}$ with one of these subscripts mean the corresponding components of X with the same subscript. Through the manipulation of these components, the final goal is to find a plausible X_{HH} from the given X_{LL} .

3.1. Image Decomposition

Image decomposition process has an important role in our scheme because it enables to exploit domain properties and to train CNN without consistent HDR ground-truth.

The reflectance is obtained as the difference between the luminance and the estimated illumination as

$$R = \log(Y) - \log(I). \quad (1)$$

The illumination I is estimated by the filtering of Y with a filter G , i.e., $I = Y * G(\cdot)$, where G is usually a normalized Gaussian in the conventional work. However, since it is known that using the Gaussian filter often makes the halo artifacts, we adopt the weighted least square (WLS) filter instead of the Gaussian. The WLS shows less halo artifacts because it preserves the edges better than the Gaussian filter. Precisely, the output image is obtained by solving an optimization problem, which is seeking the minimum of

$$\sum_p ((u_p - g_p)^2 + \lambda(a_{x,p}(g)|\nabla_x u_p|^2 + a_{y,p}(g)|\nabla_y u_p|^2)), \quad (2)$$

where p denotes the pixel location, g and u indicate the input and output respectively. Also, the smoothness weights are defined as

$$a_{x,p}(g) = (|\nabla_x l_p|^\alpha + \epsilon)^{-1}, \quad a_{y,p}(g) = (|\nabla_y l_p|^\alpha + \epsilon)^{-1}, \quad (3)$$

where l is the log of g , α is the parameter to control the sensitivity of the gradients of g , and ϵ is a small number to prevent the divide by zero. For all the training and test images, we set the parameters as $\lambda = 2$, $\alpha = 2$, and $\epsilon = 0.0001$.



Figure 3: Visualization of luminance map, estimated illumination, and reflectance component obtained by WLS filtering.

We also use linearized luminance values for decomposition. Figure 3 visualizes an example of luminance, estimated illumination, and reflectance obtained by the WLS filter, where we can see that the barely seen details in Y are revealed in R .

3.2. Illumination Enhancement

For the illumination enhancement, we first bicubically interpolate the I_{LL} to I_{LH} , and then simply compensate non-linearity to generate I_{HH} from I_{HH} . For directly generating I_{HH} from I_{LL} , we may use the CNN that is trained for the HDRI-SR of R , or we may design another CNN. However, using the CNNs do not improve the overall performance compared to the interpolation and stretching according to our experiments, because the illumination component is a smoothed image which contains little information.

To be specific with the compensation, we first bicubically interpolate the illumination values, then apply a simple gamma function $f(x) = x^{1/2.2}$ to compensate for the non-linearity. In summary, the ILL-E in Figure 2 is just the bicubic interpolation followed by gamma function.

3.3. Reflectance Enhancement

In this subsection, we present a CNN that maps R_{LL} to R_{HH} , which is named REF-Net. The overall architecture of the proposed REF-Net is shown in Figure 4 which has the stacked hourglass structure. Although U-Net [42] is shown to be effective in semantic segmentation, it may not be satisfactory for the prediction of reflectance components that have abundant textures. Hence, inspired by the work of successive stack of hourglasses [32, 39, 56], we stack two U-Nets for better prediction. For details, the size of all convolution layers are 3×3 and transposed convolution layers are 4×4 for REF-Net, respectively.

In manipulating the R component, we should note that the reflectance is the log ratio between luminance and illumination, which may ranges from $-\infty$ to $+\infty$ and hence not a suitable input to the CNN. It is found in our experiments that our CNN fails to predict high-quality reflectance when we feed R to the CNN without any preprocessing. To stabilize the CNN and to address this issue, we employ the tangent-hyperbolic function as a preprocessing step, which bounds the input to the CNN into $(-1, 1)$. In summary, our

REF-Net predicts the tanh reflectance of HDR-HR for the given tanh reflectance of LDR-LR, which can be described as

$$\tanh(\hat{R}_{HH}) = f_R(\tanh(R_{LL}); \theta_R), \quad (4)$$

$$\hat{R}_{HH} = \tanh^{-1}(f_R(\tanh(R_{LL}); \theta_R)), \quad (5)$$

where $f_R(\cdot)$ is our REF-Net with the parameter θ_R .

3.4. HDRI-SR Prediction

With the \hat{I}_{HH} and \hat{R}_{HH} , we recombine both components to build enhanced luminance \hat{Y}_{HH} as

$$\hat{Y}_{HH} = \hat{I}_{HH} \odot \exp(\hat{R}_{HH}), \quad (6)$$

where \odot denotes element-wise multiplication. Then the final HDR irradiance map is obtained by

$$\hat{X}_{HH} = (ycbcr2rgb([\hat{Y}_{HH}, Cb_{HH}, Cr_{HH}]))^\gamma, \quad (7)$$

where Cb_{HH} and Cr_{HH} are bicubic interpolation of Cb_{LL} and Cr_{LL} .

Finally, to display the \hat{X}_{HH} to an HDR display, it can be linearly stretched to the luminance value of the target HDR display. On the other hand, to display the \hat{X}_{HH} to an LDR display, the HDR irradiance map is tonemapped to be an enhanced LDR image [41].

4. Loss Function

In this section, we introduce loss functions that we use for our REF-Net.

Reconstruction Loss We adopt the reconstruction loss term as mean absolute error (MAE):

$$\mathcal{L}_{recon} = \frac{1}{N} \sum_i \|R_{HH}^i - f_R(R_{LL}^i; \theta_R)\|_1, \quad (8)$$

where i is the image index and N denotes the batch size.

Adversarial Loss For better sharpness and details, we adopt adversarial loss inspired by recent successful generative super-resolution models [29, 43, 54]. The generative scheme not only generates better sharpness and details but also enables to predict saturated regions such as washed-out areas and diminished dark pixels based on training dataset. We adopt recently proposed adversarial loss with relativistic discriminator [20], which shows great image quality with relativistic average standard GAN [10], named RaGAN. Specifically, the RaGAN loss for our scheme is described as:

$$\mathcal{L}_G = -\mathbb{E}_{x_r \sim \mathbb{P}_r} [\log(\tilde{D}(x_r))] - \mathbb{E}_{x_f \sim \mathbb{P}_g} [\log(1 - \tilde{D}(x_f))] \quad (9)$$

$$\mathcal{L}_D = -\mathbb{E}_{x_f \sim \mathbb{P}_g} [\log(\tilde{D}(x_f))] - \mathbb{E}_{x_r \sim \mathbb{P}_r} [\log(1 - \tilde{D}(x_r))], \quad (10)$$

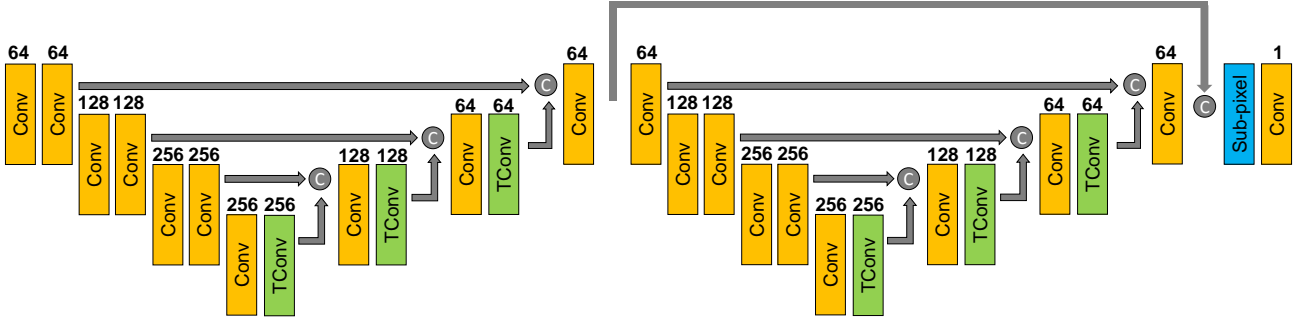


Figure 4: The overall CNN architecture of REF-Net. A stack of two hourglasses is adopted. The output of the first U-Net structure is skipped and concatenated with the second U-Net output, and then finally up-sampled via sub-pixel convolution layers.

where \mathbb{P}_r and \mathbb{P}_g are empirical distributions of R_{HH} and \hat{R}_{HH} respectively, x_r and x_f stand for real and generated data respectively, and

$$\tilde{D}(x_r) = \text{sigmoid}(C(x_r) - \mathbb{E}_{x_f \sim \mathbb{P}_g}[C(x_f)]) \quad (11)$$

$$\tilde{D}(x_f) = \text{sigmoid}(C(x_f) - \mathbb{E}_{x_r \sim \mathbb{P}_r}[C(x_r)]) \quad (12)$$

where $C(\cdot)$ denotes the output logit of discriminator. We show the discriminator architecture for the GAN training in the *supplementary material*.

Overall Loss We present two models: one is the basic model (HDRI-SR-B) without adversarial loss and the other is the complex model (HDRI-SR-C) with the adversarial loss. Formally, the overall loss for the basic model is \mathcal{L}_{recon} in eq. (8), and the loss for the complex model is

$$\mathcal{L} = \mathcal{L}_{recon} + \mu \mathcal{L}_G. \quad (13)$$

where we set $\mu = 10^{-3}$ for the training.

5. Training Strategies

For training the overall HDRI-SR, we use the MMPSG [25] dataset which consists of multi-exposure images and the HDR images constructed from the multi-exposures. Note that the input to the network is the LDR-LR and the output is the corresponding HDR-HR reflectance for the training. For constructing such set of input-output, we select 40 sets of multi-exposure images, and down-sample the standard-exposed images (not the over or under-exposed ones) as the input LDR-LR. The HDR images corresponding to the input LDR ones are selected as the output, where they are tonemapped.

To generate reflectance dataset pairs $\{R_{LL}^i, R_{HH}^i\}$, we obtain each reflectance from the linearized luminance value. To account for the variations and non-linearities caused by non-linear camera response functions (CRFs) [11], the inverse of CRF is required. However, since the CRFs are



Figure 5: MESet8. From the top left, Desktoys, House, Grandcanal, Lighthouse, Mask, Tower, Peyrou, and Landscape.

usually unknown and vary diversely, we assume it just a gamma function with $\gamma < 1$ and use the inverse of CRF for the linearization. Specifically, we use $f^{-1}(x) = x^{2.2}$ for the linearization, where $f(\cdot)$ is the assumed camera response function.

It is worth pointing out that such inconsistent characteristics of HDR datasets are alleviated by removing the illumination component from the LDR images and corresponding tonemapped HDR images. We have also tried to use the original HDR image for extracting the R (for removing illumination), however, we found that extraction from the tonmapped HDR image contains much better detail information, and eventually yields much better performance.

From the reflectance pairs so obtained, we extract the patches with the size of 48×48 in LR domain, which are augmented by rotation and flip. We only consider $\times 2$ for the scaling factor of super-resolution. The mini-batch size for training is set as 32, and the learning rate is decayed with the scale factor of 0.5, starting from 2×10^{-4} and halved once to 10^{-4} . We implement the model using Tensorflow [1] library with Titan Xp GPU.

6. Experimental Results

To evaluate the proposed method, we perform the experiments on several test sets. First, we select eight well-known

Table 1: Quantitative evaluation results on the test set. The numbers are averaged metrics for the datasets. Red indicates the best result and blue indicates the second best. In the case of HIGRADE, there are two kind of measures HIGRADE-1/HIGRADE-2.

Metric	Dataset	LDR-LR	Kovaleski-SR	Eilertsen-SR	HDRI-SR-B (Ours)	HDRI-SR-C (Ours)
NIQE(↓)	MESet8	5.157	4.280	4.544	4.230	2.746
	Part I	5.855	4.700	4.564	5.028	3.304
	Part II	5.333	4.558	4.600	4.649	3.141
	Part III	5.415	4.346	4.279	4.673	3.280
HIGRADE(↑)	MESet8	-1.148/-0.128	-0.367/0.024	-0.625/-0.107	-0.051/0.015	0.158/0.287
	Part I	-1.192/-0.578	-0.351/0.013	-0.657/-0.301	-0.253/-0.176	-0.104/0.010
	Part II	-1.115/-0.367	-0.419/0.088	-0.610/-0.056	-0.030/0.118	0.123/0.256
	Part III	-0.700/-0.238	-0.086/0.205	-0.386/-0.154	-0.109/0.135	-0.002/0.266
NQSR(↑)	MESet8	5.942	7.015	6.148	6.966	7.144
	Part I	6.839	7.704	6.967	6.988	7.175
	Part II	6.833	7.489	6.920	7.185	7.346
	Part III	6.814	7.312	6.959	6.900	7.032

images from the multi-exposure images shown in Figure 5, which will be called the “MESet8” set. We also perform the experiments on three datasets called “Part I,” “Part II,” and “Part III” from Wang *et al.*’s paper [52], which include 29, 18, and 30 images respectively.

6.1. Comparisons

We compare our method with the cascades of reverse tone mapping operators (rTMOs) and super-resolution (SR). As the reverse tone mapping operators, a conventional method (Kovaleski *et al.*’s [26]) and a CNN-based method (Eilertsen *et al.*’s [6]) are adopted. For the super-resolution algorithm, one of the recent CNN-based state-of-the-art methods, EDSR [30] is used. For these cascade implementations, we empirically selected the best combination, *i.e.*, we first apply HDRI and then super-resolution because this order gives slightly better quality than its reverse order which is also evidenced in [40]. For our proposed algorithms, two variations are demonstrated: HDRI-SR-B and HDRI-SR-C.

In summary, the following methods and input (LDR-LR) are compared.

- Kovaleski-SR: HDRI by [26] followed by EDSR [30].
- Eilertsen-SR: HDRI by [6] followed by EDSR [30].
- HDRI-SR-B: “Basic model,” trained with only \mathcal{L}_{recon} .
- HDRI-SR-C: “Complex model,” trained with \mathcal{L} based on the generative scheme.

6.2. Metrics

Since there is no reference image for the objective evaluation, we adopt three widely used NR-IQAs (no-reference image quality assessments). Specifically, we adopt natural image quality evaluator (NIQE) [37], HDR image gradient-based evaluator (HIGRADE) [27], and no-reference quality metric for single-image super-resolution (NQSR) [31].

6.3. Quantitative Measurements

Table 1 shows the overall quantitative measures for all the comparisons. As expected, the input LDR-LR shows the worst result in all metrics. In the case of NIQE [37] which is devised to reflect the naturalness, our complex model HDRI-SR-C shows the best results. Also, our proposed basic model HDRI-SR-B shows comparable results to the others. The HIGRADE is designed to measure the quality of tonemapped images [27], where two different measures HIGRADE-1 and HIGRADE-2 are defined which differ in the features they use. For these two measures, the proposed HDRI-SR-C achieves the best results, and HDRI-SR-B shows competitive results with Kovaleski-SR. In the case of NQSR [31] which is designed for the quality measure of super-resolved image, Kovaleski-SR demonstrates the best result on most sets while our HDRI-SR-C achieves the second best. The HDRI-SR-C shows the best result on MESet8. As shown, our HDRI-SR-C model shows considerable super-resolution ability comparable to cascades of the HDRI and state-of-the-art EDSR.

About the computational complexity of compared methods, our CNN for the joint HDRI and SR needs 8 M parameters, Eilertsen *et al.*’s [6] HDRI requires about 32 M, and

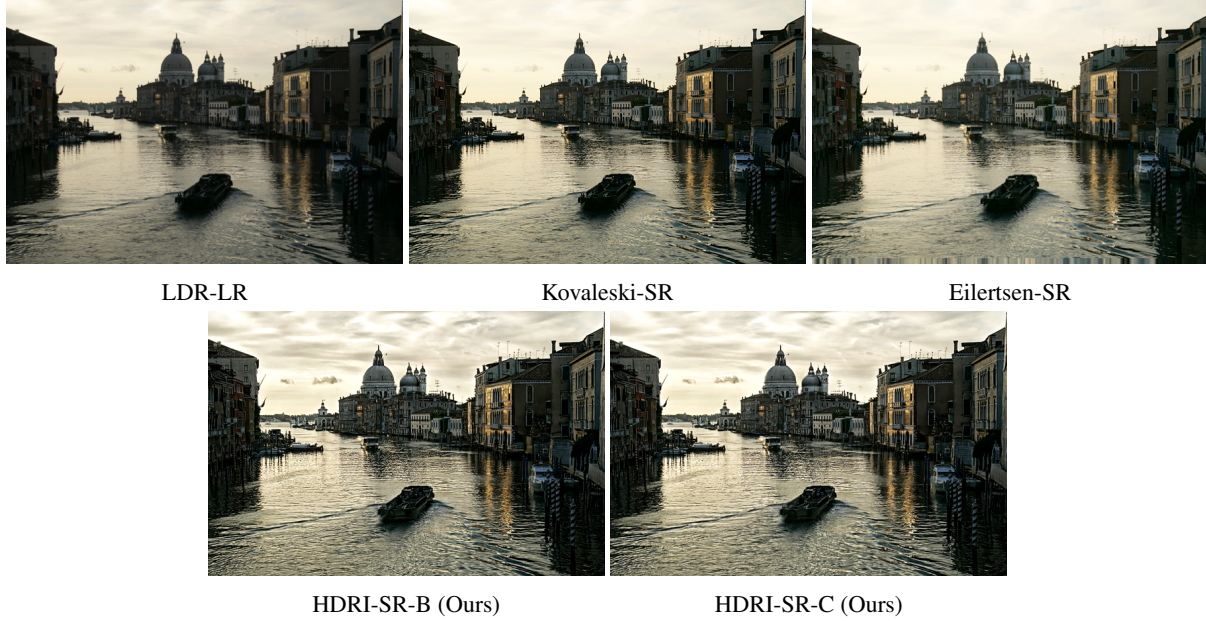


Figure 6: Visualized comparisons on “Grandcanal” of MESet8.

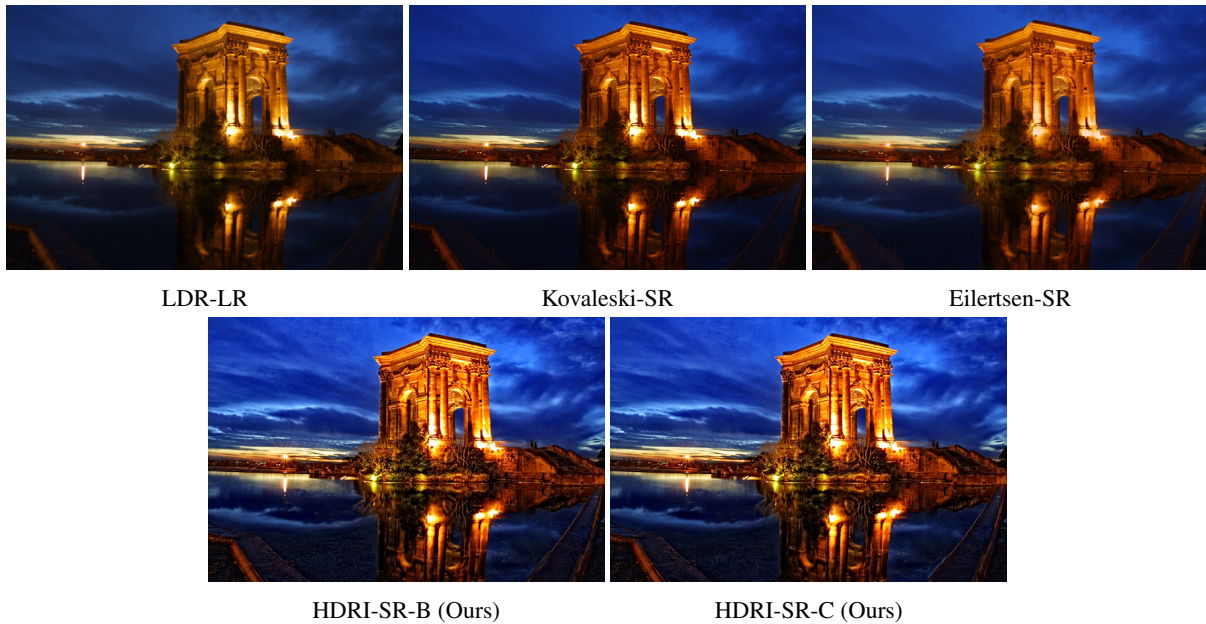


Figure 7: Visualized comparisons on “Peyrou” of MESet8.

EDSR [30] needs about 43 M parameters. Hence, the proposed method needs the least amount of CNN parameters among the compared ones while achieving comparable or better results as shown above.

6.4. Visualized Evaluations

For the qualitative evaluation, we visualize the result images by tonemapping in Figures 6, 7, and 8. It can be seen that ours show abundant textures and details while preserving natural tones in all figures. Also, the color and the global contrast are well enhanced with our algorithm.

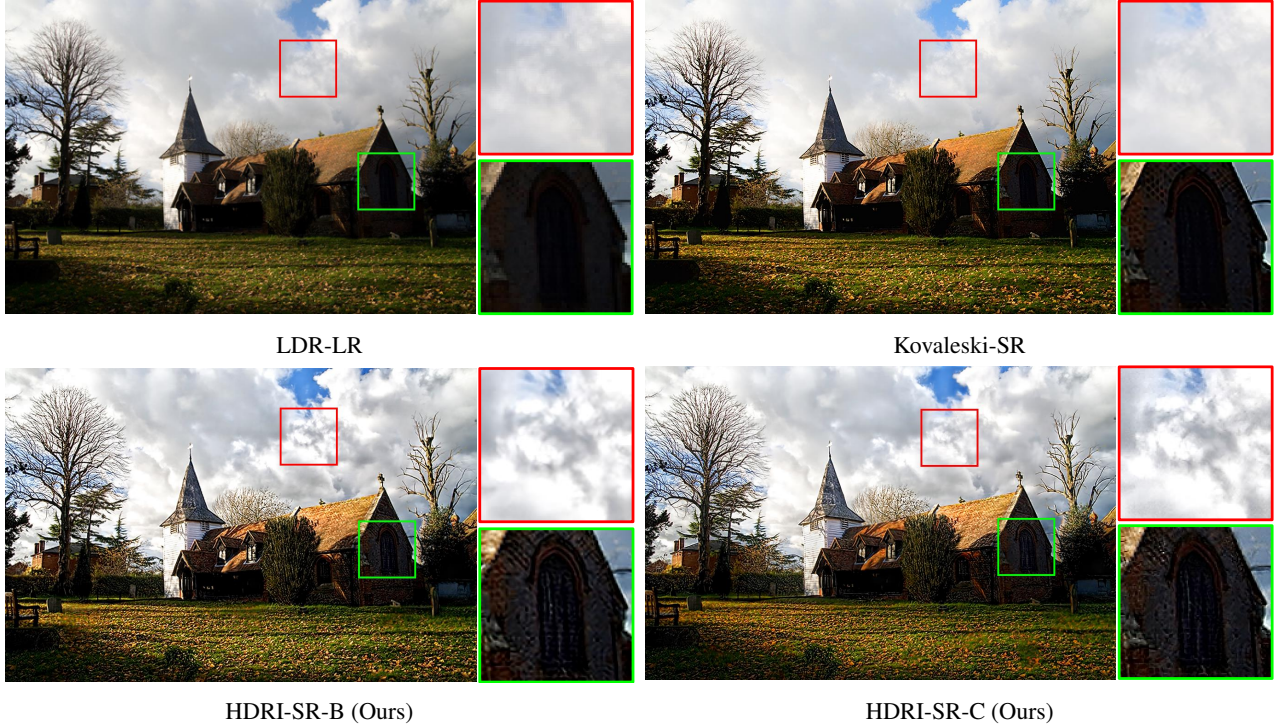


Figure 8: Visualized comparisons on “House” of MESet8.

Specifically in Figure 6, by comparing the building facades, we can see that the high frequency details are enhanced with our HDRI-SRs. Also, the texture and the volume of water flow are much enhanced compared to other algorithms. Additionally, the volume of the cloud became much realistic. In Figure 7, we also visualized “Peyrou” image of MESet8. As shown, the overall contrast is enhanced within the sky and the lake. Also, the texture and details of building and the trees are enriched with ours. Additional zoomed result is shown in Figure 1 where it is shown that the high-frequency details and edges are enhanced and the saturated details of trees are generated. In Figure 8, visualized results of “House” of MESet8 are shown. For the red boxes, our algorithms show much thicker clouds than the cascade of Kovaleski *et al.*'s [26] and the EDSR [30]. As HDRI-SR-C adopts the generative scheme, it generates much better cloud details compared to HDRI-SR-B. By comparing the green box regions, we can see that the diminished details due to low illumination are enhanced with our algorithms.

7. Conclusion

In this paper, we have proposed a CNN-based method for the joint HDRI and SR from a single image, where the domain knowledges from the existing HDRI and SR methods

are exploited in designing the framework. We also considered the issue that there is no ground-truth image for the HDRI. Specifically, we found that the key to the joint task is the reconstruction of high-frequency details. We have also found that we can get rid of inconsistent characteristics of various HDR dataset by removing the illumination. Hence, we decompose the image into the illumination and reflectance, and process the reflectance by the CNN while we just bicubic interpolate and stretch the illumination. The final output is generated by synthesizing the processed components and then linearly stretching the synthesized image according to the target display luminance. We have also proposed a generative approach and training strategies for the joint HDRI-SR task. Experiments show that the proposed methods yield better performance than the cascade of conventional CNN-based HDRI and SR.

References

- [1] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al. Tensorflow: A system for large-scale machine learning. In *OSDI*, volume 16, pages 265–283, 2016. 5
- [2] J. An, S. H. Lee, J. G. Kuk, and N. I. Cho. A multi-exposure image fusion algorithm without ghost effect. In *Acoustics, Speech and Signal Processing (ICASSP), 2011 IEEE International Conference on*, pages 1565–1568. IEEE, 2011. 1

- [3] F. Banterle, K. Debattista, A. Artusi, S. Pattanaik, K. Myszkowski, P. Ledda, and A. Chalmers. High dynamic range imaging and low dynamic range expansion for generating hdr content. In *Computer graphics forum*, volume 28, pages 2343–2367. Wiley Online Library, 2009. 1, 3
- [4] C. Dong, Y. Deng, C. Change Loy, and X. Tang. Compression artifacts reduction by a deep convolutional network. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 576–584, 2015. 2
- [5] C. Dong, C. C. Loy, K. He, and X. Tang. Learning a deep convolutional network for image super-resolution. In *European Conference on Computer Vision*, pages 184–199. Springer, 2014. 1, 2
- [6] G. Eilertsen, J. Kronander, G. Denes, R. K. Mantiuk, and J. Unger. Hdr image reconstruction from a single exposure using deep cnns. *ACM Transactions on Graphics (TOG)*, 36(6):178, 2017. 1, 2, 3, 6
- [7] Y. Endo, Y. Kanamori, and J. Mitani. Deep reverse tone mapping. *ACM Trans. Graph*, 36(6), 2017. 1, 2, 3
- [8] X. Fu, D. Zeng, Y. Huang, Y. Liao, X. Ding, and J. Paisley. A fusion-based enhancing method for weakly illuminated images. *Signal Processing*, 129:82–96, 2016. 2, 3
- [9] X. Gao, K. Zhang, D. Tao, and X. Li. Image super-resolution with sparse neighbor embedding. *IEEE Transactions on Image Processing*, 21(7):3194–3205, 2012. 1
- [10] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. 1, 2, 4
- [11] M. D. Grossberg and S. K. Nayar. What is the space of camera response functions? In *Computer Vision and Pattern Recognition, 2003. Proceedings. 2003 IEEE Computer Society Conference on*, volume 2, pages II–602. IEEE, 2003. 5
- [12] J. Guo and H. Chao. Building dual-domain representations for compression artifacts reduction. In *European Conference on Computer Vision*, pages 628–644. Springer, 2016. 2
- [13] T. Guo, H. S. Mousavi, T. H. Vu, and V. Monga. Deep wavelet prediction for image super-resolution. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2017. 2
- [14] X. Guo, Y. Li, and H. Ling. Lime: Low-light image enhancement via illumination map estimation. *IEEE Transactions on Image Processing*, 26(2):982–993, 2017. 2, 3
- [15] W. Han, S. Chang, D. Liu, M. Yu, M. Witbrock, and T. S. Huang. Image super-resolution via dual-state recurrent networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018. 1, 2
- [16] J.-B. Huang, A. Singh, and N. Ahuja. Single image super-resolution from transformed self-exemplars. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5197–5206, 2015. 1
- [17] Z. Hui, X. Wang, and X. Gao. Fast and accurate single image super-resolution via information distillation network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 723–731, 2018. 1
- [18] Y. Huo, F. Yang, L. Dong, and V. Brost. Physiological inverse tone mapping based on retina response. *The Visual Computer*, 30(5):507–517, 2014. 3
- [19] D. J. Jobson, Z.-u. Rahman, and G. A. Woodell. Properties and performance of a center/surround retinex. *IEEE transactions on image processing*, 6(3):451–462, 1997. 2, 3
- [20] A. Jolicoeur-Martineau. The relativistic discriminator: a key element missing from standard gan. *arXiv preprint arXiv:1807.00734*, 2018. 4
- [21] N. K. Kalantari and R. Ramamoorthi. Deep high dynamic range imaging of dynamic scenes. *ACM Trans. Graph*, 36(4):144, 2017. 1, 2, 3
- [22] J. Kim, J. Kwon Lee, and K. Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1646–1654, 2016. 1, 2
- [23] J. Kim, J. Kwon Lee, and K. Mu Lee. Deeply-recursive convolutional network for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1637–1645, 2016. 2
- [24] R. Kimmel, M. Elad, D. Shaked, R. Keshet, and I. Sobel. A variational framework for retinex. *International Journal of computer vision*, 52(1):7–23, 2003. 2, 3
- [25] P. Korshunov, H. Nemoto, A. Skodras, and T. Ebrahimi. Crowdsourcing-based evaluation of privacy in hdr images. In *Optics, Photonics, and Digital Technologies for Multimedia Applications III*, volume 9138, page 913802. International Society for Optics and Photonics, 2014. 5
- [26] R. P. Kovaleski and M. M. Oliveira. High-quality reverse tone mapping for a wide range of exposures. In *Graphics, Patterns and Images (SIBGRAPI), 2014 27th SIBGRAPI Conference on*, pages 49–56. IEEE, 2014. 1, 3, 6, 8
- [27] D. Kundu, D. Ghadiyaram, A. C. Bovik, and B. L. Evans. No-reference quality assessment of tone-mapped hdr pictures. *IEEE Transactions on Image Processing*, 26(6):2957–2971, 2017. 6
- [28] W.-S. Lai, J.-B. Huang, N. Ahuja, and M.-H. Yang. Deep laplacian pyramid networks for fast and accurate superresolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017. 1, 2
- [29] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4681–4690, 2017. 1, 2, 4
- [30] B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee. Enhanced deep residual networks for single image super-resolution. In *The IEEE conference on computer vision and pattern recognition (CVPR) workshops*, 2017. 1, 2, 6, 7, 8
- [31] C. Ma, C.-Y. Yang, X. Yang, and M.-H. Yang. Learning a no-reference quality metric for single-image super-resolution. *Computer Vision and Image Understanding*, 158:1–16, 2017. 6
- [32] K. Ma, Z. Shu, X. Bai, J. Wang, and D. Samaras. Docunet: Document image unwarping via a stacked u-net. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4700–4709, 2018. 4

- [33] B. Masia, S. Agustin, R. W. Fleming, O. Sorkine, and D. Gutierrez. Evaluation of reverse tone mapping through varying exposure conditions. *ACM transactions on graphics (TOG)*, 28(5):160, 2009. 1, 3
- [34] B. Masia, A. Serrano, and D. Gutierrez. Dynamic range expansion based on image statistics. *Multimedia Tools and Applications*, 76(1):631–648, 2017. 3
- [35] T. Mertens, J. Kautz, and F. Van Reeth. Exposure fusion: A simple and practical alternative to high dynamic range photography. In *Computer Graphics Forum*, volume 28, pages 161–171. Wiley Online Library, 2009. 1
- [36] S. Miller. 2017 update on high dynamic range television. *SMPTE Motion Imaging Journal*, 126(7):94–96, 2017. 1
- [37] A. Mittal, R. Soundararajan, and A. C. Bovik. Making a completely blind image quality analyzer. *IEEE Signal Processing Letters*, 20(3):209–212, 2013. 6
- [38] S. K. Nayar and T. Mitsunaga. High dynamic range imaging: Spatially varying pixel exposures. In *Computer Vision and Pattern Recognition, 2000. Proceedings. IEEE Conference on*, volume 1, pages 472–479. IEEE, 2000. 1
- [39] A. Newell, K. Yang, and J. Deng. Stacked hourglass networks for human pose estimation. In *European Conference on Computer Vision*, pages 483–499. Springer, 2016. 4
- [40] J. S. Park, J. W. Soh, and N. I. Cho. High dynamic range and super-resolution imaging from a single image. *IEEE Access*, 2018. 2, 3, 6
- [41] E. Reinhard, M. Stark, P. Shirley, and J. Ferwerda. Photographic tone reproduction for digital images. *ACM transactions on graphics (TOG)*, 21(3):267–276, 2002. 4
- [42] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 4
- [43] M. S. Sajjadi, B. Schölkopf, and M. Hirsch. Enhancenet: Single image super-resolution through automated texture synthesis. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 4501–4510. IEEE, 2017. 1, 4
- [44] F. Schubert, K. Schertler, and K. Mikolajczyk. A hands-on approach to high-dynamic-range and superresolution fusion. In *Applications of Computer Vision (WACV), 2009 Workshop on*, pages 1–8. IEEE, 2009. 2
- [45] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1874–1883, 2016. 2
- [46] J. Sun, Z. Xu, and H.-Y. Shum. Image super-resolution using gradient profile prior. In *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, pages 1–8. IEEE, 2008. 1
- [47] Y. Tai, J. Yang, and X. Liu. Image super-resolution via deep recursive residual network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017. 2
- [48] Y.-W. Tai, S. Liu, M. S. Brown, and S. Lin. Super resolution using edge prior and single image detail synthesis. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 2400–2407. IEEE, 2010. 1
- [49] J. Tang, E. Peli, and S. Acton. Image enhancement using a contrast measure in the compressed domain. *IEEE Signal Processing Letters*, 10(10):289–292, 2003. 2, 3
- [50] R. Timofte, V. De Smet, and L. Van Gool. A+: Adjusted anchored neighborhood regression for fast super-resolution. In *Asian Conference on Computer Vision*, pages 111–126. Springer, 2014. 1
- [51] T. Tong, G. Li, X. Liu, and Q. Gao. Image super-resolution using dense skip connections. In *2017 IEEE international conference on computer vision*. IEEE, 2017. 1
- [52] S. Wang, J. Zheng, H.-M. Hu, and B. Li. Naturalness preserved enhancement algorithm for non-uniform illumination images. *IEEE Transactions on Image Processing*, 22(9):3538–3548, 2013. 2, 6
- [53] T.-H. Wang, C.-W. Chiu, W.-C. Wu, J.-W. Wang, C.-Y. Lin, C.-T. Chiu, and J.-J. Liou. Pseudo-multiple-exposure-based tone fusion with local region adjustment. *IEEE Transactions on Multimedia*, 17(4):470–484, 2015. 2, 3
- [54] X. Wang, K. Yu, C. Dong, and C. C. Loy. Recovering realistic texture in image super-resolution by deep spatial feature transform. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018. 1, 4
- [55] C.-Y. Yang, C. Ma, and M.-H. Yang. Single-image super-resolution: A benchmark. In *European Conference on Computer Vision*, pages 372–386. Springer, 2014. 1
- [56] J. Yang, Q. Liu, and K. Zhang. Stacked hourglass network for robust facial landmark localisation. In *Computer Vision and Pattern Recognition Workshops (CVPRW), 2017 IEEE Conference on*, pages 2025–2033. IEEE, 2017. 4
- [57] J. Yang, Z. Wang, Z. Lin, S. Cohen, and T. Huang. Coupled dictionary training for image super-resolution. *IEEE transactions on image processing*, 21(8):3467–3478, 2012. 1
- [58] J. Yang, J. Wright, T. S. Huang, and Y. Ma. Image super-resolution via sparse representation. *IEEE transactions on image processing*, 19(11):2861–2873, 2010. 1
- [59] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE Transactions on Image Processing*, 26(7):3142–3155, 2017. 2
- [60] Y. Zhang, Y. Tian, Y. Kong, B. Zhong, and Y. Fu. Residual dense network for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018. 1, 2

Joint High Dynamic Range Imaging and Super-Resolution from a Single Image - Supplementary Material

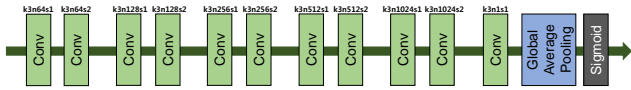


Figure 9: The discriminator architecture for training HDRI-SR-C model. On each of the convolution layers, the k , n , and s denote kernel size, the number of features, and stride of convolution operation, respectively.

1. Discriminator Architecture

Figure 9 shows the discriminator architecture of our HDRI-SR-C model. We stack several convolution layers and employ global average pooling at the latter part of the network instead of fully-connected layers. Also, we alternate max-pooling operations by 2-strided convolution layers to decrease the spatial scale of features.