

# Loss Functions for Image Restoration with Neural Networks

Hang Zhao<sup>\*,†</sup>, Orazio Gallo<sup>\*</sup>, Iuri Frosio<sup>\*</sup>, and Jan Kautz<sup>\*</sup>

**Abstract**—Neural networks are becoming central in several areas of computer vision and image processing and different architectures have been proposed to solve specific problems. The impact of the loss layer of neural networks, however, has not received much attention in the context of image processing: the default and virtually only choice is  $\ell_2$ . In this paper, we bring attention to alternative choices for image restoration. In particular, we show the importance of perceptually-motivated losses when the resulting image is to be evaluated by a human observer. We compare the performance of several losses, and propose a novel, differentiable error function. We show that the quality of the results improves significantly with better loss functions, even when the network architecture is left unchanged.

**Index Terms**—Image Processing, Image Restoration, Neural Networks, Loss Functions.

## I. INTRODUCTION

FOR decades, neural networks have shown various degrees of success in several fields, ranging from robotics, to regression analysis, to pattern recognition. Despite the promising results already produced in the 1980s on handwritten digit recognition [1], the popularity of neural networks in the field of computer vision has grown exponentially only recently, when deep learning boosted their performance in image recognition [2].

In the span of just a couple of years, neural networks have been employed for virtually every computer vision and image processing task known to the research community. Much research has focused on the definition of new architectures that are better suited to a specific problem [3], [4]. A large effort was also made to understand the inner mechanisms of neural networks, and what their intrinsic limitations are, for instance through the development of deconvolutional networks [5], or trying to fool networks with specific inputs [6]. Other advances were made on the techniques to improve the network’s convergence [7].

The loss layer, despite being the effective driver of the network’s learning, has attracted little attention within the image processing research community: the choice of the cost function generally defaults to the squared  $\ell_2$  norm of the error [8], [3], [9], [10]. This is understandable, given the many desirable properties this norm possesses. There is also a less well-founded, but just as relevant reason for the continued

popularity of  $\ell_2$ : standard neural networks packages, such as Caffe [11], only offer the implementation for this metric.

However,  $\ell_2$  suffers from well-known limitations. For instance, when the task at hand involves image quality,  $\ell_2$  correlates poorly with image quality as perceived by a human observer [12]. This is because of a number of assumptions implicitly made when using  $\ell_2$ . First and foremost, the use of  $\ell_2$  assumes that the impact of noise is independent of the local characteristics of the image. On the contrary, the sensitivity of the Human Visual System (HVS) to noise depends on local luminance, contrast, and structure [13]. The  $\ell_2$  loss also works under the assumption of white Gaussian noise, which is not valid in general.

We focus on the use of neural networks for image restoration tasks, and we study the effect of different metrics for the network’s loss layer. We compare  $\ell_2$  against four error metrics on representative tasks: image super-resolution, JPEG artifacts removal, and joint denoising plus demosaicking. First, we test whether a different local metric such as  $\ell_1$  can produce better results. We then evaluate the impact of perceptually-motivated metrics. We use two state-of-the-art metrics for image quality: the structural similarity index (SSIM [13]) and the multi-scale structural similarity index (MS-SSIM [14]). We choose these among the plethora of existing indexes, because they are established measures, and because they are differentiable—a requirement for the backpropagation stage. As expected, on the use cases we consider, the perceptual metrics outperform  $\ell_2$ . However, and perhaps surprisingly, this is also true for  $\ell_1$ , see Figure 1. Inspired by this observation, we propose a novel loss function and show its superior performance in terms of all the metrics we consider.

We offer several contributions. First we bring attention to the importance of the error metric used to train neural networks for image processing: despite the widely-known limitations of  $\ell_2$ , this loss is still the *de facto* standard. We investigate the use of three alternative error metrics ( $\ell_1$ , SSIM, and MS-SSIM), and define a new metric that combines the advantages of  $\ell_1$  and MS-SSIM (Section III). We also perform a thorough analysis of their performance in terms of several image quality indexes (Section IV). Finally, we discuss their convergence properties. We empirically show that the poor performance of some losses is related to local minima of the loss functions (Section V-A), and we explain the reasons why SSIM and MS-SSIM alone do not produce the expected quality (Section V-B). For each of the metrics we analyze, we implement a loss layer for Caffe, which we make available to the research community<sup>1</sup>.

<sup>\*</sup>NVIDIA, <sup>†</sup>MIT Media Lab

The supplementary material is available at <http://research.nvidia.com/publication/loss-functions-image-restoration-neural-networks>. The material includes additional numerical and visual comparisons, and mathematical derivations. Contact ogallo@nvidia.com for further questions about this work.

<sup>1</sup><https://github.com/NVlabs/PLANN>

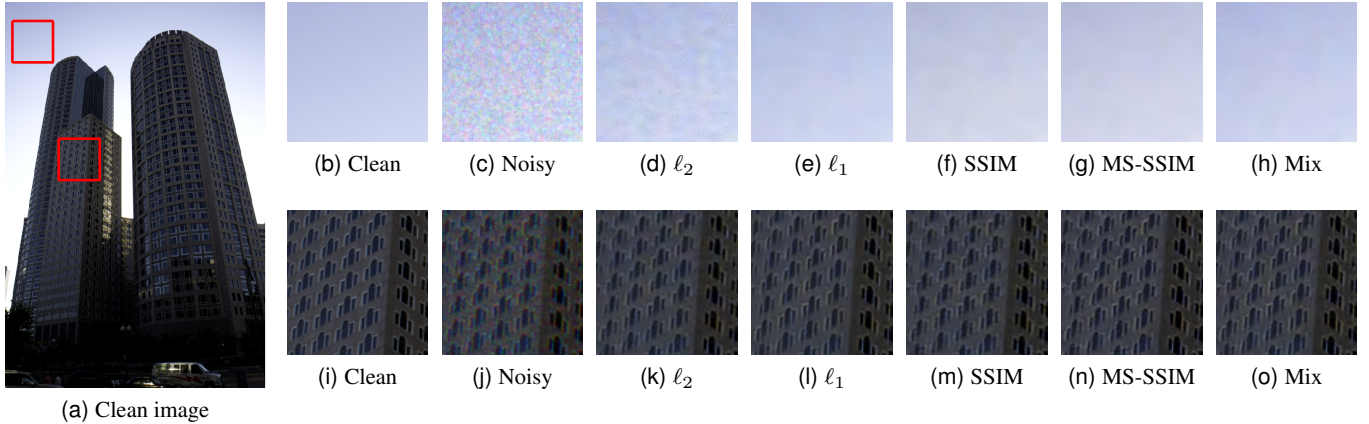


Fig. 1: Comparisons of the results of joint denoising and demosaicking performed by networks trained on different loss functions (best viewed in the electronic version by zooming in).  $\ell_2$ , the standard loss function for neural networks for image processing, produces splotchy artifacts in flat regions (d).

## II. RELATED WORK

In this paper, we target neural networks for image restoration, which, in our context, is the set of all the image processing algorithms whose goal is to output an image that is appealing to a human observer. To this end, we use the problems of super-resolution, JPEG artifacts removal, and joint demosaicking plus denoising as benchmark tests. Specifically, we show how established error measures can be adapted to work within the loss layer of a neural network, and how this can positively influence the results. Here we briefly review the existing literature on the subject of neural networks for image processing, and on the subject of measures of image quality.

### A. Neural networks for image restoration

Following their success in several computer vision tasks [15], [2], neural networks have received considerable attention in the context of image restoration. Neural networks have been used for denoising [8], [3], deblurring [4], demosaicking [10], and super-resolution [9] among others. To the best of our knowledge, however, the work on this subject has focused on tuning the architecture of the network for the specific application; the loss layer, which effectively drives the learning of the network to produce the desired output quality, is based on  $\ell_2$  for all of the approaches above.

We show that a better choice for the error measure has a strong impact on the quality of the results. Moreover, we show that even when  $\ell_2$  is the appropriate loss, alternating the training loss function with a related loss, such as  $\ell_1$ , can lead to finding a better solution for  $\ell_2$ .

### B. Evaluating image quality

The mean squared error,  $\ell_2$ , is arguably the dominant error measure across very diverse fields, from regression problems, to pattern recognition, to signal and image processing. Among the main reasons for its popularity is the fact that it is convex and differentiable—very convenient properties for optimization problems. Other interesting properties range from the fact that  $\ell_2$  provides the maximum likelihood estimate in case of

independent and identically distributed Gaussian noise, to the fact that it is additive for independent noise sources. There is longer list of reasons for which we refer the reader to the work of Wang and Bovik [16].

These properties paved the way for  $\ell_2$ 's widespread adoption, which was further fueled by the fact that standard software packages tend to include tools to use  $\ell_2$ , but not many other error functions for regression. In the context of image processing, Caffe [11] actually offers *only*  $\ell_2$  as a loss layer<sup>2</sup>, thus discouraging researchers from testing other error measures.

However, it is widely accepted that  $\ell_2$ , and consequently the Peak Signal-to-Noise Ratio, PSNR, do not correlate well with human's perception of image quality [12]:  $\ell_2$  simply does not capture the intricate characteristics of the human visual system (HVS).

There exists a rich literature of error measures, both reference-based and non reference-based, that attempt to address the limitations of the simple  $\ell_2$  error function. For our purposes, we focus on reference-based measures. A popular reference-based index is the structural similarity index (SSIM [13]). SSIM evaluates images accounting for the fact that the HVS is sensitive to changes in local structure. Wang et al. [14] extend SSIM observing that the scale at which local structure should be analyzed is a function of factors such as image-to-observer distance. To account for these factors, they propose MS-SSIM, a multi-scale version of SSIM that weighs SSIM computed at different scales according to the sensitivity of the HVS. Experimental results have shown the superiority of SSIM-based indexes over  $\ell_2$ . As a consequence, SSIM has been widely employed as a metric to evaluate image processing algorithms. Moreover, given that it can be used as a differentiable cost function, SSIM has also been used in iterative algorithms designed for image compression [16], image reconstruction [17], denoising and super-resolution [18], and even downscaling [19]. To the best of our knowledge,

<sup>2</sup>Caffe indeed offers other types of loss layers, but they are only useful for classification tasks.

however, SSIM-based indexes have never been adopted to train neural networks.

Recently, novel image quality indexes based on the properties of the HVS showed improved performance when compared to SSIM and MS-SSIM [12]. One of these is the Information Weighted SSIM (IW-SSIM), a modification of MS-SSIM that also includes a weighting scheme proportional to the local image information [20]. Another is the Visual Information Fidelity (VIF), which is based on the amount of shared information between the reference and distorted image [21]. The Gradient Magnitude Similarity Deviation (GMSD) is characterized by simplified math and performance similar to that of SSIM, but it requires computing the standard deviation over the whole image [22]. Finally, the Feature Similarity Index (FSIM), leverages the perceptual importance of phase congruency, and measures the dissimilarity between two images based on local phase congruency and gradient magnitude [23]. FSIM has also been extended to FSIM<sub>c</sub>, which can be used with color images. Despite the fact that they offer an improved accuracy in terms of image quality, the mathematical formulation of these indexes is generally more complex than SSIM and MS-SSIM, and possibly not differentiable, making their adoption for optimization procedures not immediate.

### III. LOSS LAYERS FOR IMAGE RESTORATION

The loss layer of a neural network compares the output of the network with the ground truth, i.e., processed and reference patches, respectively, for the case of image processing.

In our work, we investigate the impact of different loss function layers for image processing. Consider the case of a network that performs denoising and demosaicking jointly. The insets in Figure 1 show a zoom-in of different patches for the image in Figure 1(a) as processed by a network trained with different loss functions (see Section IV for the network’s description). A simple visual inspection is sufficient to appreciate the practical implications of the discussion on  $\ell_2$  (Section II).

Specifically, Figure 1(d) shows that in flat regions the network strongly attenuates the noise, but it produces visible splotchy artifacts. This is because  $\ell_2$  penalizes larger errors, but is more tolerant to small errors, regardless of the underlying structure in the image; the HVS, on the other hand, is more sensitive to luminance and color variations in texture-less regions [24]. A few splotchy artifacts are still visible, though arguably less apparent, in textured regions, see Figure 1(k). The sharpness of edges, however, is well-preserved by  $\ell_2$ , as blurring them would result in a large error. Note that these splotchy artifacts have been systematically observed before in the context of image processing with neural networks [3], but they have not been attributed to the loss function. In Section V-A we show that the quality achieved by using  $\ell_2$  is also dependent on its convergence properties.

In this section we propose the use of different error functions. We provide a motivation for the different loss functions and we show how to compute their derivatives, which are necessary for the backpropagation step. We also share our

implementation of the different layers that can be readily used within Caffe.

For an error function  $\mathcal{E}$ , the loss for a patch  $P$  can be written as

$$\mathcal{L}^{\mathcal{E}}(P) = \frac{1}{N} \sum_{p \in P} \mathcal{E}(p), \quad (1)$$

where  $N$  is the number of pixels  $p$  in the patch.

#### A. The $\ell_1$ error

As a first attempt to reduce the artifacts introduced by the  $\ell_2$  loss function, we want to train the exact same network using  $\ell_1$  instead of  $\ell_2$ . The two losses weigh errors differently— $\ell_1$  does not over-penalize larger errors—and, consequently, they may have different convergence properties.

Equation 1 for  $\ell_1$  is simply:

$$\mathcal{L}^{\ell_1}(P) = \frac{1}{N} \sum_{p \in P} |x(p) - y(p)|, \quad (2)$$

where  $p$  is the index of the pixel and  $P$  is the patch;  $x(p)$  and  $y(p)$  are the values of the pixels in the processed patch and the ground truth respectively. The derivatives for the back-propagation are also simple, since  $\partial \mathcal{L}^{\ell_1}(p) / \partial q = 0, \forall q \neq p$ . Therefore, for each pixel  $p$  in the patch,

$$\partial \mathcal{L}^{\ell_1}(P) / \partial x(p) = \text{sign}(x(p) - y(p)). \quad (3)$$

The derivative of  $\mathcal{L}^{\ell_1}$  is not defined at 0. However, if the error is 0, we do not need to update the weights, so we use the convention that  $\text{sign}(0) = 0$ . Note that, although  $\mathcal{L}^{\ell_1}(P)$  is a function of the patch as a whole, the derivatives are back-propagated for each pixel in the patch. Somewhat unexpectedly, the network trained with  $\ell_1$  provides a significant improvement for several of the issues discussed above, see Figure 1(e) where the splotchy artifacts in the sky are removed. Section V-A analyzes the reasons behind this behavior. Although  $\ell_1$  shows improved performance over  $\ell_2$ , however, its results are still sub-optimal (see for instance the artifacts in the sky and at the boundary of the building in Figure 3(c) and 3(g)).

#### B. SSIM

If the goal is for the network to learn to produce visually pleasing images, it stands to reason that the error function should be perceptually motivated, as is the case with SSIM.

SSIM for pixel  $p$  is defined as

$$\begin{aligned} \text{SSIM}(p) &= \frac{2\mu_x\mu_y + C_1}{\mu_x^2 + \mu_y^2 + C_1} \cdot \frac{2\sigma_{xy} + C_2}{\sigma_x^2 + \sigma_y^2 + C_2} \quad (4) \\ &= l(p) \cdot cs(p) \quad (5) \end{aligned}$$

where we omitted the dependence of means and standard deviations on pixel  $p$ . Means and standard deviations are computed with a Gaussian filter with standard deviation  $\sigma_G$ ,  $G_{\sigma_G}$ . The loss function for SSIM can be then written setting  $\mathcal{E}(p) = 1 - \text{SSIM}(p)$ :

$$\mathcal{L}^{\text{SSIM}}(P) = \frac{1}{N} \sum_{p \in P} 1 - \text{SSIM}(p). \quad (6)$$

Equation 4 highlights the fact that the computation of  $\text{SSIM}(p)$  requires looking at a neighborhood of pixel  $p$  as large as the support of  $G_{\sigma_G}$ . This means that  $\mathcal{L}^{\text{SSIM}}(P)$ , as well as its derivatives, cannot be calculated in some boundary region of  $P$ . This is not true for  $\ell_1$  or  $\ell_2$ , which only need the value of the processed and reference patch at pixel  $p$ .

However, the convolutional nature of the network allows us to write the loss as

$$\mathcal{L}^{\text{SSIM}}(P) = 1 - \text{SSIM}(\tilde{p}), \quad (7)$$

where  $\tilde{p}$  is the center pixel of patch  $P$ . Again, this is because, even though the network learns the weights maximizing SSIM for the central pixel, the learned kernels are then applied to all the pixels in the image. Note that the error can still be back-propagated to all the pixels within the support of  $G_{\sigma_G}$  as they contribute to the computation of Equation 7.

Recall that we have to compute the derivatives at  $\tilde{p}$  with respect to any other pixel  $q$  in patch  $P$ . More formally:

$$\begin{aligned} \frac{\partial \mathcal{L}^{\text{SSIM}}(P)}{\partial x(q)} &= -\frac{\partial}{\partial x(q)} \text{SSIM}(\tilde{p}) \\ &= -\left( \frac{\partial l(\tilde{p})}{\partial x(q)} \cdot cs(\tilde{p}) + l(\tilde{p}) \cdot \frac{\partial cs(\tilde{p})}{\partial x(q)} \right), \end{aligned} \quad (8)$$

where  $l(\tilde{p})$  and  $cs(\tilde{p})$  are the first and second term of SSIM (Equation 5) and their derivatives are

$$\frac{\partial l(\tilde{p})}{\partial x(q)} = 2 \cdot G_{\sigma_G}(q - \tilde{p}) \cdot \left( \frac{\mu_y - \mu_x \cdot l(\tilde{p})}{\mu_x^2 + \mu_y^2 + C_1} \right) \quad (9)$$

and

$$\begin{aligned} \frac{\partial cs(\tilde{p})}{\partial x(q)} &= \frac{2}{\sigma_x^2 + \sigma_y^2 + C_2} \cdot G_{\sigma_G}(q - \tilde{p}) \cdot \left[ (y(q) - \mu_y) \right. \\ &\quad \left. - cs(\tilde{p}) \cdot (x(q) - \mu_x) \right], \end{aligned} \quad (10)$$

where  $G_{\sigma_G}(q - \tilde{p})$  is the Gaussian coefficient associated with pixel  $q$ . We refer the reader to the additional material for the full derivation.

### C. MS-SSIM

The choice of  $\sigma_G$  has an impact on the quality of the processed results of a network that is trained with SSIM, as can be seen from the derivatives in the previous section. Specifically, for smaller values of  $\sigma_G$  the network loses the ability to preserve the local structure and the splotchy artifacts are reintroduced in flat regions, see Figure 9(e). For large values of  $\sigma_G$ , we observe that the network tends to preserve noise in the proximity of edges, Figure 9(c). See Section V-B for more details.

Rather than fine-tuning the  $\sigma_G$ , we propose to use the multi-scale version of SSIM, MS-SSIM. Given a dyadic pyramid of  $M$  levels, MS-SSIM is defined as

$$\text{MS-SSIM}(p) = l_M^\alpha(p) \cdot \prod_{j=1}^M cs_j^{\beta_j}(p) \quad (11)$$

where  $l_M$  and  $cs_j$  are the terms we defined in Section III-B at scale  $M$  and  $j$ , respectively. For convenience, we set  $\alpha = \beta_j = 1$ , for  $j = \{1, \dots, M\}$ . Similarly to Equation 7, we can

approximate the loss for patch  $P$  with the loss computed at its center pixel  $\tilde{p}$ :

$$\mathcal{L}^{\text{MS-SSIM}}(P) = 1 - \text{MS-SSIM}(\tilde{p}). \quad (12)$$

Because we set all the exponents of Equation 11 to one, the derivatives of the loss function based on MS-SSIM can be written as

$$\begin{aligned} \frac{\partial \mathcal{L}^{\text{MS-SSIM}}(P)}{\partial x(q)} &= \left( \frac{\partial l_M(\tilde{p})}{\partial x(q)} + l_M(\tilde{p}) \cdot \sum_{i=0}^M \frac{1}{cs_i(\tilde{p})} \frac{\partial cs_i(\tilde{p})}{\partial x(q)} \right) \cdot \prod_{j=1}^M cs_j(\tilde{p}), \end{aligned} \quad (13)$$

where the derivatives of  $l_j$  and  $cs_j$  are the same as in Section III-B. For the full derivation we refer the reader to the supplementary material.

Using  $\mathcal{L}^{\text{MS-SSIM}}$  to train the network, Equation 11 requires that we compute a pyramid of  $M$  levels of patch  $P$ , which is a computationally expensive operation given that it needs to be performed at each iteration. To avoid this, we propose to approximate and replace the construction of the pyramid: instead of computing  $M$  levels of the pyramid, we use  $M$  different values for  $\sigma_G$ , each one being half of the previous, on the full-resolution patch. Specifically, we use  $\sigma_G^i = \{0.5, 1, 2, 4, 8\}$  and define  $cs_i \triangleq G_{\sigma_G^i} \cdot cs_0(\tilde{p})$  and  $\partial cs_i(\tilde{p})/\partial x(q) \triangleq G_{\sigma_G^i} \cdot \partial cs_0(\tilde{p})/\partial x(q)$ , where the Gaussian filters are centered at pixel  $\tilde{p}$ , and “ $\cdot$ ” is a point-wise multiplication. The terms depending on  $l_M$  can be defined in a similar way. We use this trick to speed up the training in all of our experiments.

### D. The best of both worlds: MS-SSIM + $\ell_1$

By design, both MS-SSIM and SSIM are not particularly sensitive to uniform biases (see Section V-B). This can cause changes of brightness or shifts of colors, which typically become more dull. However, MS-SSIM preserves the contrast in high-frequency regions better than the other loss functions we experimented with. On the other hand,  $\ell_1$  preserves colors and luminance—an error is weighed equally regardless of the local structure—but does not produce quite the same contrast as MS-SSIM.

To capture the best characteristics of both error functions, we propose to combine them:

$$\mathcal{L}^{\text{Mix}} = \alpha \cdot \mathcal{L}^{\text{MS-SSIM}} + (1 - \alpha) \cdot G_{\sigma_G^M} \cdot \mathcal{L}^{\ell_1}, \quad (14)$$

where we omitted the dependence on patch  $P$  for all loss functions, and we empirically set  $\alpha = 0.84$ .<sup>3</sup> The derivatives of Equation 14 are simply the weighed sum of the derivatives of its two terms, which we compute in the previous sections. Note that we add a point-wise multiplication between  $G_{\sigma_G^M}$  and  $\mathcal{L}^{\ell_1}$ : this is because MS-SSIM propagates the error at pixel  $q$  based on its contribution to MS-SSIM of the central pixel  $\tilde{p}$ , as determined by the Gaussian weights, see Equations 9 and 10.

<sup>3</sup>We chose the value for this parameter so that the contribution of the two losses would be roughly balanced. While we tested a few different values, we did not perform a thorough investigation. We did, however, notice that the results were not significantly sensitive to small variations of  $\alpha$ .

## IV. RESULTS

For our analysis of the different loss functions we focus on joint demosaicking plus denoising, a fundamental problem in image processing. We also confirm our findings by testing the different loss functions on the problems of super-resolution and JPEG artifacts removal.

### A. Joint denoising and demosaicking

We define a fully convolutional neural network (CNN) that takes a  $31 \times 31 \times 3$  input. The first layer is a  $64 \times 9 \times 9 \times 3$  convolutional layer, where the first term indicates the number of filters and the remaining terms indicate their dimensions. The second convolutional layer is  $64 \times 5 \times 5 \times 64$ , and the output layer is a  $3 \times 5 \times 5 \times 64$  convolutional layer. We apply parametric rectified linear unit (PReLU) layers to the output of the inner convolutional layers [15]. The input to our network is obtained by bilinear interpolation of a  $31 \times 31$  Bayer patch, which results in a  $31 \times 31 \times 3$  RGB patch; in this sense the network is really doing joint denoising + demosaicking and super-resolution. We trained the network considering different cost functions ( $\ell_1$ ,  $\ell_2$ , SSIM<sub>5</sub>, SSSIM<sub>9</sub>, MS-SSIM and MS-SSIM+ $\ell_1$ )<sup>4</sup> on a training set of 700 RGB images taken from the MIT-Adobe FiveK Dataset [25], resized to a size of  $999 \times 666$ . To simulate a realistic image acquisition process, we corrupted each image with a mix of photon shot and zero-mean Gaussian noise, and introduced clipping due to the sensor zero level and saturation. We used the model proposed by Foi et al. [26] for this task, with parameters  $a = 0.005$  and  $b = 0.0001$ . The average PSNR for the testing images after adding noise was 28.24dB, as reported in Table I. Figure 1(c) shows a typical patch corrupted by noise. We used 40 images from the same dataset for testing (the network did not see this subset during training).

Beyond considering different cost functions for training, we also compare the output of our network with the images obtained by a state-of-the-art denoising method, BM3D, operating directly in the Bayer domain [27], followed by a standard demosaicking algorithm [28]. Since BM3D is designed to deal with Gaussian noise, rather than the more realistic noise model we use, we apply a Variance Stabilizing Transform [26] to the image data in Bayer domain before applying BM3D, and its inverse after denoising. This is exactly the strategy suggested by the authors of BM3D for RAW data [27], the paper we compare against for BM3D.

Figure 1 and 5 show several results and comparisons between the different networks. Note the splotchy artifacts for the  $\ell_2$  network on flat regions, and the noise around the edges for the SSIM<sub>5</sub> and SSIM<sub>9</sub> networks. The network trained with MS-SSIM addresses these problems but tends to render the colors more dull, see Section V-B for an explanation. The network trained with MS-SSIM+ $\ell_1$  generates the best results. Some differences are difficult to see in side-by-side comparisons, please refer to the additional material, which allows to flip between images.

We also perform a quantitative analysis of the results. We evaluate several image quality metrics on the output of

the CNNs trained with the different cost functions and with BM3D [27], which is the state of the art in denoising. The image quality indexes, range from the traditional  $\ell_2$  metric and PSNR, to the most refined, perceptually inspired metrics, like FSIM [23]. The average values of these metrics on the testing dataset are reported in Table I. When the network is trained using  $\ell_1$  as a cost function, instead of the traditional  $\ell_2$ , the average quality of the output images is superior for all the quality metrics considered. It is quite remarkable to notice that, even when the  $\ell_2$  or PSNR metrics are used to evaluate image quality, the network trained with the  $\ell_1$  loss outperforms the one trained with the  $\ell_2$  loss. We offer an explanation of this in Section V-A. On the other hand, we note that the network trained with SSIM performs either at par or slightly worse than the one trained with  $\ell_1$ , both on traditional metrics and on perceptually-inspired losses. The network trained on MS-SSIM performs better than the one based on SSIM, but still does not clearly outperforms  $\ell_1$ . This is due to the color shifts mentioned above and discussed in Section V-B. However, the network that combines MS-SSIM and  $\ell_1$  achieves the best results on all of the image quality metrics we consider.

For this image restoration task, we also evaluated the combination of  $\ell_2$  and MS-SSIM to test whether the use of  $\ell_1$  is critical to the performance of Mix. The results on all the indexes considered in Table I show that Mix is still superior. As expected, however,  $\ell_2$  + MS-SSIM outperforms  $\ell_2$  alone on the perceptual metrics (see supplementary material).

### B. Super-resolution

We also verify the outcome of our analysis on the network for super-resolution proposed by Dong et al. [9]. We start from the network architecture they propose, and make a few minor but important changes to their approach. First we use PReLU, instead of ReLU, layers. Second we use bilinear instead of bicubic interpolation for initialization (we do this for both training and testing). The latter introduces high-frequency artifacts that hurt the learning process. Finally we train directly on the RGB data. We made these changes for all the loss functions we test, including  $\ell_2$ , which also serves as a comparison with the work by Dong et al. We use this modified network in place of their proposed architecture to isolate the contribution of the loss layer to the results—and because it produces better or equal results than the one produced by the original architecture. Figure 2 shows some sample results. Given the results of the previous section, we only compared the results of  $\ell_1$ ,  $\ell_2$ , MS-SSIM, and Mix, see Table I. An analysis of the table brings similar considerations as for the case of joint denoising and demosaicking.

### C. JPEG artifacts removal

Our final benchmark is JPEG artifact removal. We use the same network architecture and ground truth data as for joint denoising and demosaicking, but we create the corrupted data by means of aggressive JPEG compression. For this purpose we use the Matlab's function *imwrite* with a quality setting of 25. Note that the corruption induced by JPEG compression is spatially variant and has a very different statistic than the noise

<sup>4</sup>SSIM<sub>k</sub> means SSIM computed with  $\sigma_G = k$ .



Fig. 2: Results for super-resolution. Notice the grating artifacts on the black stripes of the wing and around the face of the girl produced by  $\ell_2$ .

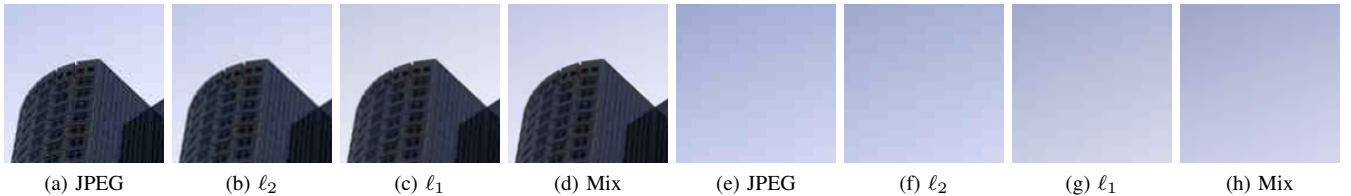


Fig. 3: Results for JPEG de-blocking. The insets are taken from the image in Figure 1. Notice the artifacts around the edges (a)-(c) and how Mix (d) removes them better than either  $\ell_1$  or  $\ell_2$ . Mix also outperforms the other metrics in the relatively flat regions, where the blocketization is more apparent, e.g., (e)-(h).

introduced in Section IV-A. We generate the input patches with two different strides. First, we use a stride of 8, which causes the boundaries of the  $8 \times 8$  DCT blocks to be aligned in each  $31 \times 31$  patch. We also use a stride of 7, which causes the grid of blocks to be in different locations in each patch. We found that the latter strategy better removes the JPEG artifacts while producing sharper images, which is why we ran all of our experiments with that configuration. Again, we only compared the results of  $\ell_1$ ,  $\ell_2$ , MS-SSIM, and Mix, see Table I. Figure 3 shows that our loss function, Mix, outperforms  $\ell_1$  and  $\ell_2$  on uniform regions and that it attenuates the ringing artifacts around the edge of the building better. More results are shown in Figure 6.

More results and comparisons, both numerical and visual, can be found in the supplementary material.

## V. DISCUSSION

In this section we delve into a deeper analysis of the results. Among other considerations, we look into the convergence properties of the different losses, and we offer an interpretation of the reasons some losses perform better than others.

### A. Convergence of the loss functions

Table I highlights an unexpected result: even after convergence, CNNs trained on one loss function can outperform another network even based on the very loss with which the second was trained. Consider, for instance, the two networks trained with  $\ell_1$  and  $\ell_2$  respectively for joint denoising and demosaicking: the table shows that the network trained with  $\ell_1$  achieves a lower  $\ell_2$  loss than then network trained with  $\ell_2$ . Note that we ran the training multiple times with different initializations. We hypothesize that this result may be related to the smoothness and the local convexity properties of the two measures:  $\ell_2$  gets stuck more easily in a local minimum, while for  $\ell_1$  it may be easier to reach better minimum, both in terms

of  $\ell_1$  and  $\ell_2$ —the “good” minima of the two should be related, after all. To test this hypothesis, we ran an experiment in which we take two networks trained with  $\ell_1$  and  $\ell_2$  respectively, and train them again until they converge using the other loss. Figure 7 shows the  $\ell_2$  loss computed on the testing set at different training iterations for either network. The network trained with  $\ell_1$  only (before epoch 1200 in the plot) achieves a better  $\ell_2$  loss than the one trained with  $\ell_2$ . However, after switching the training loss functions, both networks yield a lower  $\ell_2$  loss, confirming that the  $\ell_2$  network was previously stuck in a local minimum. While the two networks achieve a similar  $\ell_2$  loss, they converge to different regions of the space of parameters. At visual inspection the network trained with  $\ell_2$  first and  $\ell_1$  after produces results similar to those of  $\ell_1$  alone; the output of the network trained with  $\ell_1$  first and  $\ell_2$  after is still affected by splotchy artifacts in flat areas, though it is better than  $\ell_2$  alone, see Figure 8. Specifically, the network trained with  $\ell_1$  first and  $\ell_2$  afterwards achieves an  $\ell_2$  loss of  $0.3896 \cdot 10^3$ , which is the lowest across all the networks we trained, confirming that the network trained with  $\ell_2$  alone had convergence issues. However, neither of these two networks outperforms Mix on any of the perceptual metrics we use, confirming that the advantage of the proposed loss goes beyond convergence considerations. Table II shows the complete evaluation of the two networks, where bold indicates a result that is better than any of the results shown in Table I.

### B. On the performance of SSIM and MS-SSIM

Table I also reveals that SSIM and MS-SSIM do not perform as well as  $\ell_1$ . This does not have to do with the convergence properties of these losses. To investigate this, we trained several SSIM networks with different  $\sigma_G$ 's and found that smaller values of  $\sigma_G$  produce better results at edges, but worse results in flat regions, while the opposite is true for larger values, see Figure 9. This can be understood by looking at



Fig. 4: The reference images for the details shown in Figures 5 and 6.

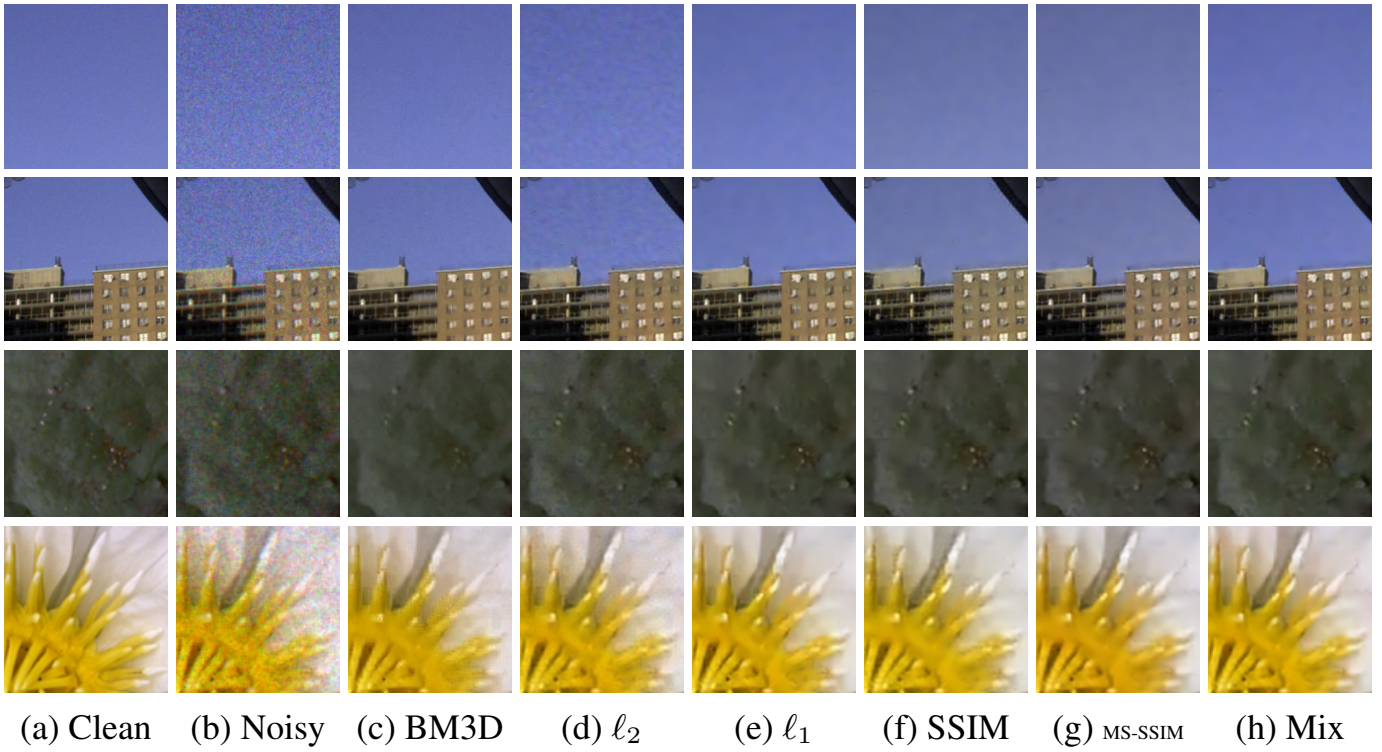


Fig. 5: Results for denoising+deblocking for different approaches. The noisy patches are obtained by simple bilinear interpolation. Note the splotchy artifacts  $l_2$  produces in flat regions. Also note the change in colors for SSIM-based losses. The proposed metric,  $MS\text{-}SSIM+l_1$  referred to as Mix, addresses the former issues. Reference images are shown in Figure 4.

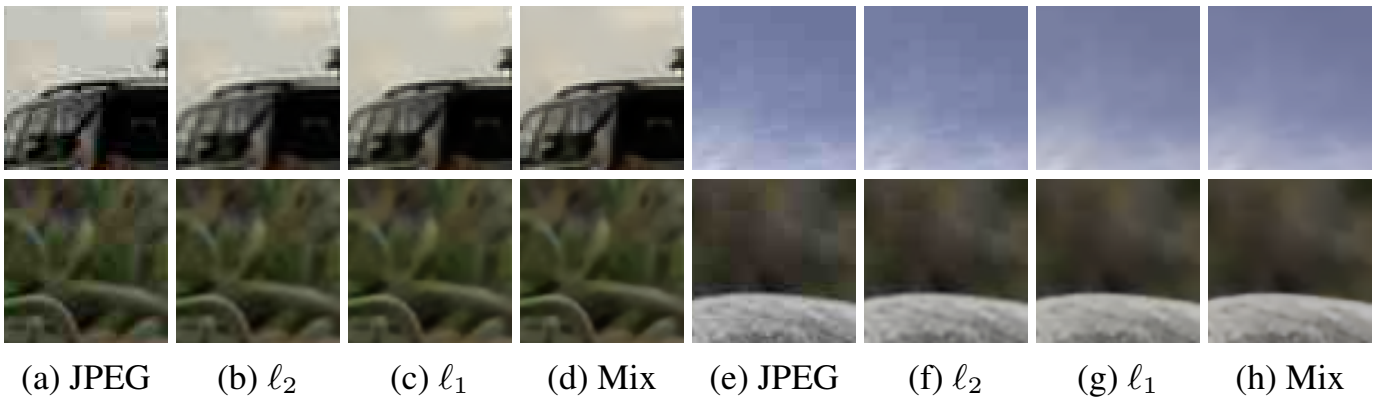


Fig. 6: Results for JPEG de-blocking for different approaches. Note that Mix outperforms both  $l_2$  and  $l_1$  on uniform regions as well as on the ringing artifacts at the in high-gradient regions. Reference images are shown in Figure 4.

Denoising + demosaicking			Training cost function					
Image quality metric	Noisy	<i>BM3D</i>	$\ell_2$	$\ell_1$	SSIM <sub>5</sub>	SSIM <sub>9</sub>	MS-SSIM	Mix
$1000 \cdot \ell_2$	1.65	0.45	0.56	0.43	0.58	0.61	0.55	<b>0.41</b>
PSNR	28.24	34.05	33.18	34.42	33.15	32.98	33.29	<b>34.61</b>
$1000 \cdot \ell_1$	27.36	14.14	15.90	13.47	15.90	16.33	15.99	<b>13.19</b>
SSIM	0.8075	0.9479	0.9346	0.9535	0.9500	0.9495	0.9536	<b>0.9564</b>
MS-SSIM	0.8965	0.9719	0.9636	0.9745	0.9721	0.9718	0.9741	<b>0.9757</b>
IW-SSIM	0.8673	0.9597	0.9473	0.9619	0.9587	0.9582	0.9617	<b>0.9636</b>
GMSD	0.1229	0.0441	0.0490	0.0434	0.0452	0.0467	0.0437	<b>0.0401</b>
FSIM	0.9439	0.9744	0.9716	0.9775	0.9764	0.9759	0.9782	<b>0.9795</b>
FSIM <sub>c</sub>	0.9381	0.9737	0.9706	0.9767	0.9752	0.9746	0.9769	<b>0.9788</b>

Super-resolution		Training cost function			
Image quality metric	Bilinear	$\ell_2$	$\ell_1$	MS-SSIM	Mix
$1000 \cdot \ell_2$	2.5697	1.2407	1.1062	1.3223	<b>1.0990</b>
PSNR	27.16	30.66	31.26	30.11	<b>31.34</b>
$1000 \cdot \ell_1$	28.7764	20.4730	19.0643	22.3968	<b>18.8983</b>
SSIM	0.8632	0.9274	0.9322	0.9290	<b>0.9334</b>
MS-SSIM	0.9603	0.9816	0.9826	0.9817	<b>0.9829</b>
IW-SSIM	0.9532	0.9868	0.9879	0.9866	<b>0.9881</b>
GMSD	0.0714	0.0298	0.0259	0.0316	<b>0.0255</b>
FSIM	0.9070	0.9600	0.9671	0.9601	<b>0.9680</b>
FSIM <sub>c</sub>	0.9064	0.9596	0.9667	0.9597	<b>0.9677</b>

JPEG de-blocking		Training cost function			
Image quality metric	Original JPEG	$\ell_2$	$\ell_1$	MS-SSIM	Mix
$1000 \cdot \ell_2$	0.6463	0.6511	0.6027	1.9262	<b>0.5580</b>
PSNR	32.60	32.73	32.96	27.66	<b>33.25</b>
$1000 \cdot \ell_1$	16.5129	16.2633	16.0687	33.6134	<b>15.5489</b>
SSIM	0.9410	0.9427	0.9467	0.9364	<b>0.9501</b>
MS-SSIM	0.9672	0.9692	0.9714	0.9674	<b>0.9734</b>
IW-SSIM	0.9527	0.9562	0.9591	0.9550	<b>0.9625</b>
GMSD	0.0467	0.0427	0.0413	0.0468	<b>0.0402</b>
FSIM	0.9805	0.9803	0.9825	0.9789	<b>0.9830</b>
FSIM <sub>c</sub>	0.9791	0.9790	0.9809	0.9705	<b>0.9815</b>

TABLE I: Average value of different image quality metrics on the testing dataset for the different cost functions. For SSIM, MS-SSIM, IW-SSIM, GMSD and FSIM the value reported here has been obtained as an average of the three color channels. Best results are shown in bold. (Lower is better for  $\ell_1$ ,  $\ell_2$ , and GMSD, higher is better for the others.)

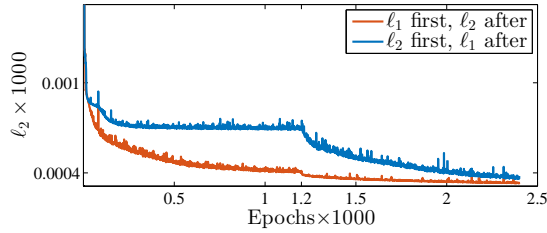


Fig. 7:  $\ell_2$  loss on the testing set for the two networks that switch loss functions during training.

Figure 10(a), where the size of the support for the computation of SSIM for the three values of  $\sigma_G$  is shown for a pixel P close to an edge. A larger  $\sigma_G$  is more tolerant to the same amount of noise because it detects the presence of the edge, which is known to have a masking effect for the HVS: in this toy example, SSIM<sub>9</sub> yields a higher value for the same pixel because its support spans both sides of the edge. Figure 10(b) shows SSIM across the edge of the same profile of (a), and particularly at pixel P. Note that SSIM<sub>9</sub> estimates a higher quality in a larger region around the step.

Despite of the size of its support, however, SSIM is not particularly sensitive to a uniform bias on a flat region. This is particularly true in bright regions, as shown in Figure 10(d), which plots the value of SSIM for the noisy signal of Figure 10(c). The same bias on the two sides of pixel S impacts SSIM's quality assessment differently. Because the term  $l(p)$

Image quality metric	Training cost function	
	$\ell_2$ first, $\ell_1$ after	$\ell_1$ first, $\ell_2$ after
$1000 \cdot \ell_2$	0.3939	<b>0.3896</b>
PSNR	34.76	<b>34.77</b>
$1000 \cdot \ell_1$	<b>12.7932</b>	12.8919
SSIM	0.9544	0.9540
MS-SSIM	0.9753	0.9748
IW-SSIM	0.9634	0.9624
GMSD	0.0432	0.0405
FSIM	0.9777	0.9781
FSIM <sub>c</sub>	0.9770	0.9774

TABLE II: Average value of different image quality metrics for the networks trained for denoising + demosaicking with alternating loss functions. Bold indicates that the network achieves a better score than any of the networks in Table I, see Section V-A.

in Equation 5 measures the error in terms of a contrast, it effectively reduces its impact when the background is bright. This is why, thanks to its multi-scale nature, MS-SSIM solves the issue of noise around edges, but does not solve the problem of the change colors in flat areas, in particular when at least one channel is strong, as is the case with the sky.

Note that, while these are important observations because they generalize to any optimization problem that uses SSIM or MS-SSIM, the definition of a novel image quality metric is beyond the scope of this paper. We limit ourselves to the analysis of their performance within the context of neural networks.

Finally, we would like to point out that using SSIM-based losses in the context of color images introduces an



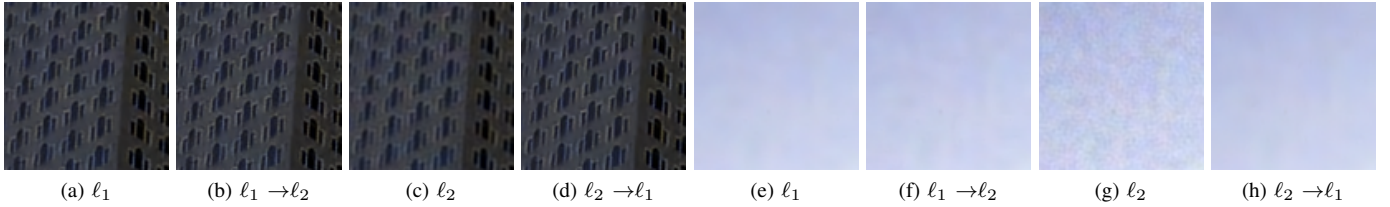


Fig. 8: Visual results of the networks trained alternating loss functions. The insets here correspond the regions marked in Figure 1. Note that  $\ell_i \rightarrow \ell_j$  indicates a network trained first with  $\ell_i$  and then with  $\ell_j$ .

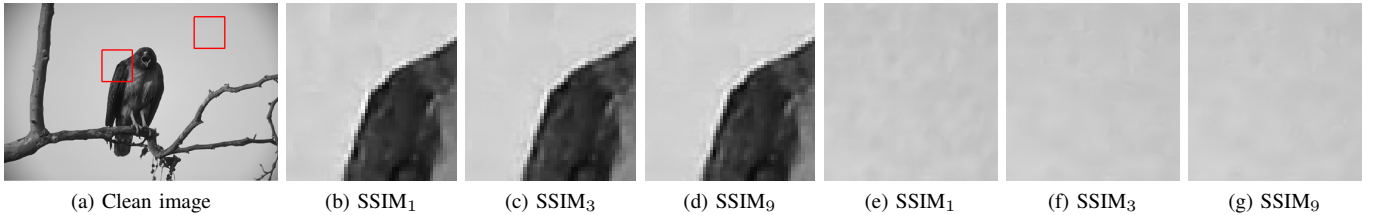


Fig. 9: Comparison of the results of networks trained with SSIM with different sigmas ( $\text{SSIM}_k$  means  $\sigma_G = k$ ). Insets (b)–(d), show an increasingly large halo of noise around the edge: smaller values of  $\sigma$  help at edges. However, in mostly flat regions, larger values of  $\sigma$  help reducing the splotchy artifacts (e)–(g). Best viewed by zooming in on the electronic copy.

approximation, since these metrics were originally designed for grayscale images. A more principled approach would be to base the loss function on a metric defined specifically for color images, such as  $\text{FSIM}_c$ . However, to the best of our knowledge, no such metric has been proposed that is also differentiable. We plan on addressing this problem in future work.

### C. Are we committing the inverse crime?

To numerically evaluate and compare the results of the different loss functions, we need access to the ground truth. For the applications described in this paper, this is only possible if we generate synthetic data. For instance, for the super-resolution experiments, we generate the low-resolution inputs by low-passing and down-sampling the original images. However, when generating synthetic data, one needs to be careful about the validity of the results. In the context of inverse problems, Colton and Kress introduce the concept of *inverse crime*, which is “committed” when synthetic data is fabricated with a forward model that is related to the inverse solver ([29], page 133), potentially hiding shortcomings of the approach. By this definition, it could be argued that we did commit the inverse crime, as we generate the synthetic data for both training and testing with the same process. However, two considerations support our choice. First, while in general the inverse crime is to be avoided, the forward model should resemble the physical system as closely as possible ([29], page 163). In the super-resolution example described above, for instance, performing down-sampling followed by low-passing mimics what the PSF of the lens and the sensor do to the incoming light. Second, we are studying the effect of the loss layer when all other conditions are equal: if we change the forward model, the quality of the results may degrade, but we expect the relationship between losses to hold. This

is exactly what we observed when we tested the networks trained with the different losses on data generated with a varying standard deviation for the Gaussian filter. Changing the standard deviation from 3 to 5 pixels, for instance, causes the average MS-SSIM yielded by the network trained on Mix to drop from 0.9829 to 0.9752. However, the networks trained with  $\ell_1$  and  $\ell_2$  experience a similar drop, yielding an average MS-SSIM of 0.9744 and 0.9748 respectively. We observed a similar behavior for the other metrics. These considerations also apply to joint denoising and demosaicking, as well as JPEG de-blocking.

## VI. CONCLUSIONS

We focus on an aspect of neural networks that is usually overlooked in the context of image restoration: the loss layer. We propose several alternatives to  $\ell_2$ , which is the *de facto* standard, and we also define a novel loss. We use the problems of joint denoising and demosaicking, super-resolution, and JPEG artifacts removal for our tests. We offer a thorough analysis of the results in terms of both traditional and perceptually-motivated metrics, and show that the network trained with the proposed loss outperforms other networks. We also notice that the poor performance of  $\ell_2$  is partially related to its convergence properties, and this can help improve results of other  $\ell_2$ -based approaches. Because the networks we use are fully convolutional, they are extremely efficient, as they do not require an aggregation step. Nevertheless, thanks to the loss we propose, our joint denoising and demosaicking network outperforms CFA-BM3D, a variant of BM3D tuned for denoising in Bayer domain, which is the state-of-the-art denoising algorithm. Table I shows the superiority of Mix over all the other losses. This can also be appreciated by visual inspection even for networks trained with  $\ell_1$ , the closest competitor, which produce stronger artifacts (see, for instance,

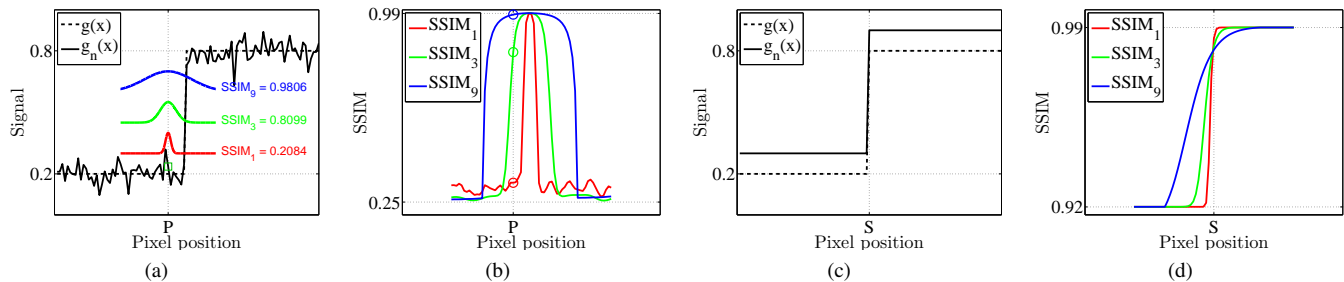


Fig. 10: Panels (a) and (b) show how the size of the support affects SSIM for a pixel P close to an edge, in the presence of zero-mean Gaussian noise. Panels (c) and (d) show that SSIM is less sensitive to a uniform bias in bright regions, such as the one to the right of pixel S (see Section V-B).

Figure 3). We also make the implementation of the layers described in this paper available to the research community.

## REFERENCES

- [1] LeCun, Y., Boser, B., Denker, J.S., Henderson, D., Howard, R.E., Hubbard, W., Jackel, L.D.: Backpropagation applied to handwritten zip code recognition. *Neural computation* **1**(4) (1989) 541–551 **1**
- [2] Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: *NIPS*. (2012) 1097–1105 **1, 2**
- [3] Burger, H., Schuler, C., Harmeling, S.: Image denoising: Can plain neural networks compete with BM3D? In: *IEEE Conference on Computer Vision and Pattern Recognition*. (June 2012) 2392–2399 **1, 2, 3**
- [4] Xu, L., Ren, J.S., Liu, C., Jia, J.: Deep convolutional neural network for image deconvolution. In: *NIPS*. (2014) 1790–1798 **1, 2**
- [5] Zeiler, M., Krishnan, D., Taylor, G., Fergus, R.: Deconvolutional networks. In: *IEEE Conference on Computer Vision and Pattern Recognition*. (2010) 2528–2535 **1**
- [6] Nguyen, A., Yosinski, J., Clune, J.: Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In: *IEEE Conference on Computer Vision and Pattern Recognition*. (2015) **1**
- [7] Hinton, G.E., Srivastava, N., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.R.: Improving neural networks by preventing co-adaptation of feature detectors. *CoRR arXiv:1207.0580* (2012) **1**
- [8] Jain, V., Seung, S.: Natural image denoising with convolutional networks. In Koller, D., Schuurmans, D., Bengio, Y., Bottou, L., eds.: *Advances in Neural Information Processing Systems* 21. Curran Associates, Inc. (2009) 769–776 **1, 2**
- [9] Dong, C., Loy, C., He, K., Tang, X.: Learning a deep convolutional network for image super-resolution. In: *European Conference on Computer Vision*. (2014) 184–199 **1, 2, 5**
- [10] Wang, Y.Q.: A multilayer neural network for image demosaicking. In: *IEEE International Conference on Image Processing*. (2014) 1852–1856 **1, 2**
- [11] Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., Darrell, T.: Caffe: Convolutional architecture for fast feature embedding. *arXiv:1408.5093* (2014) **1, 2**
- [12] Zhang, L., Zhang, L., Mou, X., Zhang, D.: A comprehensive evaluation of full reference image quality assessment algorithms. In: *IEEE International Conference on Image Processing*. (2012) 1477–1480 **1, 2, 3**
- [13] Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing* **13**(4) (2004) 600–612 **1, 2**
- [14] Wang, Z., Simoncelli, E.P., Bovik, A.C.: Multiscale structural similarity for image quality assessment. In: *Asilomar Conference on Signals, Systems and Computers*. Volume 2. (2003) 1398–1402 **1, 2**
- [15] He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *CoRR arXiv:1502.01852* (2015) **2, 5**
- [16] Wang, Z., Bovik, A.: Mean squared error: Love it or leave it? A new look at signal fidelity measures. *IEEE Signal Processing Magazine* **26**(1) (2009) 98–117 **2**
- [17] Brunet, D., Vrscay, E.R., Wang, Z.: Structural similarity-based approximation of signals and images using orthogonal bases. In: *International Conference on Image Analysis and Recognition*. (2010) 11–22 **2**
- [18] Rehman, A., Rostami, M., Wang, Z., Brunet, D., Vrscay, E.: SSIM-inspired image restoration using sparse representation. *EURASIP Journal on Advances in Signal Processing* **2012**(1) (2012) **2**
- [19] Öztireli, A.C., Gross, M.: Perceptually based downscaling of images. *ACM Trans. Graph.* **34**(4) (2015) 77:1–77:10 **2**
- [20] Wang, Z., Li, Q.: Information content weighting for perceptual image quality assessment. *IEEE Transactions on Image Processing* **20**(5) (2011) 1185–1198 **3**
- [21] Sheikh, H., Bovik, A.: Image information and visual quality. *IEEE Transactions on Image Processing* **15**(2) (2006) 430–444 **3**
- [22] Xue, W., Zhang, L., Mou, X., Bovik, A.: Gradient magnitude similarity deviation: A highly efficient perceptual image quality index. *IEEE Transactions on Image Processing* **23**(2) (2014) 684–695 **3**
- [23] Zhang, L., Zhang, D., Mou, X., Zhang, D.: FSIM: A feature similarity index for image quality assessment. *IEEE Transactions on Image Processing* **20**(8) (2011) 2378–2386 **3, 5**
- [24] Winkler, S., Susstrunk, S.: Visibility of noise in natural images. In: *Electronic Imaging 2004, International Society for Optics and Photonics* (2004) 121–129 **3**
- [25] Bychkovskiy, V., Paris, S., Chan, E., Durand, F.: Learning photographic global tonal adjustment with a database of input / output image pairs. In: *IEEE Conference on Computer Vision and Pattern Recognition*. (2011) **5**
- [26] Foi, A.: Clipped noisy images: Heteroskedastic modeling and practical denoising. *Signal Processing* **89**(12) (2009) 2609–2629 **5**
- [27] Danielyan, A., Vehvilainen, M., Foi, A., Katkovnik, V., Egiazarian, K.: Cross-color BM3D filtering of noisy raw data. In: *Intern. Workshop on Local and Non-Local Approximation in Image Processing*. (2009) 125–129 **5**
- [28] Zhang, D., Wu, X.: Color demosaicking via directional linear minimum mean square-error estimation. *IEEE Trans. on Image Processing* **14**(12) (2005) 2167–2178 **5**
- [29] Colton, D., Kress, R.: *Inverse acoustic and electromagnetic scattering theory*. Volume 93. Springer Science & Business Media (2012) **9**



**Hang Zhao** received a B.E. degree in information engineering from Zhejiang University, China in 2013. He is currently a Ph.D. student at the Massachusetts Institute of Technology. His research interests include computer vision, machine learning, and their robotic applications.



**Orazio Gallo** earned a M.S. degree in Biomedical Engineering from Politecnico di Milano in 2004, and a Ph.D. in Computer Engineering from the University of California at Santa Cruz in 2011. From 2004 to 2006 he was a research assistant at the Smith-Kettlewell Eye Research Institute. In 2011, Orazio joined NVIDIA, where he is a senior research scientist. His research interests range from computational photography, to image processing, applied visual perception, and computer vision. Orazio is also an associate editor of the IEEE Transactions of Computational Imaging, as well as Signal Processing: Image Communication.



**Iuri Frosio** got his M.S. and Ph.D. degrees in biomedical engineering at the Politecnico di Milano in 2003 and 2006, respectively. He was a research fellow at the Computer Science Department of the Università degli Studi di Milano from 2003 to 2006, and an assistant professor in the same Department from 2006 to 2013. In the same period, he worked as a consultant for various companies in Italy and abroad. He joined NVIDIA in 2014 as senior research scientist. His research interests include image processing, computer vision, inertial sensors, and machine learning. Iuri is also an associate editor of the Journal of Electronic Imaging.



**Jan Kautz** leads the Visual Computing Research team at NVIDIA, working predominantly on computer vision problems—from low-level vision (denoising, super-resolution, computational photography) and geometric vision (structure from motion, SLAM, optical flow) to high-level vision (detection, recognition, classification), as well as machine learning problems such as deep reinforcement learning and efficient deep learning. Before joining NVIDIA in 2013, Jan was a tenured faculty member at University College London. He holds a BSc in Computer Science from University of Erlangen-Nürnberg (1999), an MMath from the University of Waterloo (1999), received his Ph.D. from the Max-Planck-Institut für Informatik (2003), and worked as a post-doc at the Massachusetts Institute of Technology (2003-2006).

Jan was program co-chair of the Eurographics Symposium on Rendering 2007, program chair of the IEEE Symposium on Interactive Ray-Tracing 2008, program co-chair for Pacific Graphics 2011, and program chair of CVMP 2012. He co-chaired Eurographics 2014 and was on the editorial board of IEEE Transactions on Visualization & Computer Graphics as well as The Visual Computer.