

000  
001  
002  
003 Deep Residual Autoencoder for Real Image Denoising  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
Abstract

TODO — amet nisl suscipit adipiscing bibendum est ultricies integer quis auctor elit sed vulputate mi sit amet mauris commodo quis imperdiet massa tincidunt nunc pulvinar sapien et ligula ullamcorper malesuada proin libero nunc consequat interdum varius sit amet mattis vulputate enim nulla aliquet porttitor lacinia luctus accumsan tortor posuere ac ut consequat semper viverra nam libero justo laoreet sit amet cursus sit amet dictum sit amet justo donec enim diam vulputate ut pharetra sit amet aliquam id diam maecenas ultricies mi eget mauris pharetra et ultrices neque ornare aenean euismod elementum nisi quis eleifend quam adipiscing vitae proin sagittis nisl rhoncus mattis rhoncus urna neque viverra justo nec ultrices dui sapien eget mi proin sed libero enim sed faucibus turpis in eu mi bibendum neque egestas congue quisque egestas diam in arcu cursus euismod quis viverra nibh cras pulvinar mattis nunc sed blandit libero volutpat sed cras ornare arcu dui vivamus arcu felis bibendum ut tristique et egestas quis ipsum suspendisse

## 1. Introduction

Image denoising aims at removing the noise of a given noisy image which is an essential task in low-level computer vision. Nowadays, with the increase of digital imaging, image denoising has found many real-world use cases such as medical image denoising. In the literature there can be found considerable amount of research done for the removal of various noise models (e.g., salt and pepper noise, additive white Gaussian noise (AWGN)). Although those noise models are used to represent real-world noise, there is a considerable difference between current noise models and real-world noise [1].

In the literature, learning based methods are proven their performance. Recent state-of-the-art image denoising methods such as FDnCNN[16] and REDNet [8] which use deep Convolutional Neural Networks

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

Anonymous WACV submission

Paper ID \*\*\*\*

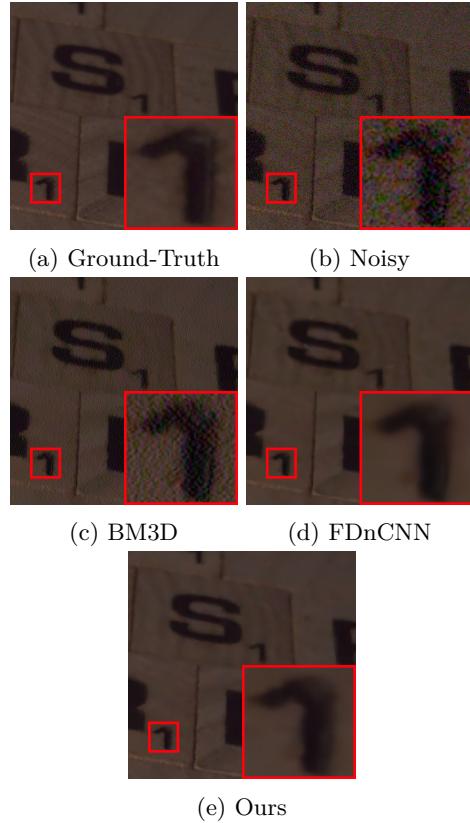


Figure 1: A and B are ground-truth and noisy image pair from SIDD dataset [1]. C through E are sample denoised images using [4], [16] and Our method respectively. (Best viewed on high-resolution display.)

(CNNs) have remarkable denoising performance on various noise models. However, since real-world noise is not fully represented by any noise model, performance of these methods are limited by the noise model used in the training of the models [1]. Even in some cases, depending on the noise model and noise level, traditional methods such as BM3D [4] can outperform the learning based methods.

Another factor which affects the performance of

108 learning based methods is the loss function used while  
109 training the model. In the literature mean square error  
110 (L1 loss) and mean absolute error (L2 loss) are  
111 widely used as loss function. However, these functions  
112 do not consider the nature of Human Visual System  
113 (HVS). As an example, L2 loss assumes that the im-  
114 pact of noise is independent from the characteristics of  
115 the image [19]. Due to the sensitivity of the HVS to lo-  
116 cal luminance, contrast and structure, this assumption  
117 causes poor result for human observers [13].  
118

119 In this paper, we propose a deep convolutional au-  
120 toencoder with skip connections in order to pass val-  
121 uable information to the deeper layers of the model.  
122 Our model also utilizes several loss functions in order  
123 to boost its denoising performance. Also, we trained  
124 our model with real-world noisy and ground truth im-  
125 age pairs provided by Smartphone Image Denoising  
126 Dataset (SIDD) [1]. Our model can achieve com-  
127 petitive results with state-of-the-art methods such as  
128 BM3D [4], FDnCNN [16] and REDNet [8] in blind im-  
129 age denoising in color images.

130 The rest of the paper is organized as follows. Section  
131 2 provides a brief summary of existing image denoising  
132 methods, noise types, datasets and image quality  
133 metrics used in the literature. Section 3 presents the  
134 proposed method and discusses about its features. In  
135 Section 4 we talk about our experimental setup and ex-  
136 tensive evaluations. We present our results in Section  
137 5 and conclude the paper by discussing our findings in  
138 Section 6.

## 139 2. Related Work

140 In this section, we will present existing traditional  
141 and state-of-the-art image denoising methods.  
142 Then, we will discuss about available image denoising  
143 datasets. Also, we will present image quality assess-  
144 ments which is necessary when measuring the quality  
145 of an image in a quantitative manner.

### 146 2.1. Traditional Methods

147 BM3D [4] is widely known as a traditional im-  
148 age denoising method in the literature. It uses ef-  
149 fective filtering in 3D transform domain by combining  
150 sliding-window transform processing with block match-  
151 ing. Burger et al. [2] showed that when the noisy image  
152 does not contain any regular structure, denoising per-  
153 formance BM3D is decreased. Also, it is shown that a  
154 plain multi-layer perceptron (MLP) can achieve simi-  
155 lar denoising performance with BM3D by Burger et al.  
156 [2].

## 157 2.2. State-Of-The-Art Methods

158 One of the recent additions to image denoising liter-  
159 ature is denoising autoencoders. Autoencoders aim to  
160 learn an approximation to identity function. By taking  
161 this property of autoencoders into account, denoising  
162 autoencoders forces the model to learn reconstruction  
163 of the input given its noisy version. [6].

164 In the literature there are many proposed methods  
165 for image denoising which uses denoising autoencoders  
166 [6, 14, 3, 15]. Gondara [6] proposed an autoencoder  
167 based denoiser for medical imaging domain. Xie [6]  
168 et al. and Ye [6] et al. proposed stacking multiple  
169 denoising autoencoders in order to better model the  
170 noise.

## 171 2.3. Datasets

172 As it is with many machine learning related research,  
173 data is one of the major factors which limits the per-  
174 formance. Image denoising is no exception. However,  
175 recently with the release of new datasets such as SIDD  
176 [1] and Darmstadt Noise Dataset (DND) [11], the bot-  
177 tleneck caused by the lack of high-quality data has de-  
178 creased. In the literature SIDD and DND are also used  
179 as a benchmarking tool for image denoising methods.

## 180 2.4. Image Quality Assessment

181 Image Quality Assessment (IQA) is another related  
182 important research topic. It aims to measure the qual-  
183 ity of the images in a quantitative way. For image de-  
184 noising domain, having a quantitative metric to mea-  
185 sure the quality of the images precisely is vital for opti-  
186 mizing learning based denoising models. Peak Signal to  
187 Noise Ratio (PSNR) and Structural Similarity (SSIM)  
188 [13] are widely used in the literature for measuring the  
189 quality of images. Feature Similarity Index (FSIM) [17]  
190 is another quality metric which measures the dissimi-  
191 larity between two images based on local information.  
192 Furthermore, Zhao [18] et al. proposed a novel method  
193 which is a combined version of Multi-Scale Structural  
194 Similarity (MS-SSIM) and mean square error (MSE)  
195 in order to eliminate each methods drawbacks.

## 196 3. Proposed Method

197 We propose a deep fully convolutional autoencoder  
198 with residual connections in order to reduce the degra-  
199 dation problem which occurs in deep networks. We  
200 benefit from Convolutional (Conv.), Transposed Con-  
201 volutional (ConvT.), Rectified Linear Unit (ReLU) [9]  
202 and Sigmoid layers in the architecture.

203 162  
204 163  
205 164  
206 165  
207 166  
208 167  
209 168  
210 169  
211 170  
212 171  
213 172  
214 173  
215 174  
216 175  
217 176  
218 177  
219 178  
220 179  
221 180  
222 181  
223 182  
224 183  
225 184  
226 185  
227 186  
228 187  
229 188  
230 189  
231 190  
232 191  
233 192  
234 193  
235 194  
236 195  
237 196  
238 197  
239 198  
240 199  
241 200  
242 201  
243 202  
244 203  
245 204  
246 205  
247 206  
248 207  
249 208  
250 209  
251 210  
252 211  
253 212  
254 213  
255 214  
256 215

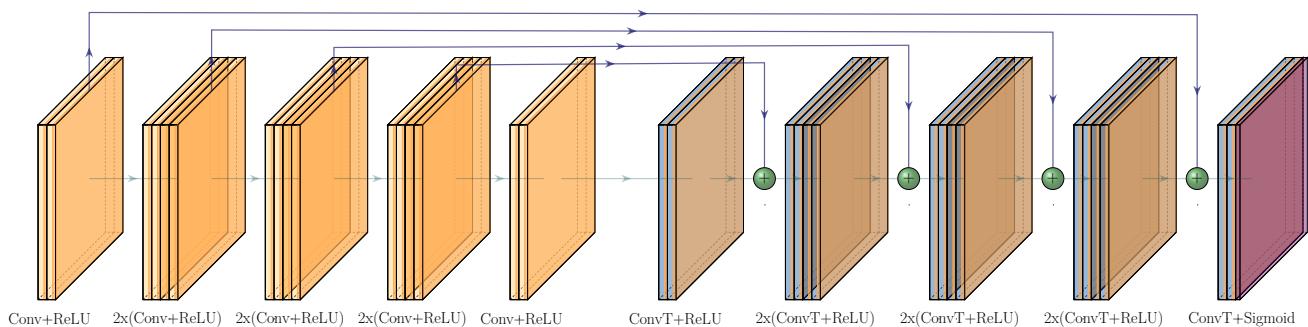


Figure 2: Proposed network architecture where "+" denotes element wise sum of feature maps.

### 3.1. Network Architecture

The network takes a color image as input with size  $3 \times M \times N$  and passes it through eight times convolution and ReLU (Conv+ReLU) layers with one padding applied to keep the input size constant. We refer to these layers as encoder layers. Output of the encoder layers are then passed to the decoder layers which are seven transposed convolution and ReLU (ConvT+ReLU) layers followed by one transposed convolution and sigmoid layer. We use Sigmoid function in output layer to clip the output values of our network between 0 and 1. Like encoder layers, decoder layers also have one padding applied to keep the input size constant. Throughout the network we placed symmetric residual connections which transfers high-level information to deeper layers in order to reduce the effects of degradation problem. Visualization of our architecture can be seen in the figure 2.

We used 3 channel in, 64 channel out convolutional layer as input layer and 64 channel in, 3 channel out Transposed Convolutional layer as output layer. All other layers have 64 channel input and 64 layer output. We also used  $3 \times 3$  kernel size for convolution and transposed convolution layers throughout the network. We used ADAM optimizer [7] with weight decay rate of 0.05 and learning rate  $10^{-4}$ . At 60% and 90% percent of the training learning rate is multiplied by  $10^{-1}$ . As loss function we used many variations of L1, L2, PSNR, SSIM, MS-SSIM and FSIM [17] functions.

### 3.2. Training Dataset

We preferred Smartphone Image Denoising Dataset (SIDD) [1] due to its quality and amount of data it provides. SIDD dataset provides high-quality real-world noisy images and their noise free ground truths which are taken with smartphones and DSLR cameras respectively. SIDD dataset has 3 versions which are small,

medium and full version. Small version of the dataset contains 160 image pairs and it is approximately 6 GB. Medium version of the dataset contains 320 image pairs and approximately 12 GB. Full version of the dataset contains 12000 image pairs and approximately 450 GB. We used medium version of the SIDD dataset due to its manageable size. As training dataset, we randomly selected 90% of medium version of the dataset. Remaining 10% of the data is used as validation set. Before the training phase we randomly cropped  $128 \times 128$  patches from images batch size x number of epochs times and saved them for fast training.

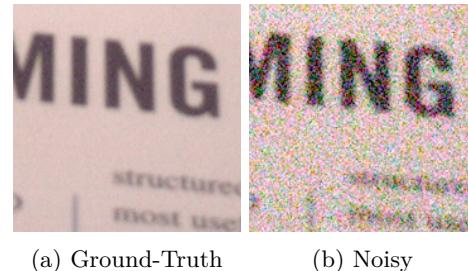


Figure 3: Sample cropped patches from SIDD dataset.

### 3.3. Validation Dataset

As we mentioned in the training dataset section, we used 10% of the SIDD medium version dataset as validation set. We recorded validation loss by using validation dataset at each epoch in order to better understand and optimize the training process of the proposed model. We calculated validation loss as mean loss throughout the validation set for each epoch.

## 4. Experiments

We focused on removing real world noise from color images. To implement, train and evaluate our proposed

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

method we used Pytorch [10]. All the experiments are conducted in Python 3.8.2 environment using PyTorch 1.6.0 running on a PC with Intel(R) Core(TM) i7-3770K CPU and Nvidia GTX 1080 GPU with 8 GB video memory. The training of a single model can be done in about 3 hours. Some demonstrations of the proposed method are available at: [https://yilmazdoga.com/image\\_denoising\\_using\\_autoencoders](https://yilmazdoga.com/image_denoising_using_autoencoders) also, the implementation is available at: [https://github.com/yilmazdoga/image\\_denoising\\_using\\_autoencoders](https://github.com/yilmazdoga/image_denoising_using_autoencoders).

#### 4.1. Performance Evaluation Criteria

In order to evaluate the performance of our proposed model, we needed a reliable metric to quantitatively measure the quality of the resulting denoised image. In the literature there are many IQA metrics which can be used for evaluating the quality of denoised images produced by our proposed method. Neural Image Assessment (NIMA) [12] and Deep Image Structure and Texture Similarity (DISTS) [5] can be given as examples of state-of-the-art image quality evaluation metrics. However, these state-of-the-art metrics are not fully adopted by the literature. For the sake of comparability, we are evaluating the performance of our proposed method with traditional metrics which are SSIM and PSNR.

#### 4.2. Test Datasets

Two test sets are used for evaluating color image denoising performance which are Darmstadt Noise Dataset (DND) and SIDD. DND dataset consists of 50 high-resolution images with realistic image noise. We used noisy images provided by DND as our first test dataset without applying any manipulation. From SIDD dataset we used 256 by 256 image patches from images which are not used in our training or validation sets as our second test dataset.

#### 4.3. Effects of Residual Learning

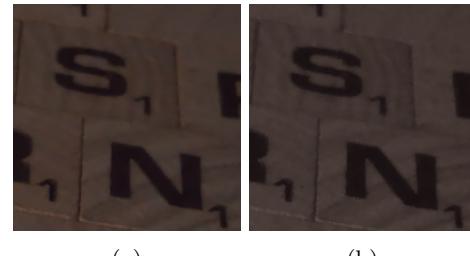
To test the effects of residual learning, we trained two networks with and without residual connections. Training of these two networks was done using the same dataset which is the medium sized version of SIDD dataset. We observed that the network with residual connections perform better than the network without residual connections. In the table 1 average PSNR and SSIM for both models can be found. Also, a sample image denoised with both networks can be found in the figure 4.

#### 4.4. Effects of Different Loss Functions

Beside optimizing the performance of our network, we also aimed to optimize the loss function we used

Table 1: Average PSNR and SSIM scores of proposed network with and without residual connections.

Method Name	Average PSNR	Average SSIM
With Residual Connections	37.75	0.898
Without Residual Connections	32.68	0.872



(a) (b)

Figure 4: Sample denoised image using proposed model with residual connections (a), without residual connections (b).

with our network. With the objective of using a better loss function we trained our model with 5 different loss functions which are L1, L2, SSIM, MS-SSIM, and L1 + MS-SSIM which is proposed by Zhao [18] et al. Performance comparison of all loss functions can be found in table 2. Also, sample images denoised with all loss functions can be found in figure 5.

Table 2: Average PSNR and SSIM scores of proposed network with L1, L2, SSIM, MS-SSIM and L1 + MS-SSIM loss functions.

Loss Function	Average PSNR	Average SSIM
L1	37.93	0.895
L2	37.71	0.891
SSIM	37.25	0.900
MS-SSIM	35.91	0.895
L1 + MS-SSIM	37.75	0.898

## 5. Results

There is a benchmarking tool provided by SIDD dataset which evaluates a given model with respect to PSNR and SSIM on benchmark dataset.

## 6. Conclusions

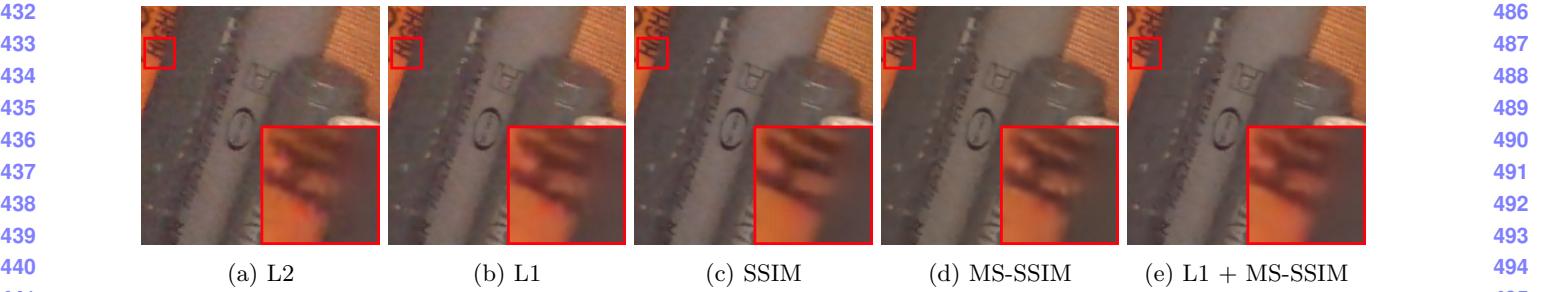


Figure 5: Sample images denoised with proposed network using L1, L2, SSIM, MS-SSIM and L1 + MS-SSIM loss functions.

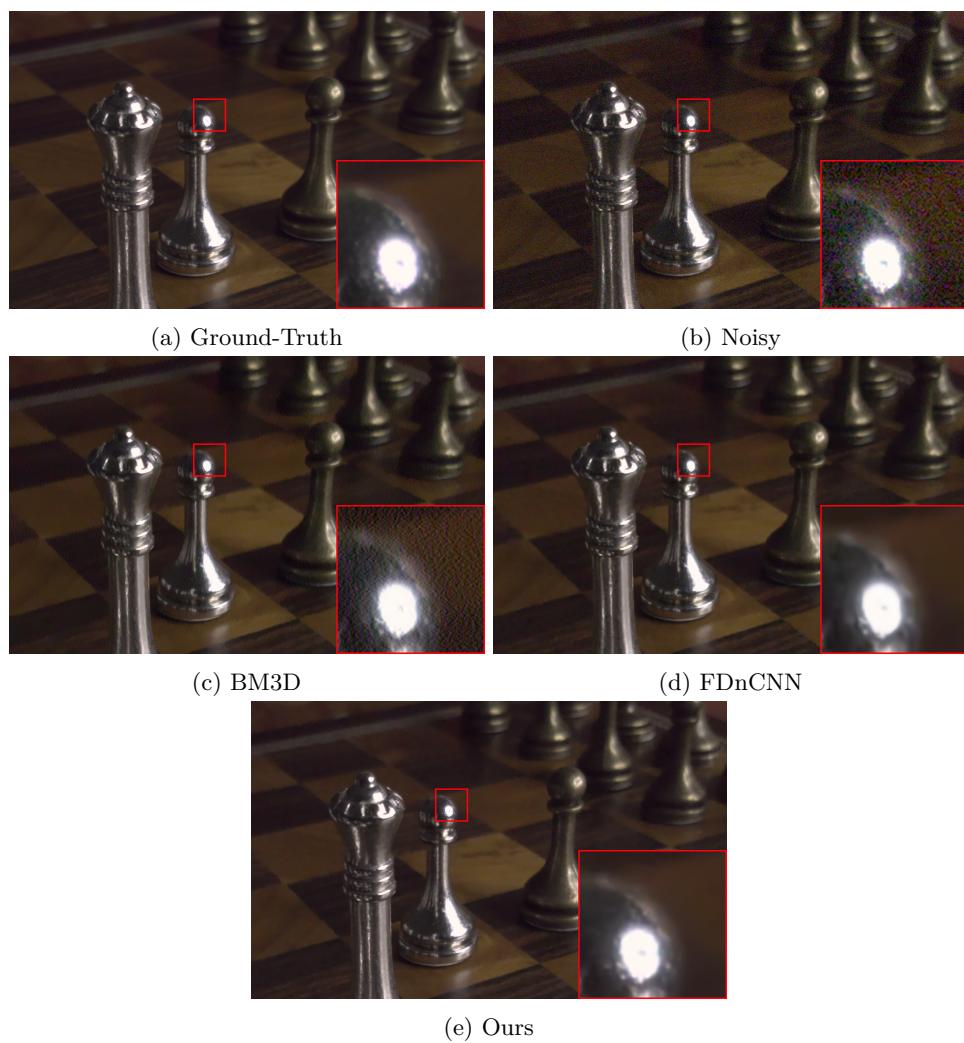


Figure 6: A and B are high-definition (1920x1200) ground-truth and noisy image pair from SIDD dataset [1]. C through E are sample denoised images using [4], [16] and Our method respectively. (Best viewed on high-resolution display.)

## References

- [1] Abdelrahman Abdelhamed, Stephen Lin, and Michael S. Brown. A high-quality denoising dataset

for smartphone cameras. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR), June 2018. 1, 2, 3, 5

- 540 [2] H. C. Burger, C. J. Schuler, and S. Harmeling. Image 594  
541 denoising: Can plain neural networks compete with 595  
542 bm3d? In 2012 IEEE Conference on Computer Vision 596  
543 and Pattern Recognition, pages 2392–2399, 2012. 2 597  
544 [3] Kyunghyun Cho. Boltzmann machines and denoising 598  
545 autoencoders for image denoising, 2013. 2 599  
546 [4] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian. 600  
547 Color image denoising via sparse 3d collaborative 601  
548 filtering with grouping constraint in luminance- 602  
549 chrominance space. In 2007 IEEE International 603  
550 Conference on Image Processing, volume 1, pages I – 604  
551 313–I – 316, 2007. 1, 2, 5 605  
552 [5] Keyan Ding, Kede Ma, Shiqi Wang, and Eero P. 606  
553 Simoncelli. Image quality assessment: Unifying 607  
554 structure and texture similarity, 2020. 4 608  
555 [6] L. Gondara. Medical image denoising using 609  
556 convolutional denoising autoencoders. In 2016 IEEE 610  
557 16th International Conference on Data Mining Workshops 611  
558 (ICDMW), pages 241–246, 2016. 2 612  
559 [7] Diederik P. Kingma and Jimmy Ba. Adam: A method 613  
560 for stochastic optimization, 2014. 3 614  
561 [8] Xiao-Jiao Mao, Chunhua Shen, and Yu-Bin Yang. 615  
562 Image restoration using very deep fully convolutional 616  
563 encoder-decoder networks with symmetric skip 617  
564 connections. CoRR, abs/1603.09056, 2016. 1, 2 618  
565 [9] Vinod Nair and Geoffrey Hinton. Rectified linear units 619  
566 improve restricted boltzmann machines vinod nair. 620  
567 volume 27, pages 807–814, 06 2010. 2 621  
568 [10] Adam Paszke, Sam Gross, Soumith Chintala, Gregory 622  
569 Chanan, Edward Yang, Zachary DeVito, Zeming Lin, 623  
570 Alban Desmaison, Luca Antiga, and Adam Lerer. 624  
571 Automatic differentiation in pytorch. In NIPS 2017 625  
572 Workshop on Autodiff, 2017. 4 626  
573 [11] Tobias Plötz and Stefan Roth. Benchmarking 627  
574 denoising algorithms with real photographs, 2017. 2 628  
575 [12] Hossein Talebi and Peyman Milanfar. Nima: Neural 629  
576 image assessment. IEEE Transactions on Image 630  
577 Processing, 27(8):3998–4011, Aug 2018. 4 631  
578 [13] Zhou Wang, Alan Bovik, Hamid Sheikh, and Eero 632  
579 Simoncelli. Image quality assessment: From error 633  
580 visibility to structural similarity. Image Processing, IEEE 634  
581 Transactions on, 13:600 – 612, 05 2004. 2 635  
582 [14] Junyuan Xie, Linli Xu, and Enhong Chen. Image 636  
583 denoising and inpainting with deep neural networks. In 637  
584 F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Wein- 638  
585 berger, editors, Advances in Neural Information Pro- 639  
586 cessing Systems 25, pages 341–349. Curran Associates, 640  
587 Inc., 2012. 2 641  
588 [15] X. Ye, L. Wang, H. Xing, and L. Huang. Denoising 642  
589 hybrid noises in image with stacked autoencoder. In 643  
590 2015 IEEE International Conference on Information 644  
591 and Automation, pages 2720–2724, 2015. 2 645  
592 [16] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang. 646  
593 Beyond a gaussian denoiser: Residual learning of deep 647  
594 cnn for image denoising. IEEE Transactions on Image 595  
595 Processing, 26(7):3142–3155, 2017. 1, 2, 5 596