

WAMnet for Image Denoising

Aaron Bies¹, Adarsh Djeacoumar¹, and Shailesh Mishra¹

Saarland University

Abstract. Photo-realistic images generated using path tracing often contain a significant amount of noise as it involves Monte-Carlo integration. Raising the total number of samples in an image reduces the amount of noise but increases the computational cost of the algorithm. It is desirable to reduce the noise and generate high-quality renders at a low computational cost. We present a lightweight novel MLP based filter architecture that denoises images rendered at 1 sample per pixel and its corresponding normal. Furthermore, this lightweight network makes use of fragment shaders and delivers a run-time performance of roughly 11 frames per second. This architecture significantly reduces the computation cost required to render high-quality images. Our network reduces noise by approximately 40% and improves the resulting structure by over 200%. We found that using the corresponding normal maps of the input image in training significantly enhances the performance. Therefore, we believe extension based on the concepts from computer graphics can improve the performance of denoising networks.

Keywords: Ray Tracing · Image Denoising · Neural Networks.

1 Introduction

Path tracing is the state of the art method for generating photo-realistic images in computer graphics and has been for many years [1]. In its most common form, it is a method based on Monte Carlo integration, which works by shooting rays from the camera into the scene, randomly bouncing them off surfaces, and sampling light sources to iteratively approximate the brightness received by each pixel of our virtual camera [1]. Intuitively it traces the path a ray of light would take in the real world on its way from a light source to the camera, only in reverse.

To produce high-quality renders, these methods perform sampling on each pixel and average the incoming radiance from all samples on that pixel. Despite being a robust method, i.e., being free from bias, it requires large computation time, and the computation explodes with the increase in resolution. Therefore, computer vision algorithms are applicable in reducing the noise for a rendered image obtained using low sample counts. In this report, we highlight the contribution that allows us to denoise the input images rendered using 1 sample per pixel (SPP) in real-time. The filter we use is a multi-layer perceptron network that is lightweight and can run on a single fragment shader.

39 2 Prior Work 39

40 Chaitanya et al. [2] described an Autoencoder with recurrent connects that improved the temporal stability for sequences of sparsely sampled input images. 41 They attempt to remove a class of noise present in Monte Carlo rendering by 42 reconstructing the global illumination in the scene, using Convolutional Neural 43 Networks. They demonstrate that their method models relationships, such as 44 depth and normals, with the input channels. 45

46 Müller et al. [4] in their paper study non-linear independent components 47 estimation (NICE) and its extensions to generate samples in Monte Carlo inte- 48 gration using deep neural networks. They demonstrate learning of joint path- 49 sampling densities in the primary sample space and importance sampling of 50 multi-dimensional path prefixes. Their approach extracted and leveraged the 51 conditional directional densities for path guiding. In a way, their experiments 52 revealed that Neural path guiding performed on-par or better than other path 53 tracing techniques at equal sample count. 53

54 Zhang et al. in [6] introduced the Denoising Convolutional Neural Network 55 (DnCNN) for denoising images. Their experiments demonstrate that the DnCNN 56 can handle Gaussian noise of unknown noise level (blind Gaussian denoising). 57 They adopt the residual learning formulation to train a residual mapping $R(y) \approx$ 58 v , and then $x = y - R(y)$ gives the predicted noise for an input given by $y = x + v$. 59 They address various problems such as Gaussian denoising, single image super- 60 resolution, and JPEG image deblocking. 60

61 Isola et al. in [3] illustrate their model pix2pix, a conditional GAN, for image- 62 to-image translation tasks. They validate their model on settings in graphics 63 (photo generation), and vision (semantic segmentation). Community access to 64 the pix2pix model further proved its robustness as a general-purpose solution 65 for image-to-image translation tasks. 65

66 We take inspiration from these works to construct a lightweight network 67 that can be a part of a graphics pipeline and fit in a fragment shader, and can 68 perform on-par or close to the state-of-the-art networks. To this end, we propose 69 the “Wait a minute network” (WAMnet). WAMnet is a multi-layer perceptron 70 network that is scene-specific, fits in a pipeline, and can generate a 3D scene in 71 real-time. Note that while we do not train our model on different orientations, 72 the pipeline also generates the scene for unknown orientations equally good as 73 the scene that was trained. 73

74 3 Architecture of WAMnet 74

75 Given an $N \times N$ patch from a noisy input image as well as the corresponding 76 normal map, our network predicts the color of the center pixel of that chunk. 77 This allows the network to take neighbouring pixels, information from the normal 78 maps and inferred surface properties into account. The network, WAMnet, is a 79 5 layer Multi-layer perceptron each with hidden layers consisting of 32 neurons. 80 These neurons are activated by a periodic activation function, inspired by the 81 SiREN network [5]. The overall architecture is depicted in Fig. 1. 81

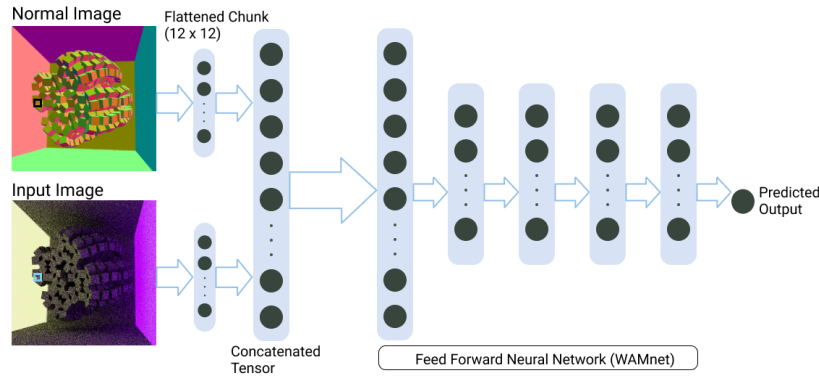


Fig. 1: Architecture of WAMnet

82 4 Experiments 82

83 4.1 Dataset 83

84 We start our experiment by designing a custom scene in a fragment shader. Using 84
 85 this custom scene, we render our training sets as well as testing sets. The training 85
 86 set consists of 1000 rendered images along with their corresponding normals. The 86
 87 input image is rendered at 1 SPP while the ground truth is rendered at 1024 87
 88 SPP. Our testing dataset consists of 100 rendered images at the same settings as 88
 89 the training set. For both testing and training set, we set up our image resolution 89
 90 to be 1024×1024 shown in Appendix A. 90

91 4.2 WAMnet 91

92 We train our model for a single scene using a single pair of input and output 92
 93 images. For each scene, we create 20,000 training samples by randomly sampling 93
 94 an input chunk of size 12×12 from 1024×1024 input image. Out of these 20,000, 94
 95 we pass 10,000 random chunks to our network in one iteration. We repeat the 95
 96 training for 20,000 iterations. This approach of training our neural net on a 96
 97 single scene requires approximately 90 secs for a given scene. We designed a 97
 98 series of experiments to evaluate the performance of our model. 98

99 First, we train WAMnet by giving 3 channel (RGB) input chunk and its corre- 99
 100 sponding output chunk (WAMnet\normals). In another experiment, we append 100
 101 the normal map generated during the pipeline to our input chunk creating a 6 101
 102 channel input chunk and compare with the same output chunk (WAMnet). We 102
 103 noticed considerable improvements in the geometry of the shapes in the scene 103
 104 when using normals. When trained without a normal map, WAMnet fails to 104
 105 maintain the sharp geometry and works closer to a blur filter. In both cases, we 105
 106 use Adam optimizer for minimizing the mean squared error (MSE) loss described 106

107 as:

$$L(y, \tilde{y}) = \frac{1}{N} \sum_{i=1}^N \|(y - \tilde{y})\|^2 \quad (1)$$

107

108 **4.3 pix2pix**

108

We split our dataset into smaller patches (128×128) and pass it to the network. We use the U-Net that is 6 layers deep and also takes Gaussian noise (z) as input to prevent deterministic outputs, as our encoder-decoder architecture for our generator (G). The discriminator (D) architecture is the PatchGAN network, which is a convnet and is 3 layers deep. Note that we are trying to optimize our architecture to run on a fragment shader and hence we design our models as light as possible. We primarily use MSE (L2 norm) and reconstruction loss (L1 norm) for our custom loss function. We observed that penalizing the MSE and reconstruction loss in the ratio of 1:100 gave the best results. Our custom loss function is described below:

$$\begin{aligned} L_D &= \frac{1}{2} (L(D(x), target_{real}) + L(D(G(x, z)), target_{fake})) \\ L_G &= L(D(G(x, z)), y) \\ L_{L1} &= \frac{1}{N} \sum_{i=1}^N \|(D(G(x, z)) - target_{real})\| \\ L_{GAN} &= L_G + 100 * L_{L1} \end{aligned}$$

109 **4.4 DnCNN**

109

110 We resize our dataset from its original 1024×1024 to 256×256 to fit in a
111 GPU. The DnCNN architecture we use is a convnet of depth 16. [6] describe the
112 DnCNN architecture to be able to learn the residual map $R(y) \approx v$, where v is
113 the noise. Then, the output image is calculated by $x = y - R(y)$, where x is the
114 clean image. We adopt the averaged MSE between the desired and estimated
115 residual mappings as the loss function:

$$L(\phi) = \frac{1}{2N} \sum_{i=1}^N \|(R(y_i, \phi) - (y_i - x_i))\|^2 \quad (2)$$

116 **5 Results**

116

117 **5.1 Evaluation Metrics**

117

118 To assess the performance of the neural network we employ the following met-
119 rics: MSE, PSNR, and SSIM. By definition, the lower the MSE, the higher the
120 PSNR, and closer the predicted image is to the ground truth. SSIM, which stands
121 for Structural Similarity Index Measure, compares the luminance, contrast, and
122 structure between the predicted image and the ground truth. The SSIM score is
123 present in $[0, 1]$ and the closer the score is to 1, the better the image.

Table 1: **Performance comparison with relevant methods.** Please note despite pix2pix showing superior metrics, the overall results introduces unnecessary artifacts.

Test Network	MSE	PSNR	SSIM
Input Images	82.93	66.67	0.21
DnCNN	79.03	67.16	0.44
WAMnet	48.04	72.17	0.68
pix2pix	39.28	74.18	0.70

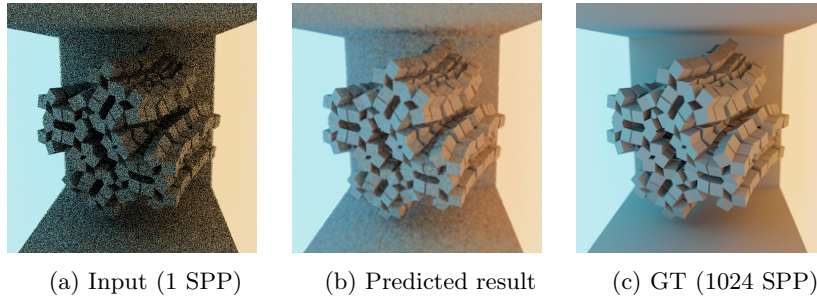


Fig. 2: Result showing prediction from WAMnet

124 5.2 Performance 124

125 Table 1 shows the performance of the networks that we trained on our dataset. 125
 126 We can observe that the input images rendered at 1 SPP have a high MSE, low 126
 127 PSNR, and poor SSIM score. Our input images fail to maintain the structure 127
 128 of the desired scene view. Given such a scene, our network can reconstruct 128
 129 the geometry (edges of the mini cubes are clear) and denoise it resulting in a 129
 130 structure that is 68% similar to the ground truth. We also observe that using normal 130
 131 maps in the input sequence reconstructs better than in the absence of normals. 131

132 Furthermore, our network outperforms DnCNN by a significant margin in 132
 133 all three metrics. Finally, pix2pix performs better than WAMnet by 2%. While 133
 134 this is the case, the complexity of the pix2pix model poses implicit constraints 134
 135 on the fragment shader and hence is unsuitable for our pipeline. Moreover, since 135
 136 we train pix2pix on patches sampled from the input images, we observe blocking 136
 137 artifacts when we reconstruct the image from the noisy input image. In this sense, 137
 138 WAMnet is not only good but also consistent and fits our purpose perfectly. 138

139 Fig. 2 shows how our network can retain the scene geometry as well as soft 139
 140 shadows. Since we feed in the normal map, we believe that the network is able 140
 141 to estimate the light contribution on the missing pixels using the information 141
 142 from noisy input as well as the normal map. 142

143 Fig. 3 compares results of WAMnet, pix2pix and DnCNN for one image from 143
 144 our test set. We observe that WAMnet and pix2pix perform extremely better 144

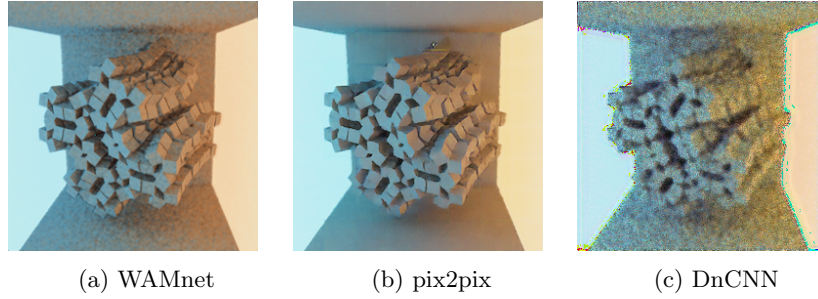


Fig. 3: Comparison of predictions with different methods

145 than DnCNN. While pix2pix looks better in low resolution, the inherent block 145
 146 artifacts that come as a result of patch training are visible in higher resolutions. 146

147 6 Further Studies 147

148 We also conducted further experiments to test out the generalization capability 148
 149 of our network. While networks such as DnCNN and pix2pix generalize for the 149
 150 entire training set and are capable of working independently on the test set, 150
 151 WAMnet is a more specialised network and can only work on the scene that it 151
 152 is trained on. Considering the fact that the model contains only 5 layers, while 152
 153 compared to the deep layers in pix2pix and DnCNN, the trade-off is justified as, 153
 154 just by training on one particular scene for 90 seconds, our pipeline can generate 154
 155 a 3D scene that varies in time at 11 frames per second. 155

156 Moreover, the resulting scene is comparable to a scene that is rendered at 30 156
 157 SPP. To put it in perspective, the renderer takes about a minute to generate an 157
 158 image at 1 SPP, while it takes about 7 minutes to generate an image at 30 SPP. 158
 159 Our pipeline aims to achieve generating a scene at higher SPP given 1 SPP (the 159
 160 lowest possible). 160

161 6.1 Different lights for training and testing 161

162 Instead of training and testing on the same scene with the same light, we changed 162
 163 the light colors on training as well as testing the image. For the same scene when 163
 164 trained on an arbitrary set of light colors \mathbf{A} , and tested on a different set of light 164
 165 colors \mathbf{B} , we noticed that our network fails to reproduce the original color. This 165
 166 limitation arises from the size of our network. The results can be observed in 166
 167 Fig. 4 167

168 6.2 Training on multiple lights 168

169 We also conducted experiments where we trained on different sets of lights. When 169
 170 trained on multiple sets of lights, we found that our model tries to produce the 170

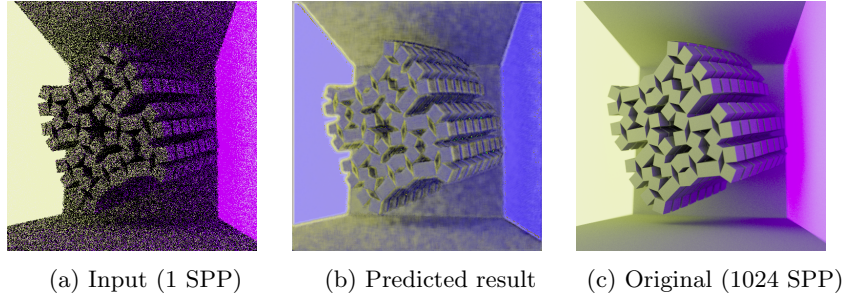


Fig. 4: Result showing effect of different colors for training and testing on same scene

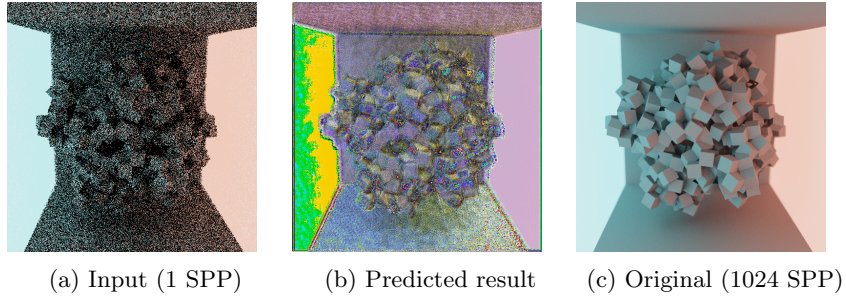


Fig. 5: Result showing effect of use of multiple colors in training phase

171 colors but is not able to reproduce them successfully. The visualization can be 171
 172 seen from Fig. 5 172

173 7 Discussion 173

174 From our experiments, we notice that passing in the normals allows the MLP 174
 175 to learn the scene geometry. Since we were constrained to run it on a fragment 175
 176 shader as a part of our graphics pipeline, we didn't go for large and deep net- 176
 177 works as they can not be compiled using the fragment shader. Therefore, we 177
 178 sacrifice generalization and work on specialization because of the budget con- 178
 179 straint imposed by the engineering design. 179

180 8 Conclusions 180

181 With our work, we have successfully designed a neural network solution that is 181
 182 capable of running inside a fragment shader as a part of a graphics pipeline. Our 182
 183 network is compact and fast enough to run inside the GPU without using any 183
 184 interface from CUDA API. Since it is not possible to address a variety of scenes, 184
 185 we have traded generalization for specialization and fast training time. Because 185

186 of the short training time, our network can fit the cases where it is expensive to 186
 187 generate the animation for a given scene at high quality. 187

188 9 Future Works 188

189 Our work is one of its kind to run a denoising network inside a fragment shader. 189
 190 From this work, we notice that running a neural network requires extensive 190
 191 optimization and we encourage the readers to pursue this path. If it is not 191
 192 possible to run on the fragment shader, one can optimize the graphics pipeline 192
 193 by introducing CUDA calls in between. By doing so, one can run large models 193
 194 that can generalize to a large variety of scenes. Further, we also noticed that to 194
 195 achieve a significant performance gain, one shouldn't just be limited to vision 195
 196 algorithms but can also introduce the concepts from graphics. One such direction 196
 197 is passing the ray direction along with normals and the input image. 197

198 References 198

- 199 1. Realistic image synthesis, <https://graphics.cg.uni-saarland.de/courses/ris-2021/> 199
 200 [index.html](https://graphics.cg.uni-saarland.de/courses/ris-2021/index.html) 200
- 201 2. Chaitanya, C.R.A., Kaplanyan, A.S., Schied, C., Salvi, M., Lefohn, A., 201
 202 Nowrouzezahrai, D., Aila, T.: Interactive reconstruction of monte carlo image se- 202
 203 quences using a recurrent denoising autoencoder. *ACM Trans. Graph.* **36**(4), 1–12 203
 204 (2017) 204
- 205 3. Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with condi- 205
 206 tional adversarial networks. *CVPR* (2017) 206
- 207 4. Müller, T., McWilliams, B., Rousselle, F., Gross, M., Novák, J.: Neural importance 207
 208 sampling. *ACM Trans. Graph.* **38**(5), 1–19 (2019) 208
- 209 5. Sitzmann, V., Martel, J., Bergman, A., Lindell, D., Wetzstein, G.: Implicit neural 209
 210 representations with periodic activation functions. *Advances in Neural Information* 210
 211 *Processing Systems* **33** (2020) 211
- 212 6. Zhang, K., Zuo, W., Chen, Y., Meng, D., Zhang, L.: Beyond a gaussian denoiser: 212
 213 Residual learning of deep cnn for image denoising. *IEEE Transactions on Image* 213
 214 *Processing* **26**(7), 3142–3155 (2017) 214

215 A Dataset 215

216 We show a sample of our input (Fig. 6), normal (Fig.7) and output (Fig. 8) 216
 217 images rendered using our ray tracing algorithm. 217

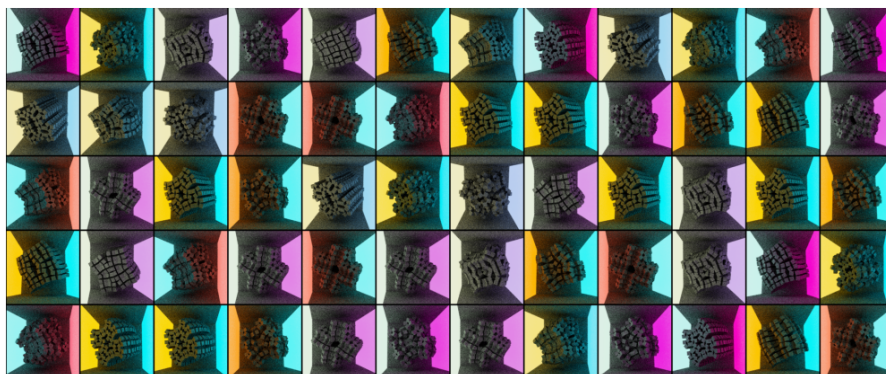


Fig. 6: Input Images

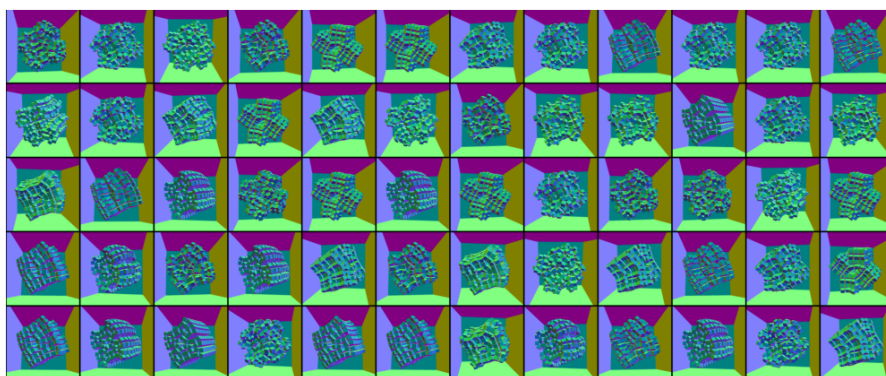


Fig. 7: Normal Images

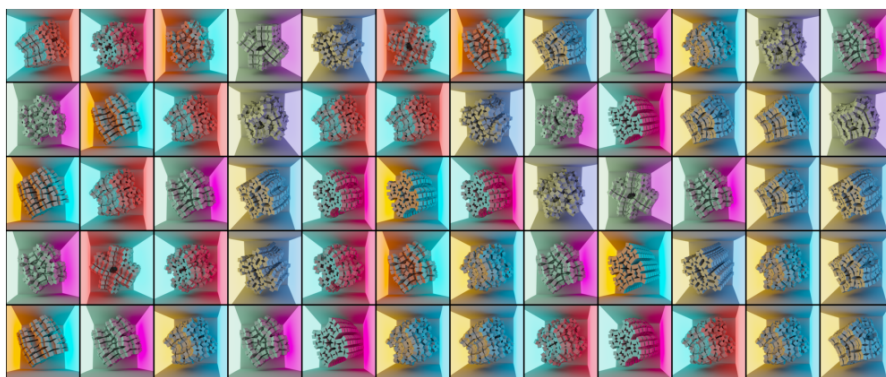


Fig. 8: Output Images