# Robust Average Networks for Monte Carlo Denoising

# JAVOR KALOJANOV, WētāFX, New Zealand KIMBALL THURSTON, WētāFX, New Zealand

We present a method for converting denoising neural networks from spatial into spatio-temporal ones by modifying the network architecture and loss function. We insert Robust Average blocks at arbitrary depths in the network graph. Each block performs latent space interpolation with trainable weights and works on the sequence of image representations from the preceding spatial components of the network. The temporal connections are kept live during training by forcing the network to predict a denoised frame from subsets of the input sequence. Using temporal coherence for denoising improves image quality and reduces temporal flickering independent of scene or image complexity.

## CCS Concepts: $\bullet$ Computing methodologies $\rightarrow$ Ray tracing; Image processing.

Additional Key Words and Phrases: denoising, neural networks, ray tracing, monte carlo

#### **ACM Reference Format:**

Javor Kalojanov and Kimball Thurston. 2025. Robust Average Networks for Monte Carlo Denoising. 1, 1 (April 2025), 7 pages. https://arxiv.org/abs/2310. 04080

# 1 INTRODUCTION

Deep Convolutional Neural Networks (CNNs) have been recently established as state-of-the-art techniques for image filtering in various areas. In this paper, we discuss denoising images generated via Monte Carlo path tracing in Visual Effects (VFX) production. Informally, denoising techniques reduce error by averaging pixel values in space and time, and the challenging aspects are balancing noise reduction with loss of detail due to blurring across image features. The main benefit of using CNNs for denoising is the large amount of filters with trainable weights that can provide very good feature preservation. This is especially important for VFX, where visual realism and image fidelity is often driven by the amount of detail in the rendered images. In this paper, we propose a method that improves temporal averaging while preserving image detail. We show that the technique applies to denoising complex imagery that contains interactions between rendered elements such as skin, water and hair.

Our method is targeted for production, where flickering, either introduced or left untreated by the denoising method, is not acceptable. We incorporate temporal information into our neural networks via recurrent components that interpolate between the latent representations of consecutive frames. These components consider a sequence of preceding and succeeding frames and convert spatial

https://arxiv.org/abs/2310.04080

kernel-predictive neural networks into spatio-temporal ones. The main benefits of kernel-predictive networks for our target application are the robustness against color shifts, and the ability to reuse denoising weights across multiple image buffers (AOVs). We furthermore train without ground truth reference images, which can be prohibitively expensive to generate in large enough quantities in a production setting.

The contributions of this paper revolve around a set of convolutional network layers that we call Robust Average (RA) blocks. These blocks perform interpolation and outlier removal in the time dimension. They operate on the network's internal state and can be inserted at arbitrary depth. The resulting networks have increased (but fixed) receptive field in time and better incorporate information from the entire frame sequence. The balance between spatial and temporal information used for denoising is learned from the data. This can lead to the networks degenerating to spatial ones during training. To prevent this, we warp each frame inside the temporal window of the network to convert it, as much as possible, to an estimate of the frame being denoised (the central frame). In order to further increase temporal contributions in the absence of temporal references during training, we propose a loss reformulation which forces the network to estimate the output frame from subsets of the input sequence that exclude the central frame. We thereby convert spatial loss functions into temporal ones without losing the ability to penalize warping artifacts.

We also propose to use thresholds on the activations of the kernel predictive layers. On one hand, this allows for user control of the influence of the denoising network on the input images. On the other, it prevents artifacts created by overestimated kernel weights in high-intensity regions of the image.

## 2 RELATED WORK

Kernel-predictive neural networks Our work is closely related to [Bako et al. 2017; Vogels et al. 2018; Zhang et al. 2021, 2023], since we use kernel-predictive denoising networks, i.e., the output of the network are per-pixel 2D (or 3D) kernels. Instead of predicting the colors the denoised image directly, the kernels are applied (via dot product) on the input image to denoise it. The main difference between our method and these related works is in the neural network components and loss for temporal denoising, but the remaining techniques are compatible with our method. Mildenhall et al. [2018] employ a kernel predictive network for denoising burst photographs (instead of rendered images). Their loss generalization is similar to ours in that it introduces loss terms for the individual images in the burst. We instead pair images before and after the central frame to combat occlusion and disocclusion, and do not downscale temporal terms while training. Instead of filtering the image after it is rendered, sample-based methods for error reduction such as [Gharbi et al. 2019] correct pixel samples before aggregating them into a framebuffer. We scope this paper on using temporal information for

Authors' addresses: Javor Kalojanov, jkalojanov@wetafx.co.nz, WētāFX, New Zealand, Wellington; Kimball Thurston, kthurston@wetafx.co.nz, WētāFX, New Zealand, Wellington.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s). © 2023 Copyright held by the owner/author(s).



Fig. 1. A sketch of the recurrent ResNet network model used for evaluation. The network consists of 24 residual blocks, RRA blocks added after blocks 3,6,9,12,15 and 24, and skip connections after the first 4 RA blocks. The internal dimension of all convolutional layers is 80, except for the kernel predictive layers, which are  $5 \times 5 \times 5 = 125$  large.

denoising, and leave out discussion on techniques that run as the individual frames are rendered like [Cho et al. 2021; Yu et al. 2021].

**Recurrent neural networks for denoising** Introducing recurrent connections in denoising neural networks has been done, for example by Chaitanya et al. [2017] and Hasselgren et al. [2020]. The purpose of these connections is the same as in our work: to force the network to use temporal information when denoising. Our approach differs in that it is bidirectional, i.e., we take advantage of succeeding as well as previous frames. The remaining differences with these works are due to the target applications. We trade off network complexity to achieve higher image quality.

**Non-neural denoising** Denoising techniques which do not rely on neural networks can also be used in production. Methods like Non-Local Means Filtering [Buades et al. 2005] or Guided Filtering [He et al. 2013], use model-based approaches to derive per-pixel weights used for spatial denoising. Other similar techniques were extended to include temporal information [Dabov et al. 2007] in addition to the larger body of work for film restoration which can be repurposed for Monte Carlo denoising [Kokaram 2014]. We use such a denoiser to generate pseudo-reference images, but leave further discussion on non-neural denoising outside the scope of this paper. Finally, we recommend the survey on adaptive sampling and denoising by Zwicker et al. [2015] for the readers looking for an overview of the field.

## 3 ROBUST AVERAGE BLOCKS



Fig. 2. Robust Average block for a sequence of length 5. We exclude the first frame, average the remaining frames, and interpolate between the average and the excluded frame. This is repeated (recurrently) until each frame of the sequence has been interpolated with the robust average of the other frames.

In the following, we introduce the Robust Average (RA) blocks. We use them to force the network to denoise in time via latent space

Publication date: April 2025.

interpolation. Intuitively, averaging the internal state across several noisy estimates of the same image will result in similar network output and reduce temporal noise. However, we did not obtain satisfactory results using a straight forward implementation that just averages network state at various depths. The networks we tried were prone to creating artifacts (e.g. light leaks). Our understanding is that averaging in latent space reduced the ability of the network to remove outliers caused by noise or warping errors.

We address this problem in two ways. First, we let the network independently predict interpolation weights and use them to blend each frame of the sequence with an average of the remaining frames. This enables outlier removal and preserves the ability of the network to represent an identity transformation. The latter is important for networks like ResNet [He et al. 2016]. Second, given enough frames, we use an averaging operation that is robust to outliers: given a set of values  $S \subset \mathbb{R}$  with  $|S| \ge 3$ , we call *robust average* computing the mean after discarding the minimum and maximum:

$$ravg(S) = \frac{1}{|S| - 2} \left( \sum_{S} x - \max(S) - \min(S) \right)$$
(1)

The Robust Average block shown in Figure 2 operates on a fixed length window of frames centered at the current frame. First, we exclude one frame from the sequence, and compute robust average of the remaining frames. Then a set of convolutional layers is used to estimate per-pixel weights, which are then used to blend the excluded frame with the average. This is repeated once for each frame in the sequence. Note that the operation of the RA block is recurrent since each successive step uses the result of the previous step.

The main advantage of the RA blocks compared to averaging, is that the balance between using temporal and spatial information for each pixel is learned instead of fixed. The network can reduce or increase the amount of temporal interpolation per pixel at different depths. In addition to noise reduction, this helps with correcting errors due to inaccurate warping.

We use motion vectors generated during rendering to warp frames on each side of the central frame. This includes auxiliary output buffers such as the albedo and normals. We also compute confidence in the warp from the motion vectors the warped albedo and normal buffers. In order to avoid temporal artifacts without decreasing weights along temporal connections, we mix in pixel values from the central frame inversely proportional to the confidence score for each side frame. This helps prevent the network from degenerating to a spatial only model during training.

The number of frames used as input is fixed because the blocks are bidirectional. However, we are able to extend the temporal range by running multiple denoising passes. This is possible for networks with high detail preservation and low bias such as sharpening.

As the blocks can be added at any depth, they allow filtering in time in the very first layers. This is in contrast with approaches that denoise in space first, and only then introduce temporal connections. In the evaluation, we show evidence that our method helps retain small scale image features that are blurred by a network with a deep spatial component followed by a temporal part. Furthermore, adding RA blocks at multiple depths reduces the size and expressive power of the spatial-only components of the network. This favors denoising in the time dimension, and hinders the model from degenerating into a spatial-only during training.

#### 4 THRESHOLDED KERNEL PREDICTION

A common choice for activation for the output layer for kernelpredictive networks is the softmax function:

$$\sigma(w_i) = \frac{e^{w_i}}{\sum_j e^{w_j}} \tag{2}$$

This has the advantage that the weights  $w_i$  are normalized to sum up to 1, however the predicted weights are strictly positive. With non-zero kernel weights, every pixel in the kernel support contributes to the final estimate. We found that this can lead to artifacts when high intensity pixels are present in the kernel support. The high intensity pixels can be fireflies, lightsources, or specular reflections of lightsources, and they amplify inaccuracies in the predicted weights.

For intuition, consider a value  $y \in \mathbb{R}$  that is reconstructed from a set of estimates  $x_i$ , using kernel weights  $w_i$ . If there are weights  $w_j = 0$  participating in the reconstruction, it holds that y can be reconstructed without additional error even if  $x_j$  becomes an outlier, e.g., a firefly. Note that in the case of Monte Carlo denoising, the objective is to reconstruct the correct pixel values from random estimates. It is therefore likely that the estimates for the same pixel regions sometimes contain fireflies, i.e., pixels with comparatively larger error.

We enable zero kernel weights by using use thresholded ReLUs followed by normalization in the kernel outputting layer of our networks. More specifically, our activations are:

$$\eta_t(w_i) = \frac{\max(0, w_i - t)}{\sum_{i=1}^K \max(0, w_j - t)}$$
(3)

with  $t < \frac{1}{K}$  to prevent zeroing out normalized kernels with equal weights throughout the domain. If not mentioned otherwise, we report results using  $t = \frac{1}{2K}$ . If all weights are below *t*, we replace the resulting kernel with identity. Note that thresholds can be used with any activation, including softmax, at the expense of an additional normalization if it is desirable for the kernel weights to sum up to 1.

Increasing the threshold can be used to correct loss of detail in denoised images at the expense of reintroducing noise from the input (see Figure 3). This accommodates fine-tuning the model after training based on artist feedback, if done once before the trained model is deployed. Alternatively, the threshold value can be user controlled and adjusted per image. Finally, note that in our case increasing the kernel threshold usually results in less temporal



Fig. 3. Increasing the value of the kernel threshold t adjusts the influence of the denoiser on the image. The figure shows denoised images and scaled difference to the noisy input image. ©20th Century Studios / Walt Disney Studios Motion Pictures

utilization. This is because for each weight for the central frame, we have multiple weights for remaining frames.

## 5 SPATIAL-TO-TEMPORAL LOSS CONVERSION

Denoising via convolutional neural networks is commonly formulated as follows. Let  $X := \{x^1, \ldots, x^N\}$  and  $Y := \{y^1, \ldots, y^N\}$  be two corresponding sets of noisy and reference images, and  $l(\cdot, \cdot) :$  $\mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}$  a per-pixel image loss (e.g. mean absolute error). Let  $f_{\theta}(\cdot) : \mathbb{R}^n \to \mathbb{R}^n$ , be a convolutional neural network parameterized by  $\theta$  that denoises one image at a time. The optimal set of parameters for the network is determined via supervised learning:

$$\underset{\theta}{\arg\min} \frac{1}{N} \sum_{i=1}^{N} l\left(f_{\theta}(x^{i}), y^{i}\right) \tag{4}$$

In the following, we set N = 1 to simplify notation.

We propose to generalize the procedure from a spatial into a temporal one by asking the neural network to predict the reference image from any subset of the input image sequence. Let  $X := \{x_{-k}, \ldots, x_k\}$  be a set of consecutive noisy frames and  $y_0$  the reference for the center frame. Assume that for each  $i \in [-k, k]$ , there is a motion compensating transformation (image warp) that aligns features from  $x_i$  to  $x_0$ , and let  $X' := \{x'_{-k}, \ldots, x'_k\}$  be the set of transformed images. Consider a family of neural networks  $G_{\theta} := \{g^i_{\theta}(\cdot)\}$  such that for each  $\hat{X} \subset X'$ , there is a  $g^i_{\theta}(\hat{X}) : \cdot \to \mathbb{R}^n$ . The minimization problem then becomes:

$$\arg\min_{\theta} \sum_{\hat{X} \in \mathcal{P}(X')} l\left(g_{\theta}^{i}(\hat{X}), y_{0}\right)$$
(5)

where  $\mathcal{P}(X')$  is the power set of X'. Given  $g_{\theta}(\cdot) : \mathbb{R}^{n(2k+1)} \to \mathbb{R}^n$ , a spatio-temporal kernel-predictive neural network for X', a set  $G_{\theta}$ is easy to construct by zeroing out predicted kernel weights for the excluded frames and re-normalizing the remaining weights. In other words, we assume that the neural networks predict an image from a set of estimates instead of a single one. We treat noise, warping error, and other biases as deviations from the target image and rely on the network to correct for them.



Fig. 4. Spatial to temporal loss conversion for a kernel-predictive neural network. Here, we add two temporal loss terms that force the network to predict the center frame from pairs of one previous and one subsequent frame. We use the same reference frame  $ref_0$  in each term.

In practice, we do not construct a network and loss term for each subset of the input sequence to reduce memory and computational complexity. Instead, as shown in Figure 4, we train a network with temporal window of 5 with loss terms for:

- the center frame
- a weighted average of the first and fourth frame
- a weighted average of the second and fifth frame

we compare against the center reference frame in each case. We select pairs of frames before and after the center frame, to improve the chance of retaining image features that are occluded after or revealed during the central frame. We only use one term containing the center input frame in order to promote the use of temporal information.

Note that our approach works with a spatial loss function instead of depending on measured deviations in time. Because of that, we are able to train without temporal reference images, while still penalizing error using the center frame. As a result, the network learns to reduce error coming from the off-center frames that is caused by either noise, motion, or warping. Temporal artifacts can be further penalized by adding a loss term that compares the combined result to the reference frame. Note that Vogels et al. [2018] train the temporal component of their denoiser using just the combined loss term. In our case: end to end training of a spatio-temporal network, this steered the network to use less temporal information. Therefore, we only used a downweighted global loss term to post-train models if they introduce temporal artifacts.

Our conversion works by adding spatial loss terms for output images generated via temporal averaging, independent of the choice of spatial loss  $l(\cdot, \cdot)$ . We tested this by training with various losses including mean absolute error ( $L_1$ ), symmetric mean absolute percentage error [Vogels et al. 2018], and the perceptual loss in [Gatys et al. 2016; Johnson et al. 2016]. We did not observe substantial differences in the learning behavior in the different cases. However, networks trained with a perceptual loss had better detail preservation compared to the remaining losses. We therefore train using

$$l(x,y) := \left( L_{VGG_{5,4}}(x,y) + L_{VGG_{3,4}}(x,y) \right) L_{SMAPE}(x,y)$$
(6)

where  $L_{VGG_{a,b}}(\cdot, \cdot)$  is a  $L_1$  distance in the latent space of block a, layer b of VGG19 [Simonyan and Zisserman 2014] trained on ImageNet [Deng et al. 2009]. The choice of addition and multiplication for the separate losses is to take into account the implicit scaling of the terms. We sum terms with similar magnitude and multiply otherwise. We opted for SMAPE instead of  $L_{VGG_{2,2}}(\cdot, \cdot)$  as a "feature" loss, because the latter led to aliasing artifacts. A simpler  $L_1$  loss can be used instead of SMAPE depending on color space and pixel value range and distribution.

#### 6 TRAINING DATA

Creating large bodies of training data for machine learning methods can be challenging in a VFX environment because of the effort required to produce imagery. It is often the case that generating a fully converged reference images or image sequences for denoising is prohibitively expensive. Based on the results demonstrated by Lehtinen et al. [2018], we do not train using ground truth reference images. Instead, we render images at different noise levels and use them to construct input and target pairs for training. In addition to the noisy references we also use pseudo-references generated by an existing denoiser. The latter can be a previous version of the denoising network, or in our case, a separate method based on using a modified TV-L1 [Wedel et al. 2009] optical flow method initialization from the renderer-provided motion vectors, confidence based robust averaging, Non-Local Means Filtering [Buades et al. 2005], and Guided Filtering [He et al. 2013].

The training data we use for evaluation has been generated by instrumenting a production path tracer to render images at 3 noise levels:

- *noisy* RGB image using *n* samples per pixel
- half quality RGB image using  $\left|\frac{n}{e}\right|$  samples per pixel<sup>1</sup>
- low quality RGB image using 4 samples per pixel

The number of samples taken varies per image and per pixel. The noise level of the images is typically low enough to deliver final quality after denoising and compositing. During training, we pick the pairs of noise levels randomly (including pseudo-references) and denoise from more noisy to less noisy.

We train on data from approximately 11000 rendered sequences of length of up to several hundred frames. From these, we randomly sample 10 disjoint sets of 16000 to 20000 tiles of dimensions  $128 \times 128 \times 11$ .

### 7 EVALUATION

In this section, we compare our network extension via RA blocks to an alternative spatio-temporal architecture used in [Vogels et al. 2018] which is, to our knowledge the most similar alternative in terms of target application and network structure. We show that it is possible to obtain at least similar results for a single image (Figure 5, Table 2), while using larger temporal weights (Figure 7). Note that our networks do not contain temporal connections other than the

<sup>&</sup>lt;sup>1</sup>We use an irrational fraction to ensure that subsequences of quasi-random samples remain low discrepancy.



Fig. 5. Our denoiser delivers similar image quality with slightly better image details compared to a baseline (tKPCN) network without RA blocks. A spatial UNet keeps detail by preserving noise, and sacrificing temporal stability. ©20th Century Studios / Walt Disney Studios Motion Pictures

ones in the RA blocks (even in the final kernel-predictive layers). The network being able to balance weights between the central and non-central frames is thanks to these components.



Fig. 6. Temporal breakdown of denoising weights. Noisy, denoised image, and contributions for the input sequence of 5 frames. The outermost frames have significantly smaller total contribution, however the minimum and maximum weights for each frame are 0 and 1. ©20th Century Studios / Walt Disney Studios Motion Pictures

The Robust Average blocks can be used together with a number of convolutional neural network types. We tested this on ResNet and UNet architectures, as well as the multiscale ResNet variant in [Vogels et al. 2018]. The trained models performed sufficiently well to be used for noise reduction in a commercial rendering pipeline. The key contributing factors for this were the absence of introduced temporal flickering and the detail/edge preservation (see Figure 5). Compared to noise reduction via sampling, the denoised static images compared favorably to frames with  $\frac{1}{4}$  estimated pixel error. For final quality renders, we were not able to achieve lower temporal noise than a denoised sequence of frames by increasing pixel samples.

From the topologies we tested, the best performing network is a ResNet (Figure 1) without upscaling and downscaling components. The network denoises from a 5 frame temporal window using  $5 \times 5 \times 5$  kernels (height, width, time), and has 24 million trainable parameters.

The input for the network is pre-warped using motion vectors and consists of color, albedo, normals, and confidence. The latter is a single channel image representing confidence in the warp derived from the motion vectors.

Figure 6 is an example from a scene with relatively slow camera or character movement. The average pixel contribution from the central frame is smaller than the combined weight of the remaining frames, meaning that the denoiser uses more temporal than spatial information. The outermost frames have significantly smaller average contribution. This is expected, given that in a large random training set the temporal coherency and image warp quality decreases fast with the distance to the central frame. On the other hand, the minimum and maximum weights for each frame are 0 and 1. In other words, the network can completely replace pixel values from the central frame from the surrounding frames.



Fig. 7. Spatial vs temporal contributions for our denoiser compared to a baseline (tKPCN) using a temporal loss and network like [Vogels et al. 2018] instead of our loss and RA blocks. ©20th Century Studios / Walt Disney Studios Motion Pictures

When compared against a baseline, the RA blocks and temporal loss improve temporal utilization. We remove all RA blocks from the network, and instead use the temporal module and loss proposed in [Vogels et al. 2018]. As with their work, the spatial parameters are still shared across the sequence of input frames. This network, which we call tKPCN, can train to denoise in time, however strongly favors the center frame if trained with a standard loss. Figure 7 shows a crop from the same scene as Figure 6, which exhibits good temporal coherency. The baseline method uses smaller weights for pixels from the non-center frames both on average: (0.0063, 0.1158, 0.1503, 0.007) compared to (0.0418, 0.2486, 0.2420, 0.0436); and in terms of maximum weights: (0.1163, 0.9636, 0.8337, 0.1122) compared to 1.

We break down the influence of the network structure and the loss function in Table 1. Both our network and the tKPCN baseline can be trained to favor the central frame when using a "spatial" loss term measuring the difference to the reference image. This is not surprising, since there is no incentive for the network to use the pixel values from the surrounding frames. The RA blocks slow down the rate at which the network converges to a mostly spatial network. On the other hand, the temporal loss reformulation can steer both our network and the baseline to use larger temporal weights, however we were not able to train a tKPCN network with a good temporal

						frame id
		-2	-1	0	1	2
RA +	avg	0.0197	0.1671	0.6409	0.1514	0.0209
spatial loss	max	0.2013	0.3223	0.9966	0.3199	0.1950
R∆ ⊥	21/0	0.0580	0 1402	0 5694	0 1531	0 0792
101 +	avg	0.0050	0.1402	1 0000	0.1331	0.0772
our loss	max	0.9959	0.9982	1.0000	0.9925	0.9582
tKPCN +	avg	0.0021	0.1387	0.7200	0.1379	0.0013
spatial loss	max	0.1505	0.9956	1.0000	0.9975	0.3985
tKPCN +	avg	0.0387	0.2994	0.4076	0.2123	0.0421
our loss	max	0.9831	0.9999	1.0000	0.9987	0.9964
tKPCN + spatial loss tKPCN + our loss	avg max avg max	0.0021 0.1505 0.0387 0.9831	0.1387 0.9956 0.2994 0.9999	0.7200 1.0000 0.4076 1.0000	0.1379 0.9975 0.2123 0.9987	0.0013 0.3985 0.0421 0.9964

Table 1. Ablation tests on a frame from the Payakan scene (see Figure 8). RA blocks provide better temporal utilization with a spatial-only loss compared to the non-RA network trained for 11 epochs (35k training samples each). Both networks converge to using predominantly spatial information with further training with the spatial loss. Our loss reformulation does improve temporal utilization for the tKPCN network.

utilization, no artifacts, and better detail preservation compared to our RA models (see Figure 8). In our experience, it was easier to train a tKPCN model to remove more noise, while RA networks are better at edge preservation. For this comparison, we trained the both networks with identical losses and procedure: an initial pass with the temporal loss followed by pass with added downscaled spatial loss term. We report results using the checkpoints with the best validation loss for both networks.



Fig. 8. Image comparison using our network and a tKPCN model trained with our temporal loss reformulation (from Table 1). The tKPCN model removes more noise, while our network preserves more detail, especially on parts of the image where the creature is submerged underwater. ©20th Century Studios / Walt Disney Studios Motion Pictures

The test in Figure 6 shows the network predicting predominantly spatial weights for the hand of the character, which moves fast, and predominantly temporal weights on the character face and body which remain mostly static. Reverting to spatial denoising in the absence of temporal coherency reduces or avoids image artifacts such as ghosting and light leakage which occur when interpolating between non-corresponding pixel values in time.



Fig. 9. Denoising the specular component of a render with the kernel weights from the RGB output can work well in practice. Some small details in the specular component of the image are still recognizable after denoising. ©20th Century Studios / Walt Disney Studios Motion Pictures

Figure 9 shows the level of detail we are able to obtain when predicting weights from an image and applying them only to the specular component of the image. While not mathematically correct, this works well enough in practice and recuces computational effort when denoising renders with a lot of auxiliary buffers.

In Table 2 we compare peak signal-to-noise ratio (PSNR) and structural similarity index measure (SSIM) on our datasets. We aggregate statistics from a set of test images that have not been used for training or validation. In the test we denoise from the half quality RGB image to a pseudo-reference denoised with a production denoiser. In VFX, the quality of the denoised images is judged by expert users, which makes the assessment perceptual rather than statistical. The test results here do match our assessment of the relative image quality: both models deliver similar results, and our network is consistently better at keeping image detail.

	min	avg	PSNR max	min	avg	SSIM max
input	42.4898	44.6739	46.8580	0.9880	0.9900	0.9920
tKPCN	42.0127	44.9809	47.9492	0.9877	0.9910	0.9944
ours	45.3110	48.0516	50.7922	0.9946	0.9957	0.9969

Table 2. Aggregate SSIM and PSNR comparison between our model and the temporal baseline network. The statistics are from 768 test samples from scenes that did not participate in training.

#### 8 LIMITATIONS

We show a reconstruction from a low sample count input in Figure 10. While the reconstruction is plausible, details like the shadow of the railing are missing. In this case, the network is not able to predict image detail revealed by further sampling. With our data and training, we favor preserving detail and sacrifice the ability to reconstruct from very low initial sample counts. Note that the reconstruction quality from lower sample counts can be accetable for certain use cases, e.g., previews.

The main disadvantage of extending networks via Robust Average blocks is the increased network graph complexity. Mixing in temporal components at different network depths makes it hard to



Fig. 10. Validation image with low sample count. Our network fails to reconstruct small scale shadows present in the pseudo-reference image.

train spatial and temporal parts of the network separately. In our case, we had to compensate for this by training with batch size of 1 or 2 on a GPU with 24GB of memory. Memory restrictions can also hinder increasing the temporal window of the neural network. Our results suggest that there is not much potential in expanding beyond 5 frames without improving the warping quality. We instead increase the temporal window by running multiple denoising passes if necessary.

The large network graph also impacts inference time negatively. Our implementation in Python using Tensorflow 2.5 [Abadi et al. 2015] needs approximately 2 min to denoise a 2k stereo image on a high-end GPU (Nvidia RTX A5000). This includes loading the network from a file and compiling the network graph, warping the input frames and auxiliary buffers, kernel prediction on tiles of size  $728 \times 548$ , kernel thresholding, warping output buffers, and kernel application. In practice, we amortize the network initialization by denoising in batches of frames, and optimize for efficiency by running on CPUs with low thread counts, which results in longer inference times.

Finally, denoising in time assumes temporal coherency in the input sequence. If coherency in time is not present, we rely on the network to revert to using spatial information. Our networks can, for example, preserve shadows moving over static surfaces, however we currently use a spatial denoiser for renders of fire and explosions.

#### 9 CONCLUSION

In this paper, we introduce a method for designing and training temporal denoising neural networks. We introduced a new recurrent block pattern, an improved training technique, and an alternate method for using temporal information in image sequences. We believe that our work helps make a small step towards higher image fidelity and rendering efficiency in VFX, and hope that this paper will help and inspire further research in the field.

#### ACKNOWLEDGMENTS

The authors thank Mihail Moskov for consulting on deep learning and image processing. We thank Joe Letteri, Sam Cole, and Francois Sugny for guidance on imaging and compositing, and Youngbin Park for pipeline integration. This work has received a Technical Achievement Award from the Academy of Motion Picture Arts and Sciences.

### REFERENCES

- Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martín Wicke, Yuan Yu, and Xiaoqiang Zheng. 2015. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems. https://www.tensorflow.org/
- Steve Bako, Thijs Vogels, Brian McWilliams, Mark Meyer, Jan Novák, Alex Harvill, Pradeep Sen, Tony Derose, and Fabrice Rousselle. 2017. Kernel-predicting convolutional networks for denoising Monte Carlo renderings. ACM Trans. Graph. 36, 4 (2017), 97–1.
- A. Buades, B. Coll, and J.-M. Morel. 2005. A non-local algorithm for image denoising. In CVPR'05, Vol. 2. 60–65 vol. 2. https://doi.org/10.1109/CVPR.2005.38
- Chakravarty R. Alla Chaitanya, Anton S. Kaplanyan, Christoph Schied, Marco Salvi, Aaron Lefohn, Derek Nowrouzezahrai, and Timo Aila. 2017. Interactive Reconstruction of Monte Carlo Image Sequences Using a Recurrent Denoising Autoencoder. ACM Trans. Graph. 36, 4 (2017).
- In-Young Cho, Yuchi Huo, and Sung-Eui Yoon. 2021. Weakly-Supervised Contrastive Learning in Path Manifold for Monte Carlo Image Reconstruction. ACM Trans. Graph. 40, 4, Article 38 (jul 2021), 14 pages. https://doi.org/10.1145/3450626.3459876
- Kostadin Dabov, Alessandro Foi, Vladimir Katkovnik, and Karen Egiazarian. 2007. Image Denoising by Sparse 3-D Transform-Domain Collaborative Filtering. IEEE Transactions on Image Processing 16, 8 (2007), 2080–2095. https://doi.org/10.1109/ TIP.2007.901238
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition. IEEE, 248–255.
- Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. 2016. Image Style Transfer Using Convolutional Neural Networks. In 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). 2414–2423. https://doi.org/10.1109/CVPR.2016.265
- Michaël Gharbi, Tzu-Mao Li, Miika Aittala, Jaakko Lehtinen, and Frédo Durand. 2019. Sample-based Monte Carlo Denoising Using a Kernel-splatting Network. ACM Trans. Graph. 38, 4 (2019), 125:1–125:12.
- J. Hasselgren, J. Munkberg, M. Salvi, A. Patney, and A. Lefohn. 2020. Neural Temporal Adaptive Sampling and Denoising. *Computer Graphics Forum* 39, 2 (2020), 147–155.
- Kaiming He, Jian Sun, and Xiaoou Tang. 2013. Guided Image Filtering. IEEE transactions on pattern analysis and machine intelligence 35 (06 2013), 1397–1409.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In CVPR. 770–778. https://doi.org/10.1109/CVPR.2016.90
- Justin Johnson, Alexandre Alahi, and Li Fei-Fei. 2016. Perceptual losses for real-time style transfer and super-resolution. In European conference on computer vision. 694– 711. http://arxiv.org/abs/1603.08155
- A.C. Kokaram. 2014. Motion Picture Restoration: Digital Algorithms for Artefact Suppression in Degraded Motion Picture Film and Video. Springer London.
- Jaakko Lehtinen, Jacob Munkberg, Jon Hasselgren, Samuli Laine, Tero Karras andMiika Aittala, and Timo Aila. 2018. Noise2Noise: Learning Image Restoration without Clean Data. CoRR abs/1803.04189 (2018). arXiv:1803.04189 http://arxiv.org/abs/ 1803.04189
- Ben Mildenhall, Jonathan T. Barron, Jiawen Chen, Dillon Sharlet, Ren Ng, and Rob Carroll. 2018. Burst Denoising with Kernel Prediction Networks. In CVPR. 2502–2510. https://arxiv.org/abs/1712.02327
- Karen Simonyan and Andrew Zisserman. 2014. Very Deep Convolutional Networks for Large-Scale Image Recognition. https://arxiv.org/abs/1409.1556
- Thijs Vogels, Fabrice Rousselle, Brian McWilliams, Gerhard Röthlin, Alex Harvill, David Adler, Mark Meyer, and Jan Novák. 2018. Denoising with kernel prediction and asymmetric loss functions. *ACM Trans. Graph.* 37, 4 (2018), 1–15.
- Andreas Wedel, Thomas Pock, Christopher Zach, Horst Bischof, and Daniel Cremers. 2009. An Improved Algorithm for TV-L1 Optical Flow. Springer-Verlag, Berlin, Heidelberg, 23–45.
- Jiaqi Yu, Yongwei Nie, Chengjiang Long, Wenju Xu, Qing Zhang, and Guiqing Li. 2021. Monte Carlo Denoising via Auxiliary Feature Guided Self-Attention. ACM Trans. Graph. 40, 6, Article 273 (dec 2021), 13 pages. https://doi.org/10.1145/3478513. 3480565
- Xianyao Zhang, Marco Manzi, Thijs Vogels, Henrik Dahlberg, Markus Gross, and Marios Papas. 2021. Deep Compositional Denoising for High-quality Monte Carlo Rendering. Computer Graphics Forum (2021). https://doi.org/10.1111/cgf.14337
- Xianyao Zhang, Gerhard Röthlin, Marco Manzi, Markus Gross, and Marios Papas. 2023. Deep Compositional Denoising on Frame Sequences. In Eurographics Symposium on Rendering (EGSR).
- M. Zwicker, W. Jarosz, J. Lehtinen, B. Moon, R. Ramamoorthi, F. Rousselle, P. Sen, C. Soler, and S.-E. Yoon. 2015. Recent Advances in Adaptive Sampling and Reconstruction for Monte Carlo Rendering. *Comput. Graph. Forum* 34, 2 (may 2015), 667–681.