# Deep Radiance Caching: Convolutional Autoencoders Deeper in Ray Tracing

Giulio Jiang[a], Bernhard Kainz[a,*]

[a]*Imperial College London, SW7 2AZ, London*

## ABSTRACT

Rendering realistic images with global illumination is a computationally demanding task and often requires dedicated hardware for feasible runtime. Recent research uses Deep Neural Networks to predict indirect lighting on image level, but such methods are commonly limited to diffuse materials and require training on each scene. We present Deep Radiance Caching (DRC), an efficient variant of Radiance Caching utilizing Convolutional Autoencoders for rendering global illumination. DRC employs a denoising neural network with Radiance Caching to support a wide range of material types, without the requirement of offline pre-computation or training for each scene. This offers high performance CPU rendering for maximum accessibility. Our method has been evaluated on interior scenes, and is able to produce high-quality images within 180 seconds on a single CPU.

## 1. Introduction

Ray Tracing is capable of producing photo-realistic images virtually indistinguishable from real pictures. Progressive refinements on rendering algorithms, such as Bi-Directional Path Tracing (BDPT) [1] and Metropolis Light Transport [2] have increased the efficiency of rendering engines in scenarios in which light paths are difficult to evaluate due to the high amount of indirect lighting and Global Illumination. Complex lighting conditions are, however, still highly expensive to resolve, and most algorithms require long rendering times to reduce the noise from Monte Carlo sampling.

Biased methods have been implemented to produce convincing quality images at a fraction of the cost required by a ray tracer. An early method, Instant Radiosity [4], exploited the low rate of illumination change over diffuse surfaces to approximate global illumination by rendering the same scene many times using Virtual Point Lights sampled at locations reached by the main light sources to simulate secondary bounces.
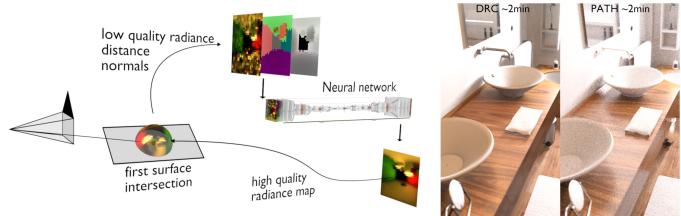


**Fig. 1. Our proposed Deep Radiance Caching (DRC) method evaluates a low quality radiance map, distance and normals at the first intersection, and uses a neural network to obtain a higher quality radiance map containing all indirect lighting. The *Bathroom* scene from PBRTv3's scenes [3] shows DRC compared to a same-time path traced image. DRC produces a noise-free result with convincing global illumination after 30 seconds, and progressively refines quality afterwards.**

Radiance Caching [5] focuses on accelerating rendering of glossy materials, by caching an optimized representation of the radiance received on a surface. This method enables view-dependant reflections to be rendered correctly.

Biased algorithms have evolved and recent research has attempted to use machine learning techniques to accelerate rendering of global illumination effects. The *Deep Illumination* [6] approach uses a GAN to translate diffuse albedos, normals and depth maps to a global illumination component layer and

---
*Corresponding author:
 e-mail:* giuliojiang@gmail.com (Giulio Jiang),
b.kainz@imperial.ac.uk (Bernhard Kainz)

obtained good results at predicting indirect illumination in real time for diffuse materials. The network requires specific training for each scene to be rendered, but is able to extrapolate and adapt to dynamic objects and newly introduced shapes.

To overcome the limitations of diffuse-only materials of current machine learning renderers, and to make the renderer easily usable without the need for offline training, we attempt a more general approach to resolving global illumination. Based on the estimation and caching of radiance maps for indirect illumination only, we combine Radiance Caching and a neural network that reconstructs and denoises radiance maps. Indirect lighting is the primary cause of strong noise in Monte Carlo rendering. This noise is much more difficult to clear than first bounce lighting due to the high dimensionality of the paths. Rendering can therefore be accelerated by predicting approximate but noise-free radiance maps that can be interpolated and used to obtain the indirect illumination component of the final image. Thanks to the slowly changing nature of indirect illumination, artifacts and bias are not very visible in the general case, while the overall predicted global illumination looks convincing and noise-free.

In contrast to full image level methods, Deep Radiance Caching (DRC) as outlined in Figure 1 uses path tracing to find the first intersection in the scene. A Convolutional Autoencoder predicts at this intersection a high quality, geometry-dependent radiance map from a path traced map rendered at just 1 sample per pixel, a depth and a normal map.

Using high quality radiance maps we can approximate all indirect illumination contributed from path tracing bounces beyond the first one. We use a neural network to predict high quality radiance maps for these approximations. Consequently, we accelerate the rendering of global illumination effects significantly. Our approach works with a wide range of material types, and does not require any offline pre-computation or per-scene training.

## 2. Related Work

### 2.1. Advanced denoising systems

Recent research uses image denoising and machine learning techniques for rendering and image generation purposes. The most immediate way to apply machine learning to graphics and ray tracing is to operate on the final rendered image level. These approaches take as input a scene rendered from the final camera's viewpoint, and attempt to output a transformation that results in higher visual quality, removal of noise, or the addition of effects.

Denoising of natural images is a core research area in Computer Vision. Numerous approaches have been proposed using, *e.g.*, Total Variation [7, 8], Non-local means filtering [9], dictionary learning [10] or recently, wavelet transformations of the contracting path of a convolutional neural auto-encoder with skip connections to gain a higher perceptive field with minimal computational costs [11].

Monte Carlo rendering methods like path tracing produce images with stochastic noise artifacts, thus image denoising algorithms are a natural fit for improving the final image quality. Denoising of Monte Carlo renderings has been intensely studied in Computer Graphics literate. For example, [12] use a first-order model with auxiliary buffers and a Non-local means nonlinear regression kernel to predict optimal filter parameters. [13] proposes an adaptive rendering method, which fits local polynomial functions to approximate the image and predicts the local optimal polynomial order with a multistage error estimation process.

One of the first neural network-based approaches to improve rendering quality [14] used a network to obtain filter parameters to denoise path traced images. The method relies on collecting primary features such as world coordinates, surface normals, texture values and illumination visibility from a ray tracer, from which secondary features are computed: gradients, mean and deviation. A Multi Layer Perceptron (MLP) uses the secondary features to output filter parameters. The method achieves good denoising quality, but requires the use of a modified ray tracer-based on PBRTv2 [15]. *Kernel-predicting convolutional networks for denoising Monte Carlo renderings* [16] and the recent extension *Denoising with Kernel Prediction and Asymmetric Loss Functions* [17] improve denoising networks by separating the processing of diffuse and specular layers with Convolutional Neural Networks (CNN). This approach is also highly applicable to existing ray tracers, as many can natively output separate layers for diffuse, specular, z-buffer, mist and others. The separation of layers allows the network to better handle the high dynamic range of input values. The introduction of logarithmic transforms to the input data further increases performance. The CNN shows excellent performance with moderate noise.

A method proposed by Chaitanya et al [18] differs from the other examples as it does not take inputs with a small number of samples per pixel, but with exactly one sample. By retaining information from previous frames in an animation, the used Convolutional Autoencoder is capable of reconstructing missing details given a very small amount of information. This method is particularly suited for fast, interactive animation previews. A variant of this approach has been proposed by [19] to tackle the problem of image reconstruction from sparsely, non-uniformly sampled renderings.

Very recent approaches focus on different aspects of image-level denoising, like multi-variate pixel values with color information from different scene depths [20], novel denoising architecture to relate permutation variant samples directly to the output image through splatting [21], and [22] who use the power of generative adversarial networks to introduce more realistic high-frequency details and global illumination by learning the distribution from a set of high-quality Monte Carlo path tracing images. The latter is related to popular super-resolution techniques in the Computer Vision community, *e.g.* [23].

The above approaches are advanced denoising systems, which are limited regarding the details they can create that were not captured by the rendering engine or missing due to the low quality of the output. Missing details that were not captured in the current frame can only be reconstructed from information present in previously rendered frames, *e.g.*, [18].

### 2.2. Global illumination with machine learning

In the following, we highlight some past research that focuses on adding Indirect Illumination and global illumination effects

using neural networks by operating on the image plane level.

A CNN implementation in *Deep Shading* [24] predicts screen-space effects such as Ambient Occlusion, Motion Blur, Anti-Aliasing and Diffuse Indirect Illumination. The CNN takes OpenGL rasterization primary features, and produces outputs that can match real-time algorithms typically used in video games, while being able to generate any combination of effects in a single pass. This approach does not attempt to achieve photorealistic results.

A Conditional Generative Adversarial Network (CGAN) is trained in *Deep Illumination* [6] on specific scenes, with diffuse surfaces, to generate realistic and accurate global illumination effects. Despite being trained on specific scenes, the network is able to deal with moving objects and newly introduced shapes. Its temporal stability makes it suitable for producing animations, given a limited scene variability.

*Global illumination with radiance regression functions* [25] focuses on realtime indirect illumination rendering. This approach uses a neural network that learns the relationship between local and contextual attributes such as vertices and light position, to the indirect illumination value. This method shows very good performance and quality, although it is limited to point light sources and requires pre-baking of the Radiance Regression Function for each scene.

### 2.3. Machine learning for integration

A different integration point for machine learning is to accelerate the convergence of existing path tracers and bidirectional path tracers. Reinforcement Learning (RL) [26] exploits the similarity between the rendering equation and RL. The 3D scene is subdivided into Voronoi cells, each with an importance map dome that is updated using RL to remember the most efficient light path directions. The system is a path tracer that learns to dynamically update local importance sampling maps to increase its path efficiency and to reduce the number of zero-contribution paths. A similar idea is presented in *Machine Learning and Integral Equations* [27], where an approximate solution to the integral equation is learned.

*On-line Learning of Parametric Mixture Models for Light Transport Simulation* [28] optimizes the learning of difficult sampling distributions to guide standard ray tracing algorithms. This significantly improves the importance accuracy of light rays.

*Practical Path Guiding* [29, 30, 31] improves the efficiency of path tracing by adjusting importance sampling of direct and indirect bounces based on the radiance received on surfaces of the scene. This method accelerates initial convergence of ray tracing in difficult lighting situations and has been adopted in production engines.

*Deep Scattering* [32] shows the power of neural networks when applied to a very specific use case: rendering of clouds. This method is able to achieve extremely high quality in the generation of complex lighting effects such as cloud silverlining by using a hierarchical 3D representation of the cloud structure.

## 3. Contribution

We propose Deep Radiance Caching (DRC), an efficient variant of radiant caching, which separates the processing of direct and indirect lighting. DRC is able to produce convincing indirect illumination, including glossy reflections and ambient occlusion, without any preprocessing on the scene. Our main contributions are:

- Use of a Convolutional Autoencoder to predict high quality radiance maps from low quality data for indirect illumination;

- Reintegration of direct illumination using a separate integrator;

- Progressive sampling and smooth interpolation of indirect illumination values without pre-computation and additional storage.

## 4. Method

DRC renders the final image by computing direct and indirect illumination independently. The direct illumination pass uses a standard ray tracer with depth fixed at 1. The indirect illumination subdivides the image into tiles, and computes several radiance maps within each tile. Each radiance map is obtained by ray tracing a low resolution intensity, depth and normals map from the first intersection point into the scene. These maps are augmented by the Convolutional Autoencoder. The radiance maps are used to compute the final indirect illumination contribution of the pixels.

### 4.1. Illumination components

DRC splits rendering of direct illumination and indirect illumination. The indirect illumination component evaluates cached radiance maps, obtained efficiently using a deep neural network that reads local geometrical information of the scene. A standard path tracer evaluates the Rendering Equation [33] recursively:

$$
\begin{aligned}
L_o(x, \vec{\omega}_0) = & L_e(x, \vec{\omega}_0) + \\
& + \int_{\Omega_{2\pi}} f_r(x, \vec{\omega}_0, \vec{\omega}_i) cos(\theta_i) L_i(x, \vec{\omega}_i) d\omega_i.
\end{aligned}
\tag{1}
$$

$L_e$ is the radiance emitted by the surface, $f_r$ is the BRDF or BSDR of the material, the *cos* term is used to compute the irradiance accounting for the incident angle of the incoming radiance $L_i$ from the next bounce. A ray in path tracing bounces on surfaces in the scene until a light is hit, or the path is terminated. The recursively evaluated radiance, the BRDF $f$ of the material and the viewing direction $\vec{\omega}_i$ are used to obtain the final pixel value. DRC approximates the evaluation of the Rendering Equation by collapsing all the light bounces beyond the first into a single step. A ray tracer shoots a ray to find the first intersection point, and evaluates on its hemisphere a predicted

radiance map, a depth map, and a normal map. These three images are the inputs for our Convolutional Autoencoder, which returns a higher quality radiance map via

$$L_o(x, \vec{\omega}_0) = L_e(x, \vec{\omega}_0) +$$
$$\int_{\Omega_{2\pi}} f_r(x, \vec{\omega}_0, \vec{\omega}_i) cos(\theta_i) R_i(\vec{\omega}_i) d\omega_i. \tag{2}$$

We replace the recursive radiance term with $R$, the high quality radiance map obtained from the neural network, and evaluate it in the same way according to the BRDF $f$ of the surface in the viewing direction $\vec{\omega}_i$.

Previous research has achieved excellent degrees of success in denoising and reconstructing missing parts of the image from the output of a path tracer using Convolutional Neural Networks. Our network is particularly suited for extracting high level information from the locality of the image data. In the context of DRC, we chose to implement a Convolutional Autoencoder, as the noisy input can be associated to a highly lossy and compressed version of the same image when rendered at a higher number of samples per pixel. The output radiance map is placed back on the hemisphere around the first intersection point, and contains an approximation of all the recursive indirect radiance that path tracing would evaluate over time.

In practice, $R$ in Equation 2 does not have infinite resolution. DRC is able to preserve the high level of detail of the first-bounce illumination by excluding the direct light source contribution from $R$, and reintroducing it using a separate direct illumination pass. The final DRC Rendering Equation is

$$L_o(x, \vec{\omega}_0) = L_e(x, \vec{\omega}_0) +$$
$$\int_{\Omega_{2\pi}} f_r(x, \vec{\omega}_0, \vec{\omega}_i) cos(\theta_i) R_i(\vec{\omega}_i) d\omega_i +$$
$$\int_{\Omega_{2\pi}} f_r(x, \vec{\omega}_0, \vec{\omega}_i) cos(\theta_i) D_i(\vec{\omega}_i) d\omega_i. \tag{3}$$

$D_i$ is the direct illumination contribution from light sources visible at the first bounce. DRC can use a low resolution radiance map to approximate the $R$ term, while the direct illumination rendering pass preserves high frequency details such as edges and hard shadows.

Evaluating the entire radiance map $R$ would still be expensive. Our implementation uses Multiple Importance Sampling [34] with a small fixed number of samples to obtain indirect illumination values. We alternate sampling from $R$ as if it were an environment map with perfect visibility, and from the probability distribution of the BRDF. Ambient occlusion is encoded into these local environment maps because they are taken from the point of intersection.

To include direct illumination, the first intersection point also generates shadow rays to light sources (not shown in Figure 1). It is important to exclude direct light information from the radiance maps, as they do not have a high enough resolution to produce accurate shadows.

### 4.2. Progressive refinement and interpolation

Instead of evaluating a radiance map at each pixel of the image, DRC accelerates the rendering process by sampling a few pixels, interpolating their indirect illumination values, and progressively refining the final picture by adding samples. Sampling and interpolation only applies to indirect illumination, the direct illumination passes are performed over all pixels.

Progressive refinement allows the user to start rendering without setting a target number of passes, but to watch the image as its quality improves and stop the process once satisfied. Additional progressive refinement passes are also effective for smoothing small inconsistencies in indirect illumination.

The interpolation method of DRC does not use any precomputation to determine the best sampling locations. We compute a grid of sampling points on the image plane, evaluate a radiance map at each point, and interpolate all other values. After each complete pass on the image, the distance between points on the grid is reduced by a constant factor, allowing the algorithm to progressively refine the resolution of the indirect illumination component.

The interpolation strategy for points that are not sampled is based on both the distance from the sampled radiance maps and the surface normal at the primary intersection. The primary intersection is the closest non-specular intersection point of a ray shot from the camera. If a ray hits a transmission or mirror material, we follow the bounce recursively, and place the primary intersection point when the bounced ray hits again. The weight of each radiance map for pixel $i$ is based on the simple heuristic

$$w = w_p \cdot w_n + w_p + \epsilon. \tag{4}$$

Position and normal weights are multiplied to increase the weight when both are high, and position weight is added again to generate smoother weight values when the normal weight changes rapidly in noisy areas. $\epsilon$ is a small term used to avoid computation errors when all other terms are zero. The position weight is

$$w_p = dist(p_i, p)/r. \tag{5}$$

$dist(p_i, p)$ measures the distance in pixels on the image plane between the pixel being evaluated and the cached radiance map, and $r$ is the distance between two points in the grid.

The normal weight $w_n$ is maximal when the pixel's intersection point and the cached radiance map's normal ($n$ and $n_i$ respectively) point in the same direction, and becomes zero when their dot product is negative:

$$w_n = 1 - max\left(0, \left(\frac{n_i}{|n_i|} \cdot \frac{n}{|n|}\right)\right) \tag{6}$$

The weights $w$ of the radiance maps are normalized and converted into sampling probabilities for the radiance sampling.

### 4.3. Neural network

A neural network is at the core of DRC. The network predicts a smooth and accurate radiance map from primary features that can be computed quickly by a first hit ray tracer.

Each input or output layer has 32x32 floating point values. The input is composed of 7 layers: predicted radiance map
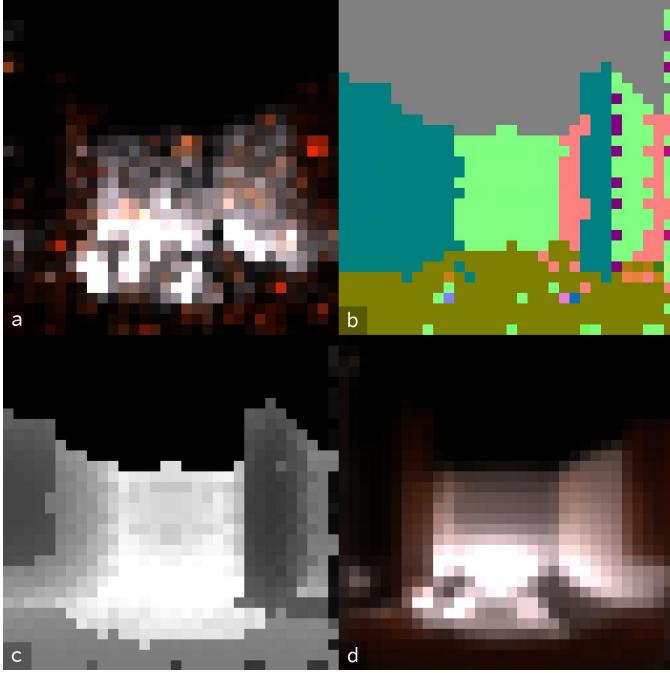
**Fig. 2. Examples of latitude-longitude hemispherical maps generated by DRC. This sample was generated from one intersection in the scene `bathroom`, part of the PBRTv3 example scenes [3]. Image (a) is a path traced intensity map at 1 sample/pixel. Image (b) is a normals map. Image (c) is the distance map. Image (d) is the same intensity map, traced at 4096 samples/pixel, used as reference ground truth during training.**

(RGB), normals (XYZ), and distance (Z). The output has 3 layers: Radiance (RGB). Figure 2 shows an example for input and output maps.

The maps are encoded in hemispherical latitude-longitude equirectangular projection. The coverage of a hemispherical map is not the full sphere, but only the upper hemisphere visible from a scene intersection point, with the surface normal aligned towards the Z direction of the latitude-longitude map. The up vector used in the hemispherical maps is the same as the one specified in the global scene, and it is rotated by 90 degrees towards the global Y vector only when the surface normal also points towards Z. In our final model we use the L1 loss $L1 = \sum_{i=0}^{n} |y_i - h(x_i)|$, where $y_i$ is a ground truth data point, and $h(x_i)$ is a predicted value.

Our network architecture is shown in Figure 3 and is similar to the U-Net architecture [35], a convolutional auto-encoder with skip connections and regular dimensionality progression. The encoder and decoder have symmetrical structure: each encoder stage uses two 3x3 convolutions and doubles the dimensional depth, while each decoder stage has two 3x3 deconvolutions and reduces the depth by half.

All intermediate stages use batch normalization and LeakyReLU activation functions. The output stage has two 3x3 deconvolutions with LeakyReLU, and a final 1x1 convolution with ReLU activation to output the final data. Each downsampling stage uses a 2x2 Max Pooling, and the upsampling stages use 2x2 Bilinear Upsampling. We do not use any Dropout layer.
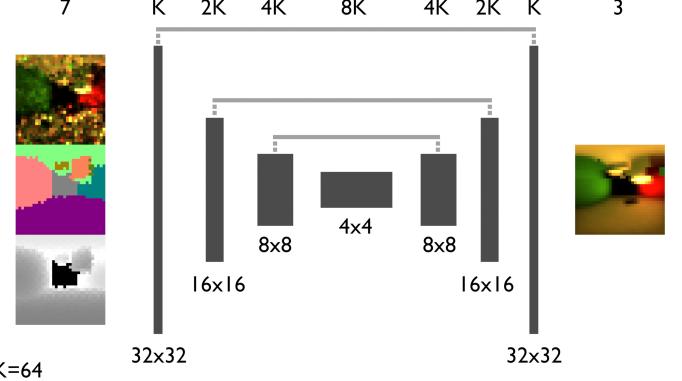


**Fig. 3. DRC Network layout**

## 5. Implementation

**Model training:** A set of 43 scenes has been selected from the PBRTv3 example scenes[1], Benedikt Bitterli's resources[36], and our own custom created scenes. A subset of the scenes is shown in Figure 4. As DRC focuses on rendering indirect illumination, the selection of the training scenes has been weighted to include mostly interior scenes with difficult illumination, and scenes with strong global illumination.



**Fig. 4. Selected scenes from the training set (cropped to fit)**

Each of the scenes has been processed to generate a set of hemispherical maps. Similarly to the approach taken by Kalantari [14], we change sampling algorithms and seeds to prevent the network from overfitting to specific noise patterns rather than higher level features. We use multiple samplers, including Sobol and Random, to produce varied noise patterns. We collected 16000 individual training examples. This dataset was sufficient in our tests to achieve good training results.

We implemented the neural network model and training process in Python using Pytorch[2]. Network training has been performed with GPU acceleration on an Nvidia GTX 1060 6GB card hosted by an Intel i7 4770 processor. All images and performance metrics have been acquired by using a single CPU, a Ryzen 2700. Network inference is done in DRC without any GPU acceleration.

We chose a base layer size of K=64, and a minimum hidden layer dimension of 4x4. These parameters result in 3 downsizing steps and 3 corresponding upsampling steps.

---

[1] http://pbrt.org/scenes-v3.html
[2] http://pytorch.org/

We used the Adam optimizer included in PyTorch, with a learning rate of $6e-5$ and implemented data augmentation in a Pytorch data loader to include combinations of rotation and flipping. We achieved good network performance after 20 minutes of training with mini-batches containing 32 examples each. Our training covers 3 epochs over the augmented dataset: more than half a million individual examples.

Our full dataset is 0.5GB gzipped, and can be easily loaded into the memory of a single GPU when uncompressed during network training. The trained model is 20MB, which can be efficiently evaluated on modern CPUs, whose caching mechanisms often exceed this size.

**PBRT Extension:** The DRC method and algorithm has been implemented by modifying PBRTv3 [37], a C++ ray tracer implementing many popular algorithms.

Starting from PBRTv3, we implement a new `integrator` which evaluates direct and indirect illumination separately and communicates to the neural network to obtain radiance map predictions. PBRT-DRC can operate in two different modes: *Reference* and *Rendering*.

While the *Rendering* mode can be used by the end-user to render a scene, the *Reference* mode computes training examples for given scenes. The *Reference* mode automatically generates a large number of examples suitable for both training and validation purposes. The number of tiles can be defined, which determines the way the rendered image is subdivided into a regular grid. Each intersection in the grid is used to shoot a single camera ray into the scene, and the intersection found becomes the viewpoint of one set of examples. The high quality path tracing render's number of samples can be specified by the user.

In the rendering process we use two independent steps to compute indirect illumination and direct illumination. The direct illumination pass is handled using a standard direct illumination integrator implemented in the original PBRT software. The indirect illumination uses the *Path* integrator to collect the neural network's input maps. The machine learning evaluation runs in a separate Python process running Pytorch. Each rendering thread has its own child process and communication is handled through standard Unix pipes. Each Python evaluation process loads the saved model parameters, and reads flattened input data from `stdin`. The predicted output is written back to `stdout`.

## 6. Evaluation and Results

DRC is evaluated in terms of its network performance, final image quality, and overall rendering performance. We analyze the output of the Convolutional Autoencoder in Section 6.1, perform an ablation study in Section 6.2 and evaluate image quality using established metrics like SSIM and Entropy in Section 6.3. SSIM is preferred in our experiments over, *e.g.*, PSNR, since we are analyzing Monte Carlo rendered images [38].

### 6.1. Neural network evaluation

Figure 5 shows training and test loss values during training. Our model reaches stable loss values over a 20 minutes learning period, without overfitting to the training dataset.
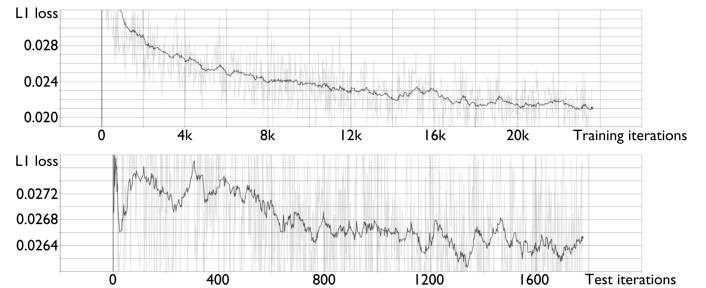


**Fig. 5. Training and Test L1 Loss in Model 11**

The network was trained on a collection of 3D scenes, and the evaluation is conducted on three selected scenes that were never part of the training dataset: *White Room Daytime* is part of the original PBRTv3 example scenes [3], *Veach Ajar* by Veach and adapted by B. Bitterli [36], and *Mbed1*, custom created for the purpose of this evaluation.

The difference in amount and depth of indirect illumination paths across the test scenes allows us to verify that our network is able to deal with different input data quality seamlessly, and that it is able to extrapolate and adapt to geometry not encountered during training.

As with the training set generation, we obtain input maps and ground truth path traced radiance maps for several viewpoints located at primary surface intersection points. While the PBRT-DRC implementation re-normalizes the output images to match the expected intensity level of the entire radiance map, in our tests we compare directly using the final intensities obtained after the upstream processing to provide an unbiased view of the network performance.

Figure 6 shows a random selection of radiance maps predicted using our neural network model. The first three columns show the inputs of the network: normals are tonemapped from XYZ to RGB via $n_i = ((n_i * 0.5) + 0.5) * 255$, distances from the camera and 1 sample/pixel path tracing radiance maps. The result is compared to a simple Gaussian blur and a reference path traced ground truth, rendered at 1024 samples per pixel. As many of the initial path tracer outputs comprise mostly of black pixels and a few very bright ones, standard image denoising filters such as Bilateral and Median Filter do not perform well. The Gaussian blur, although yielding softer outputs, is more capable at filling the black gaps.

The Gaussian blur uses a one standard deviation filter width, chosen for its good compromise between smoothing power and blurriness. A smaller radius would not be able to remove much noise, while a larger radius would yield considerably blurrier results compared to either ground truth or network predicted output. Due to the very sparse samples generated by the path tracer, a simple Gaussian blur is a significant improvement over the 1spp (1 sample per pixel) intensity map.

The examples show that the network predicted output can consistently outperform the Gaussian blur by producing sharper images, preserving more detail, and removing more noise. We can observe that:

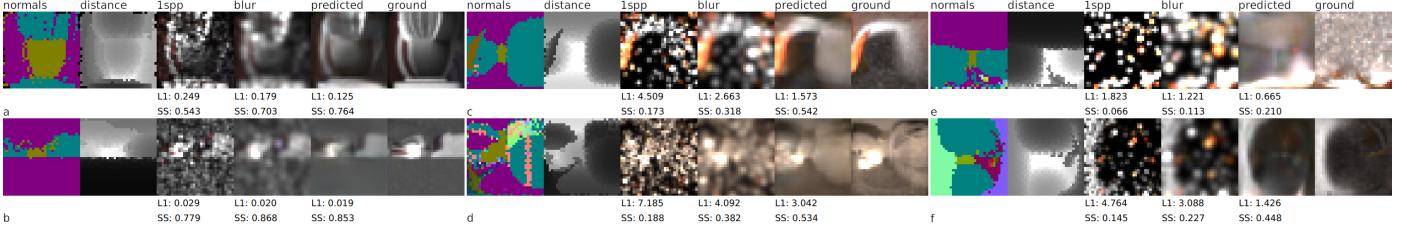- The predicted result is generally sharper than the Gaussian

**Fig. 6. A random selection from the test scenes *Veach Ajar*, *Custom Mbed*, *White Room Daytime*.** The normals are mapped via $n_i = ((n_i * 0.5) + 0.5)$ and then scaled to RGB space. All examples show excellent noise removal when comparing *1spp* to *predicted*, with the predictions occasionally presenting less noise than the reference ground truth. In e) there is some color inaccuracy, with the predicted image showing some slight tint that is not present in the ground truth. e) and f), part of *Veach Ajar*, are very difficult due to the very sparse sampling of the input intensity map. They show excellent noise removal, producing less noise than the ground truth. c) shows a slight shape distortion on the orange part, and a), b) and d) are less sharp than the ground truth, but still considerably more usable than the Gaussian blur versions.

blur. This confirms that we would not be able to use a larger radius for the Gaussian blur for difficult cases without losing more details.

- The network is able to reconstruct some details that are missing from the 1spp intensity map, although it can cause artifacts.

- The network has effectively learned to use adaptive blurring filters. In simple cases it behaves very similarly to a planar Gaussian blur, but in the general case it is able to adapt to different sampling densities more effectively.

- The network eliminates all visible noise, and can produce images that look smoother with less noise when compared to the high quality path traced reference images. This shows the ability of the network to extrapolate from the training dataset, which still contains a small amount of noise in its ground truth examples. This is in line with findings from Noise2Noise [39].

To verify our claims, we evaluate similarity metrics on a large number of radiance evaluation points collected from the testing scenes, totaling over 1000 distinct testing examples that were not seen during training. For each example set of images, we compute the L1 absolute difference between predicted and reference images

$$L1 = \frac{1}{n} \sum_{i=0}^{n} |p_i - g_i|, \qquad (7)$$

where $p_i$ is a predicted pixel value, $g_i$ is a reference pixel value. Figure 7 shows L1 and Structural Similarity Index [40]. The Gaussian blur does considerably increase the quality of the radiance maps compared to using raw 1spp images. The structural similarity increases, as the blur fills many of the black gaps present in the path traced map.

The network output considerably outperforms the Gaussian blur in both L1 and structural similarity metrics. We confirm that for both metrics the samples belong to different probability distributions by computing the Kruskal-Wallis null hypothesis testing [41] p values for L1 and Structural Similarity. We verified the p values to be less than 0.00001.

### 6.2. Neural network ablation study

The ablation study aims to verify hypotheses about the data that the neural network uses to produce accurate predictions.

Training is conducted using the entire dataset of intensity, normals, and distance maps, and repeated with the normals, distance, or both components removed from the network's inputs. The structure of the network is left unchanged, and the ablation is implemented by setting the target components to be zero arrays. The evaluation is conducted in a similar way, by setting the ablated channels to zero arrays.

Figure 8 shows loss values during training and testing using the four different configurations. The use of only first hit path tracing maps performs better than a simple blur, but is improved by the use of normal maps. Using distances without normals does not yield better performance, but the combination of all three inputs outperforms the other configurations.

The Structural Similarity Index test on the entire test set in Figure 9 confirms that the addition of normal or distance maps helps the network to predict more accurate results, with the normals being slightly more effective than distances alone. The combination of both additional layers further increases accuracy, and in particular improves the lower quartile notably.

### 6.3. Final rendering evaluation

The evaluation of the final images is based on both qualitative and quantitative metrics. We render each of our test scenes at different levels of quality exclusively on a CPU (Ryzen 2700 without GPU acceleration), and present observations, statistics and comparisons.

DRC is not based on a per-pixel quality setting, but on two distinct settings: the number of samples per pixels for the direct illumination integrator, and the number of tasks to be processed in the indirect illumination pass. A *task* is a single computational job in the indirect illumination pass that corresponds to a rectangular region of the rendered image. A task is responsible for collecting the input maps in the region, evaluating them through the neural network, and adding the final light contributions to the image. DRC heuristically splits indirect illumination rendering into tasks to distribute the computational load onto multiple CPU cores.

We chose a target quality setting by rendering the scene at a fixed number off direct lighting passes, while changing the number of indirect lighting tasks. For this comparison, we use both DRC and the bidirectional path tracer integrator from PBRTv3 in CPU-only mode, without using any dedicated accelerator to evaluate the neural network. We used a PC with a Ryzen 2700 and 16GB of RAM. We measure the output's Structural Similarity compared to a high quality path traced image rendered at 2048 samples/pixel, and the PNG file size
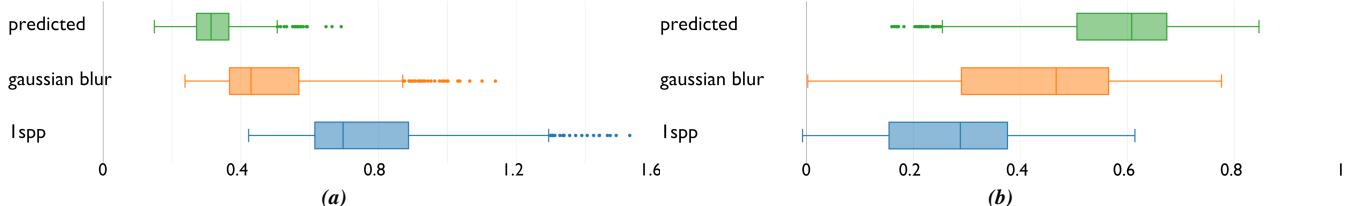
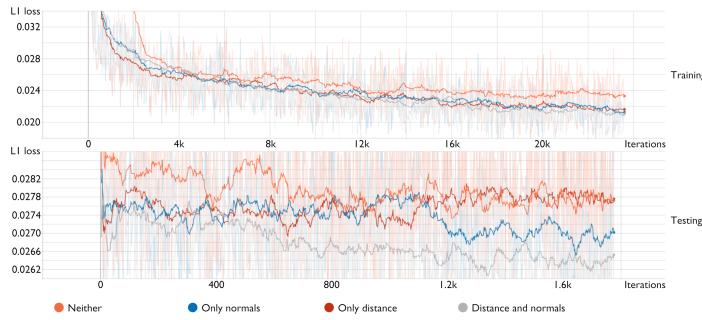**Fig. 7. (a) L1 absolute difference, (b) Structural Similarity**



**Fig. 8. Training loss ablation testing. Removing information from the input data slows down convergence of the loss values. The effect is visible on the training set, but becomes more evident on the test set.**

as an indication of the amount of noise. Table 1 shows that DRC's overhead on render time is very small until 16 indirect tasks, with Structural Similarity values increasing. By using more than 16 indirect tasks, quality improves only marginally, although rendering takes a much longer time. As DRC is a biased method, Structural Similarity does not converge towards 1, and any artifacts present in the output need to be judged subjectively. The amount of noise continues decreasing as the indirect tasks increase, although the noise in the output image is mostly produced by the direct illumination passes. We therefore use 16 indirect samples for our comparisons, and 8 direct lighting samples as it is a good compromise between noise and speed.

We compare the DRC result at 8 direct and 16 indirect passes, which took 43 seconds, with path tracing at 16 samples/pixel
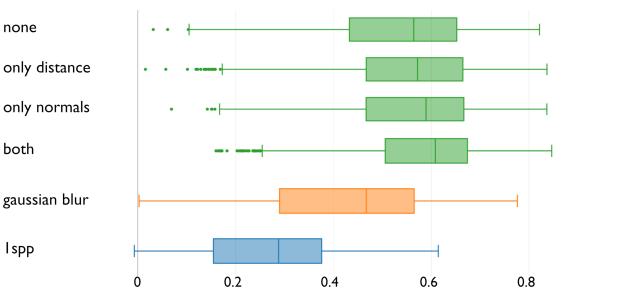


**Fig. 9. The box plots show structural similarity values in our ablation testing. None: predicted results with the network trained on intensity maps only. Only distance: network trained with intensity and distance maps. Only normals: network trained with intensity and normal maps. Both: network trained with intensity, distance and normal maps. Gaussian blur: Gaussian blur applied to the input intensity map. 1spp: input intensity map path traced at 1 sample per pixel.**

**Table 1. Quality and noise metrics on the Bathroom Green scene [42] with an increasing number of indirect illumination tasks. We compare DRC to BDPT [1] integration numerically. We evaluate *Render time* as total rendering time in seconds; *SSIM*: the Structural Similarity Index; and the *PNG File Size* as a measure of the amount of noise in the final image. Better or equal values are highlighted in bold.**

| | BDPT | | | DRC | | |
|---|---|---|---|---|---|---|
| Tasks | Time (s) | SSIM | PNG (KB) | Time (s) | SSIM | PNG (KB) |
| 1 | **3** | 0.48 | 2497 | 36 | **0.77** | **2482** |
| 2 | **7** | 0.57 | 3094 | 36 | **0.84** | **2691** |
| 4 | **12** | 0.66 | 3176 | 38 | **0.87** | **2549** |
| 8 | **23** | 0.75 | 3050 | 39 | **0.88** | **2486** |
| 16 | 44 | 0.82 | 2885 | **43** | **0.89** | **2412** |
| 32 | 95 | 0.87 | 2710 | **59** | **0.89** | **2354** |
| 64 | 181 | **0.90** | 2549 | **88** | **0.90** | **2313** |
| 128 | 351 | **0.92** | 2392 | **160** | 0.91 | **2278** |
| 256 | 693 | **0.93** | 2258 | **322** | 0.91 | **2250** |

completed in a similar amount of time, 46 seconds, and bidirectional path tracing completed in 44 seconds. A reference path traced image at 2048 samples/pixel is also included (99m 30s).

In the comparison, we also include results post processed by the Intel® Open Image Denoiser [43], a production-quality image denoiser built to be used with path traced images. We used the Open Image Denoiser with all settings at default, in LDR mode.

Figure 10 shows a comparison between path tracing, DRC, bidirectional path tracing and each of them after applying the Open Image Denoiser. The DRC output appears less noisy than an image from path tracing after the same amount of rendering time. The generated global illumination effects are visually improved in the DRC output, including the green bleeding onto the slightly-glossy sink, and the ambient occlusion in the darker areas. The crops show that DRC is capable of a higher level of detail, especially in poorly illuminated areas. The Open Image Denoiser framework works well for removing all noise, although the filtered path tracing output at 16 samples/pixel has a notable loss of detail in the area below the sink, while the filtered DRC output retains more details.

It is important to point out that the noise that is still visible in the DRC output is not caused by our method, but by the re-integration of the direct illumination layer.

**Glossy and specular materials:** One of the key properties of

with Open Image Denoise     without Open Image Denoise

DRC 8 direct passes
16 indirect passes

unidirectional
path tracing
16 samples/pixel

unidirectional
path tracing
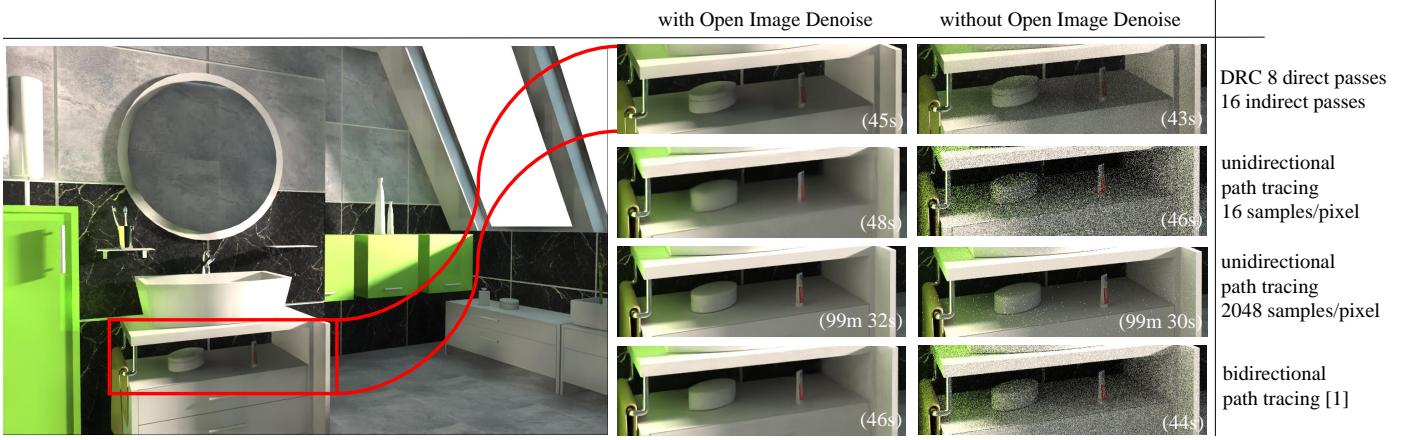2048 samples/pixel

bidirectional
path tracing [1]

**Fig. 10.** Comparison between DRC, path tracing, and Bi-Directional path tracing on Bathroom Green [42]. path tracing at 2048 samples/pixel should be considered the reference image. We have adjusted render quality settings for other methods to yield similar computation times. Note that the camera angle and some materials and textures are not identical to [42] since we had to manually convert the scene to PBRT.
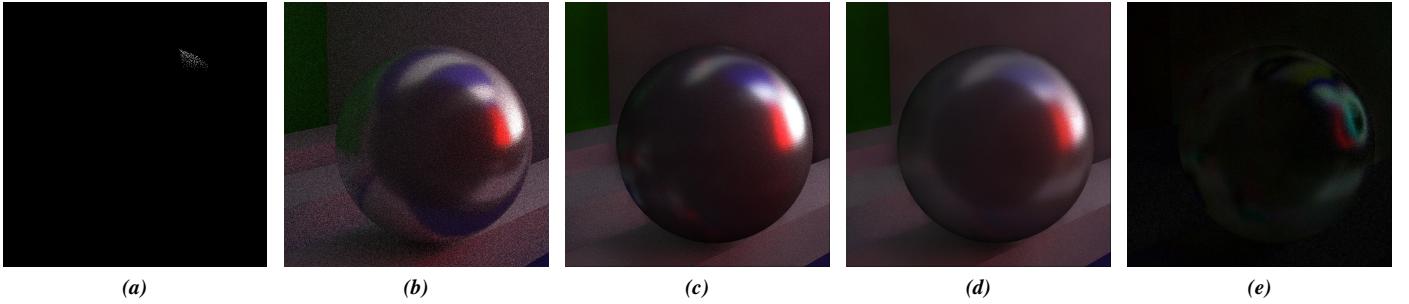


*(a)*           *(b)*           *(c)*           *(d)*           *(e)*

**Fig. 11.** A direct comparison between **(a)** Direct illumination, **(b)** Path tracing, **(c)** DRC Gaussian blur, **(d)** DRC Neural Network, and **(e)** the difference between between (c) and (d). The baseline method (c) misses several details on the lower left part of the sphere that have been correctly predicted in (d) when compared to (b). MSE: (c), (b) = 500.31; (d), (b) = 384.81. SSIM: (c), (b) = 0.21; (d), (b) = 0.23.

DRC is its ability to work with a wide range of materials, including glossy surfaces. Global illumination reflections heavily depend on the quality of the radiance maps. The neural network evaluation (Section 6.1) shows that a Gaussian blur applied to the intensity maps offers a good level of similarity compared to the ground truth. Figure 11 shows a sphere with a highly glossy surface under difficult lighting conditions. The direct illumination integrator in Figure 11(a) is almost completely black, confirming that all reflections are indirect. Blurred radiance maps may be sufficient to handle materials with a low level of glossiness or on low importance scene objects. The neural network in DRC brings improvements in the rendering quality of shinier surfaces. For example, blurred radiance maps do not have sufficient resolution in the dark areas (compare Figure 11(c) and (d)) and the neural network produces more stable and detailed output so that reflected shapes are better preserved. In contrast to denoising and blur filters, DRC can further be tuned by increasing the size of the predicted radiance maps, which directly reduces resolution-induced blur in the final result. The computational overhead of (d) compared to (c) is 40%, which can be improved in future extensions by adaptively choosing a radiance map enhancement strategy depending on the material properties.

Specular highlights, reflections and transmission materials are handled by following the original ray, and evaluating the radiance maps after the bounce. The threshold between the first and secondary ray bounce in DRC has therefore slightly different semantics from most ray tracers as perfect mirrors and transmission rays do not increase the depth counter. These specular rays are also ignored in the depth count of the direct illumination pass of DRC to provide a correct complementary integrator. In the details of the Bathroom Green scene in Figure 10, the metallic objects preserve detailed reflections and realistic structural details are predicted.

**Interpolation:** Figure 12 shows two images rendered at a limited number of tasks. The simple interpolation strategy linearly interpolates using pixel distances on the image plane. Light bleeding is clearly visible on the objects. This image also shows hard edges, artifacts resulting from individual threads rendering different tiles of the image, and the lack of interpolation among adjacent tiles. We resolve some major artifacts by including surface normals in the weighting system (Section 4.2), and by improving the way tiles are overlapped in the image. Some artifacts remain near the edges due to the very low number of samples, while the overall image looks smooth.

Increasing the number of indirect illumination passes in DRC
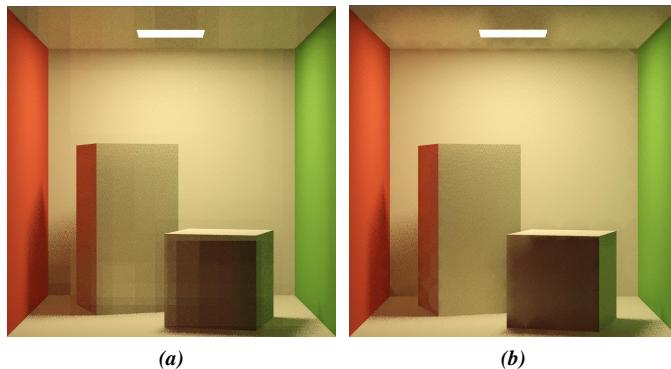
*(a)*                                   *(b)*

**Fig. 12. Comparison between simple distance-based interpolation (a) and interpolation based on both distance and surface normals (b) Using surface normals improves smoothness across surfaces, and removes large artifacts and light bleeding on the edges of the cubes.**

only affects computational time. Memory usage remains constant, as each interpolation pass is independent and does not require any data from previous passes.

## 7. Discussion

The DRC method offers a new biased rendering method for global illumination capable of producing high quality results efficiently, while supporting a wide range of material types and scene compositions.

Compared to other methods, DRC does not require any form of pre-computation on the target scene, and accelerates the computation of global illumination at a virtually unbounded number of indirect bounces at the cost of a small amount of bias introduced. The typical use case scenario is rendering an interior 3D scene where indirect light is significant. Most interior scenes have strong indirect lighting, making light paths computationally intensive.

Faster estimation of indirect light through a neural network can drastically reduce the amount of noise present in the scene, and the typical low-frequency variation of indirect illumination easily masks away minor artifacts and interpolation imprecisions that originate from the machine learning subsystem. Thanks to the preservation of a standard direct illumination pass in the final image, most of the details and textures present in the scene are preserved in the final output.

DRC splits the processing of direct and indirect illumination, allowing the user to better tune the amount of computation to be used for each component.

### 7.1. Limitations

Even though the progressive refinement algorithm permits the user to set an arbitrary number of passes, DRC remains a biased approach to global illumination. The neural network, although very accurate in typical scenarios, can produce artifacts, inaccurate results, and loss of detail, causing hard to correct errors.

A second important limiting factor is the fixed resolution of the radiance maps, set for DRC at a default of 32x32 pixels. Small details are inevitably lost.

These two limitations become particularly visible when computing globally illuminated glossy reflections, which appear less defined, although less noisy, than the path traced reference.

Temporal stability is another limitation of DRC. The rendering tiles can have large variations from one sample to another, and while being perceived as smooth transitions in a single frame, animations would display them as low frequency flickering of large areas in the scene. The lack of temporal stability indicates that DRC would not be adequate for animation rendering. To solve the temporal stability issue, a persistent caching system for radiance maps can be implemented, allowing animations to be rendered at faster rates and with less temporal noise caused by flickering of radiance maps. Precomputing and caching radiance maps would however not be sufficient to achieve temporal stability in scenes with dynamic objects.

Due to the low resolution of the indirect radiance maps, DRC is not capable of producing accurate caustics, which require a much denser sampling of indirect light.

The progressive refinement algorithm used in DRC requires no preprocessing and very little additional memory to run. However in some scenarios predetermined and optimized radiance interpolation points can be a more effective solution, and implementing the possibility to choose between the two strategies would be a useful future extension.

## 8. Conclusion

We have presented *Deep Radiance Caching* – DRC, a novel ray tracing method that uses Convolutional Autoencoders to obtain high performance and accurate global illumination effects. We show that our method achieves competitive performance on a CPU, while being able to produce higher quality images than same-time path tracing, with a significantly smaller amount of noise. Our method supports a wide range of material types and does not need offline pre-computation or per-scene training. If required, both, the used deep learning network and the path tracing core map naturally to parallel hardware like GPUs, which has been shown to achieve real-time performance for ray-tracing [44, 45] and denoising [46]. Recent developments like Nvidia's RTX technology [47] opens up further avenues to exploit deep neural networks deeper in image generation algorithms.

**Resources:** A website including plugin for Blender and a source code repository on github are available and will be added here after the anonymous reviewing phase.

## References

[1] Lafortune, EP, Willems, Y. Bi-directional path tracing. In: Compugraphics' 93. 1993, p. 145–153.

[2] Veach, E, Guibas, LJ. Metropolis light transport. In: Proceedings of the 24th annual conference on Computer graphics and interactive techniques. ACM Press/Addison-Wesley Publishing Co.; 1997, p. 65–76.

[3] Pharr, M, Jakob, W, Humphreys, G. Scenes for pbrt-v3. http://pbrt.org/scenes-v3.html; 2018. Accessed: 25/04/2018.

[4] Keller, A. Instant radiosity. In: Proceedings of the 24th Annual Conference on Computer Graphics and Interactive Techniques. SIGGRAPH '97; New York, NY, USA: ACM Press/Addison-Wesley Publishing Co. ISBN 0-89791-896-7; 1997, p. 49–56. URL: https://doi.org/10.1145/258734.258769. doi:10.1145/258734.258769.

[5] Krivanek, J, Gautron, P, Pattanaik, S, Bouatouch, K. Radiance caching for efficient global illumination computation. IEEE Transactions on Visualization and Computer Graphics 2005;11(5):550–561.

[6] Thomas, MM, Forbes, AG. Deep illumination: Approximating dynamic global illumination with generative adversarial network. In: arXiv preprint arXiv:1710.09834. 2017,.

[7] Chambolle, A, Pock, T. A first-order primal-dual algorithm for convex problems with applications to imaging. Journal of mathematical imaging and vision 2011;40(1):120–145.

[8] Beck, A, Teboulle, M. Fast gradient-based algorithms for constrained total variation image denoising and deblurring problems. IEEE transactions on image processing 2009;18(11):2419–2434.

[9] Buades, A, Coll, B, Morel, JM. A non-local algorithm for image denoising. In: 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05); vol. 2. IEEE; 2005, p. 60–65.

[10] Elad, M, Aharon, M. Image denoising via sparse and redundant representations over learned dictionaries. IEEE Transactions on Image processing 2006;15(12):3736–3745.

[11] Liu, P, Zhang, H, Zhang, K, Lin, L, Zuo, W. Multi-level wavelet-cnn for image restoration. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops. 2018, p. 773–782.

[12] Bitterli, B, Rousselle, F, Moon, B, Iglesias-Guitián, JA, Adler, D, Mitchell, K, et al. Nonlinearly weighted first-order regression for denoising monte carlo renderings. In: Computer Graphics Forum; vol. 35. Wiley Online Library; 2016, p. 107–117.

[13] Moon, B, McDonagh, S, Mitchell, K, Gross, M. Adaptive polynomial rendering. ACM Transactions on Graphics (TOG) 2016;35(4):1–10.

[14] Kalantari, NK, Bako, S, Sen, P. A machine learning approach for filtering monte carlo noise. ACM Trans Graph 2015;34(4):122:1–122:12. URL: http://doi.acm.org/10.1145/2766977. doi:10.1145/2766977.

[15] Pharr, M, Humphreys, G. Physically Based Rendering: From Theory To Implementation. Morgan Kaufmann, 2nd edition; 2010.

[16] Bako, S, Vogels, T, McWilliams, B, Meyer, M, Novák, J, Harvill, A, et al. Kernel-predicting convolutional networks for denoising monte carlo renderings. In: Proceedings of ACM SIGGRAPH 2017. ACM; 2017,.

[17] Vogels, T, Rousselle, F, McWilliams, B, Röthlin, G, Harvill, A, Adler, D, et al. Denoising with kernel prediction and asymmetric loss functions. ACM Transactions on Graphics (Proceedings of SIGGRAPH 2018) 2018;37(4):124:1–124:15. doi:10.1145/3197517.3201388.

[18] Chaitanya, CRA, Kaplanyan, AS, Schied, C, Salvi, M, Lefohn, A, Nowrouzezahrai, D, et al. Interactive reconstruction of monte carlo image sequences using a recurrent denoising autoencoder. ACM Trans Graph 2017;36(4):98:1–98:12. URL: http://doi.acm.org/10.1145/3072959.3073601. doi:10.1145/3072959.3073601.

[19] Kuznetsov, A, Kalantari, NK, Ramamoorthi, R. Deep adaptive sampling for low sample count rendering. In: Computer Graphics Forum; vol. 37. Wiley Online Library; 2018, p. 35–44.

[20] Vicini, D, Adler, D, Novák, J, Rousselle, F, Burley, B. Denoising deep monte carlo renderings. In: Computer Graphics Forum; vol. 38. Wiley Online Library; 2019, p. 316–327.

[21] Gharbi, M, Li, TM, Aittala, M, Lehtinen, J, Durand, F. Sample-based monte carlo denoising using a kernel-splatting network. ACM Transactions on Graphics (TOG) 2019;38(4):1–12.

[22] Xu, B, Zhang, J, Wang, R, Xu, K, Yang, YL, Li, C, et al. Adversarial monte carlo denoising with conditioned auxiliary feature modulation. ACM Transactions on Graphics (TOG) 2019;38(6):1–12.

[23] Ledig, C, Theis, L, Huszár, F, Caballero, J, Cunningham, A, Acosta, A, et al. Photo-realistic single image super-resolution using a generative adversarial network. In: Proceedings of the IEEE conference on computer vision and pattern recognition. 2017, p. 4681–4690.

[24] Nalbach, O, Arabadzhiyska, E, Mehta, D, Seidel, HP, Ritschel, T. Deep shading: Convolutional neural networks for screen-space shading. In: Computer Graphics Forum. 2017,.

[25] Ren, P, Wang, J, Gong, M, Lin, S, Tong, X, Guo, B. Global illumination with radiance regression functions. ACM Transactions on Graphics (TOG) 2013;32(4):130.

[26] Dahm, K, Keller, A. Learning light transport the reinforced way. In: International Conference on Monte Carlo and Quasi-Monte Carlo Methods in Scientific Computing. Springer Proceedings in Mathematics and Statistics Vol 241; 2016,.

[27] Dahm, K, Keller, A. Machine learning and integral equations. CoRR 2017;abs/1712.06115. URL: http://arxiv.org/abs/1712.06115. arXiv:1712.06115.

[28] Vorba, J, Karlík, O, Šik, M, Ritschel, T, Křivánek, J. On-line learning of parametric mixture models for light transport simulation. ACM Transactions on Graphics (Proceedings of SIGGRAPH 2014) 2014;33(4).

[29] Herholz, S, Elek, O, Vorba, J, Lensch, H, Křivánek, J. Product importance sampling for light transport path guiding. In: Computer Graphics Forum; vol. 35. Wiley Online Library; 2016, p. 67–77.

[30] Müller, T, Gross, M, Novák, J. Practical path guiding for efficient light-transport simulation. In: Computer Graphics Forum; vol. 36. Wiley Online Library; 2017, p. 91–100.

[31] Müller, T. "practical path guiding" in production. In: ACM SIGGRAPH Courses: Path Guiding in Production, Chapter 10. New York, NY, USA: ACM; 2019, p. 18:35–18:48. doi:10.1145/3305366.3328091.

[32] Kallweit, S, Müller, T, McWilliams, B, Gross, M, Novák, J. Deep scattering: Rendering atmospheric clouds with radiance-predicting neural networks. ACM Transactions on Graphics (TOG) 2017;36(6):231.

[33] Kajiya, JT. Seminal graphics. chap. The Rendering Equation. New York, NY, USA: ACM. ISBN 1-58113-052-X; 1998, p. 157–164. URL: http://doi.acm.org/10.1145/280811.280987. doi:10.1145/280811.280987.

[34] Veach, E, Guibas, L. Bidirectional estimators for light transport. In: Photorealistic Rendering Techniques. Springer; 1995, p. 145–167.

[35] Ronneberger, O, Fischer, P, Brox, T. U-Net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical image computing and computer-assisted intervention. Springer; 2015, p. 234–241.

[36] Bitterli, B. Rendering resources. https://benedikt-bitterli.me/resources/; 2018. Accessed: 24/05/2018.

[37] Pharr, M, Jakob, W, Humphreys, G. Physically Based Rendering: From Theory To Implementation. Morgan Kaufmann, 3rd edition; 2016.

[38] Whittle, J, Jones, MW, Mantiuk, R. Analysis of reported error in monte carlo rendered images. The Visual Computer 2017;33(6-8):705–713.

[39] Lehtinen, J, Munkberg, J, Hasselgren, J, Laine, S, Karras, T, Aittala, M, et al. Noise2noise: Learning image restoration without clean data. arXiv180304189 2018;.

[40] Wang, Z, Bovik, AC, Sheikh, HR, Simoncelli, EP. Image quality assessment: from error visibility to structural similarity. IEEE Transactions on Image Processing 2004;13(4):600–612.

[41] Kruskal, WH, Wallis, WA. Use of ranks in one-criterion variance analysis. Journal of the American statistical Association 1952;47(260):583–621.

[42] Cenobi, U. Bathroom scene. https://www.blendswap.com/blends/view/52486; 2012. Accessed: 29/08/2018.

[43] Intel®, . Open image denoiser. Tech. Rep.; Intel (R); 2020. URL: https://openimagedenoise.github.io/.

[44] Parker, SG, Bigler, J, Dietrich, A, Friedrich, H, Hoberock, J, Luebke, D, et al. Optix: a general purpose ray tracing engine. ACM Transactions on Graphics (TOG) 2010;29(4):66.

[45] Goodfellow, I, Bengio, Y, Courville, A, Bengio, Y. Deep learning; vol. 1. MIT press Cambridge; 2016.

[46] Mara, M, McGuire, M, Bitterli, B, Jarosz, W. An efficient denoising algorithm for global illumination. High Performance Graphics 2017;10:3105762–3105774.

[47] Burgess, J. RTX on—The NVIDIA Turing GPU. IEEE Micro 2020;40(2):36–44.