

CS 496 BTP II

Report

Coherent Rendering for Augmented Reality

*Submitted in partial fulfillment of
the requirements for the award of the degree of*

**Bachelor of Technology
in
Computer Science and Engineering**

Submitted by

Sneha Bhakare
150050040

Under the guidance of
Prof. Parag Chaudhuri



Department of Computer Science and Engineering

INDIAN INSTITUTE OF TECHNOLOGY BOMBAY

Mumbai, Maharashtra, India – 400076

Spring Semester 2018-19

Acknowledgments

I would like to express my gratitude to Prof. Parag Chaudhuri for his valued guidance and continuous support during the entire course of work.

Sneha Bhakare
IIT Bombay
5 May, 2019

Abstract

Coherent Rendering for Augmented Reality is concerned with seamlessly blending the virtual world and real world in real time. We are interested in applying real-world light to virtual objects. This implies the measurement of the real-world lighting, also known as photometric registration. In this project we implement a method to estimate high quality illumination using convolutional neural networks (CNNs). Our representation of light model of real world is based on spherical harmonics representation. The aim of this project is to make augmentations photorealistic while removing any constraints in the environment.

Contents

Acknowledgements	i
Abstract	ii
1 Introduction	1
1.1 Augmented Reality	1
1.2 Coherent Rendering	1
1.3 Spherical Harmonics	2
1.4 Problem Definition	2
2 Related Work	3
2.1 Active lightprobes	3
2.2 Inverse Rendering Equation	3
2.3 Deep outdoor illumination estimation	3
2.4 Learning Lightprobes	4
3 Method	5
3.1 Synthesis of training images	5
3.2 Training the CNN	5
3.3 Implementation Details	6
4 Experiments and Results	7
4.1 Estimation for a Fixed Pose of known object	7
4.2 Estimation independent of object pose	8
4.3 Performance for different network parameters	9
4.4 Failed attempt: Approach with planes	10
5 Conclusion and Future Work	11
References	12

Chapter 1

Introduction

1.1 Augmented Reality

Augmented reality is an attractive and exciting way to present virtual content in a real context for various application domains, like architectural visualizations, virtual prototyping, marketing and sales of not yet existing products. With such wide range of applications, AR has developed into a topic of broad and current interest in computer graphics. With the advent of GPUs and other improvements in computation speeds, Augmented Reality has become feasible in realtime.

1.2 Coherent Rendering

Many of the applications require AR to be convincing enough for a spectator to not be able to differentiate between real and virtual. This requirement demands both the real and virtual worlds to be lit with same light models. This is what is also known as visual coherence. A key step for coherent rendering is to estimate the incident lighting of a scene. Common approaches use active lightprobes to capture the light model. The major disadvantage of having a lightprobe is that it interferes with the user's experience.

Gruber et al. [3] demonstrated a probe-less method that estimates the incident light purely from reflections in the unmodified scene, scanned in real-time with a commodity RGB-D sensor. This method solves inverse rendering equation to get spherical harmonics coefficients and thus render the scene using the obtained light model. This approach forgoes artificial lightprobes, but suffers from high computational demands and cannot be used on mobile devices.

Another viable learned lightprobes method is proposed by Mandl et al. [7] The learned lightprobes work in an unmodified scene, providing unobtrusive user experience, while still retaining the computational benefits of a normal lightprobe. The incident lighting is estimated with a pre-trained convolutional neural network (CNN), which is orders of magnitude cheaper to evaluate at runtime than inverse rendering.

1.3 Spherical Harmonics

The Spherical Harmonic[10] (SH) series of functions is the analogue of the Fourier Series for functions on the surface of a unit sphere centered at origin. The complete set of functions is an infinite-dimensional basis for functions on the sphere, but in practical use the series is truncated to give an approximation of an arbitrary function by a finite weighted sum of basis functions.

In computer graphics, spherical harmonics are used as a form of compression, which in turn greatly accelerates computations. For instance, the incoming light at a point in space is a spherical function (since it varies with direction), but using SH its approximation can be compactly represented by a handful of coefficients (the weights for the first few basis functions). This allows computations to be performed on a vector of these SH coefficients, rather than the spherical functions themselves.

The Spherical Harmonic Series is composed of separate bands of functions, usually denoted L_0 , L_1 , L_2 etc. The L_0 band is a single constant function; the L_1 band consists of three linear functions; the L_2 band contains five quadratic functions; and so on. Higher band indices account for finer details.

1.4 Problem Definition

The problem being attempted is to estimate the light model of real world without introducing a lightprobe, the input is constrained to be just RGB images and output rendering is expected to be calculated in realtime without any constraints in the environment. The main challenge is to use the objects in real world to act like a lightprobe. Thus, we need to learn a lightprobe from the scene. A great advantage of the learning lightprobe is that practically anything can be a lightprobe. The only practical requirement is that the lightprobe should have a good distribution of surface normals.

Chapter 2

Related Work

A large body of work on both photorealistic and non-photorealistic rendering for AR exists. Coherence requires estimation of light models. Many other issues need to be seen for photorealistic AR like occlusion, shadows due to virtual objects which are not discussed as part of this report. Here, we focus on methods that provide an estimate of the current lighting in the users environment only. Work related to major light model estimation approaches relevant to our problem statement is discussed here.

2.1 Active lightprobes

The seminal work of Debevec demonstrated how to make use of a light probe in order to directly measure the incident light [1]. Light probes are placed inside an environment and provide light measures as high-dynamic range images, which are subsequently used to compute an environment map. Active light probes obtain the environment map directly. A camera with fish-eye lens or an omni-directional camera is placed directly in the scene to acquire images of all directions in one step [2]. Passive light probes make use of a reflective object, commonly a sphere, which is placed within the scene and observed by a camera. While light probes deliver environment maps at full frame rates, they are invasive and not suitable for casual use.

2.2 Inverse Rendering Equation

As an alternative Gruber et al. [3] demonstrated a probe less method that estimates incident light purely from reflections scanned in real time using a RGB-D sensor. This method estimates radiance transfer and solves inverse rendering equation to get spherical harmonics coefficients and thus render the scene using the obtained light model. Further, this approach was improved to Gruber et al. [4] [5]. These methods forgoe artificial light probes, but suffer from high computational demands that currently cannot be run on mobile devices.

2.3 Deep outdoor illumination estimation

Hold-Geoffroy et al. [6] wherein they described use of natural available data to train a CNN to predict lighting model in outdoor environment. The idea was that to use parameters from sky luminance distribution model as the output to be learned from

images via a CNN. So the neural network served as a complex mapping between input image pixels and output illumination parameters. But this approach was constrained on outdoor environments. Such sky models wont be very much useful for indoor illumination.

2.4 Learning Lightprobes

Mandl et al. [7] described an innovative approach to use CNNs to approximate the scene lighting in form of spherical harmonics (SH). The idea was that of using any known object in scene as a light probe and then creating synthetic data for the neural net using that probe and predict spherical harmonic coordinates to light the scene. The learned lightprobes work in an unmodified scene, providing unobtrusive user experience, while still retaining the computational benefits of a normal lightprobe. The incident lighting is estimated with a pre-trained convolutional neural network (CNN), which is orders of magnitude cheaper to evaluate at runtime than inverse rendering.

This approach is the closest to the problem we are attempting. In this approach, using a known object adds constraints on the environment. The only practical requirement is that the lightprobe should have a good distribution of surface normals. Our aim is to get rid of this known object and use the features in environment as features.

In this work, a separate CNN has to be trained for each pose. The model is illuminated with different spherical harmonics variations and the resulting images are used to train a set of CNN instances. The CNN results for all camera poses are stored in a database. During AR rendering a live camera frame is taken to detect the camera pose. Subsequently, the camera pose is used to select the CNN which was trained from the most similar pose in the database. Furthermore, the selected CNN then outputs estimated Spherical Harmonics, which are used to illuminate augmented 3D objects. A lot of pre-training is required in this method.

Chapter 3

Method

3.1 Synthesis of training images

A key aspect for CNN training is rich training data. Using image synthesis to supply sizeable amounts of training data without additional cost in human labor is an obvious approach. First, a known object in the scene is selected as a lightprobe. The only practical requirement is that the lightprobe should have a good distribution of surface normals. The Spherical Harmonics (SH) representation is used for modeling the light.

A mesh model of hulk is used as the known object and SH representation consisting of four bands is chosen for the light model. This SH representation consists of 16 coefficients $\mathbf{I} = [C_0, C_1, \dots, C_{15}]$ in the range $[-1, 1]$. Thus, 2^{15} illumination variations are created by setting the ambient light $C_0 = 1$ and quantizing the higher-order SH coefficients C_1 to C_{15} with two different values in the range of -1 to 1. To provide the training data, the known object is repeatedly rendered, while systematically varying the incident illumination. For a fixed camera pose, we render a set of training images \mathbf{T} , each with size 224×224 . The resulting dataset consists of 32768 images.

Note, selecting only two different values for each coefficient keeps the size of the dataset small enough for training the CNNs in an acceptable amount of time. Since only a subset of possible light conditions is used for training, the detection of very different light settings might suffer from missing training data.

3.2 Training the CNN

The CNN is directly trained on the SH coefficients, using \mathbf{T} as training data and \mathbf{I} lighting as reference for the corresponding loss function. The layout of the network is as shown in Figure 3.1.

Each image passes through five convolution-and-pooling layers with Rectified Linear Unit (ReLU) activation. The first layer is using a 5×5 filter kernel, while the remaining four layers use a 3×3 kernel size. After each convolution, we downsample the image to half of its size with max-pooling, i.e., max-pooling with a 2×2 kernel with a stride of 2, which is then followed by ReLU. At the end of the pipeline, we use a standard fully connected layer with 500 neurons and tanh activation. We use ReLUs in the convolution layers to

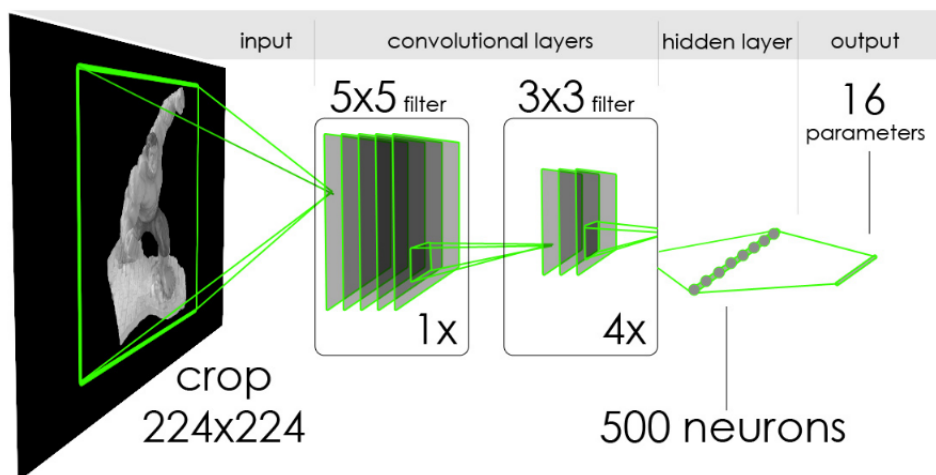


Figure 3.1: The CNN input are images with 224x224 pixels. The input images pass through five convolution layers, one hidden layer with 500 neurons and arrive at an output layer, which combines the result to 16 SH coefficients used to render the input image.

enhance the training speed, and use tanh at the end to keep the output of the network within the range of SH coefficients.

The loss function is a standard mean squared error function, thus solving a linear regression problem. For validation during training, one sixth of the training data as the validation set to have both rich training and validation data. The metric for validation is the Euclidean distance of the estimate to the ground-truth. The optimizer is Stochastic Gradient Descent (SGD) and the CNN is iterated with a maximum of 200 epochs, using a learning rate = 0.1.

3.3 Implementation Details

All the coding has been done in C++. Synthesis of training dataset is done using OpenGL methods by varying illumination based on SH representation. The CNN model is implemented using PyTorch framework. Training is performed on a desktop PC using an i5-7400 processor, 32 GB RAM and a GTX 1080Ti GPU with 11 GB VRAM.

Chapter 4

Experiments and Results

4.1 Estimation for a Fixed Pose of known object

The network layout and hyperparameters are the same as mentioned in Mandl et al. [7] and described in the method section. The **validation loss** is **0.016579**. Some of the results are as follows.

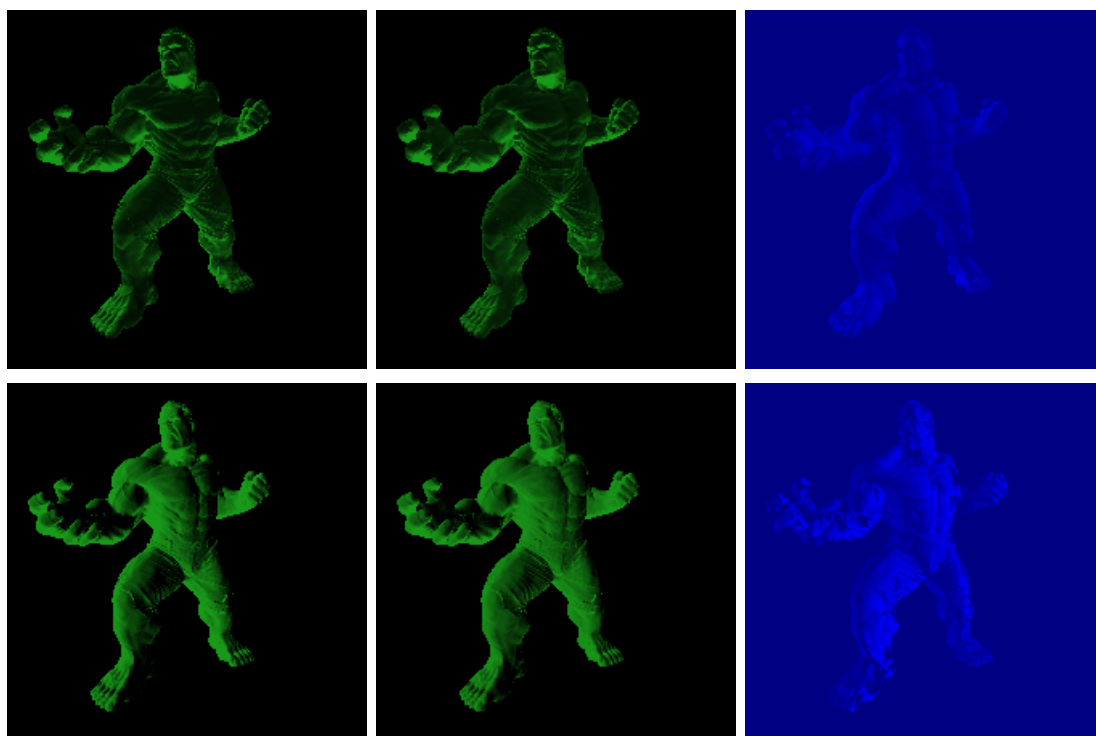


Figure 4.1: (a) Input Image, (b) Output Image (c) Absolute pixel-wise color coded difference



Figure 4.2: The colorbar for the difference image

4.2 Estimation independent of object pose

The SH representation of light model uses the normal value at each point to calculate illumination. Hence the illumination of the object is majorly dependent on the distribution of normals of the known object. The above approach uses a separate CNN for each pose, so the distribution of normals is fixed. The variation in pose causes change in the distribution of normals. A normal map of an image is the color coded image of the normals at every point. The proposed method uses the normal map as an input along with the illuminated image. For every point the normal information is available which makes this method independent of the pose of known object.

The SH representation of three bands is chosen for the light model and 36 poses are included by rotating the object by 10° . The SH representation consists of 9 coefficients. These coefficients were quantized with 2 different values of 0.5 and -0.5 to generate a training set of 2^8 images for each pose. The resulting dataset consists of 9216 images.

Each image passes through nine convolution layers with BatchNorm[11] followed by Rectified Linear Unit(ReLU) activation. At the end of the pipeline, we use a standard fully connected layer with 72 neurons and tanh activation. The optimizer is Adam[9] and cyclic learning rate with linear variation from 10^{-3} to 10^{-6} every 1000 timesteps is used. The MSE loss for validation dataset is **0.002730**. Some of the results are as follows.

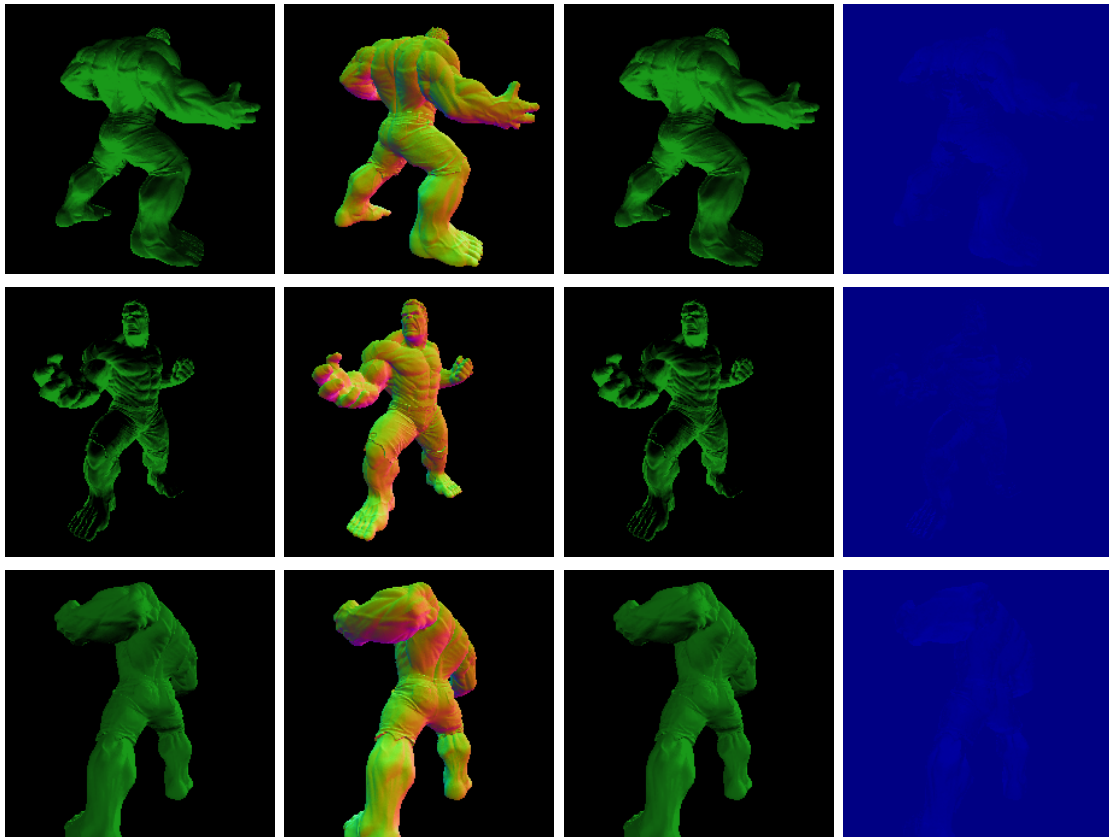


Figure 4.3: (a) Input Image, (b) Normal Map, (c) Output Image (d) Absolute pixel-wise color coded difference

4.3 Performance for different network parameters

The custom CNN network was experimented by varying the number of layers and using training strategies like batch normalization, dropout regularization[12] and SGD with warm restarts [8]. The validation loss observed is as shown in the following table.

#Layers	BatchNorm	Dropout	Cyclic Learning Rate	Validation Loss
9	✓	×	✓	0.002730
6	×	×	×	0.018124
4	✓	×	✓	0.039008
4	×	×	×	0.044377
4	×	×	×	0.0
9	✓	✓	×	0.075487
9	✓	✓	✓	0.070386

Table 4.1: Validation Loss Variation

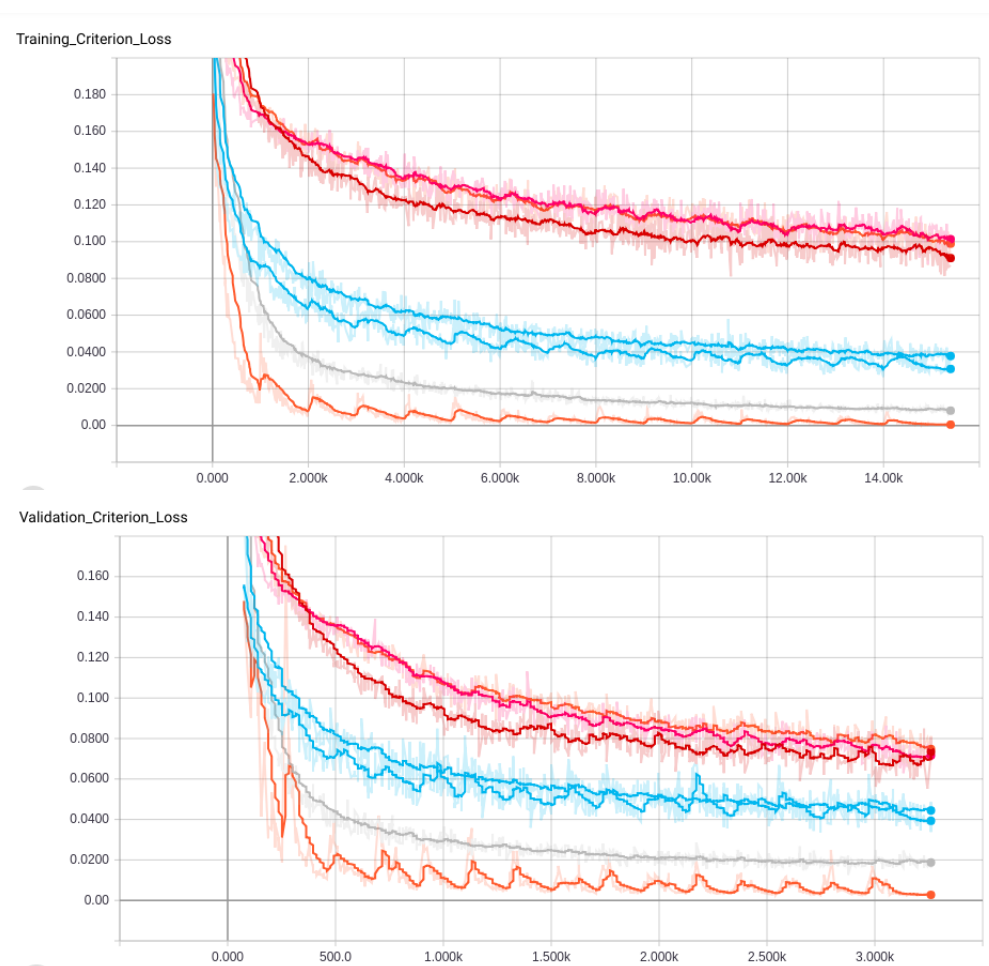


Figure 4.4: (a) Training Loss vs Timesteps, (b) Validation Loss vs Timesteps

4.4 Failed attempt: Approach with planes

The problem we are attempting involves removal of the known object from the environment. This can be achieved if the illumination information and distribution of normals is recovered from the already existing features in the environment. One way to approach this is detect small planar patches such they account to a good distribution of normals. The illumination and normals can then be used to create a known geometry/object.

The CNN is trained on different poses of a single plane with varying illumination. However, the validation loss was too high to use for coherent rendering.

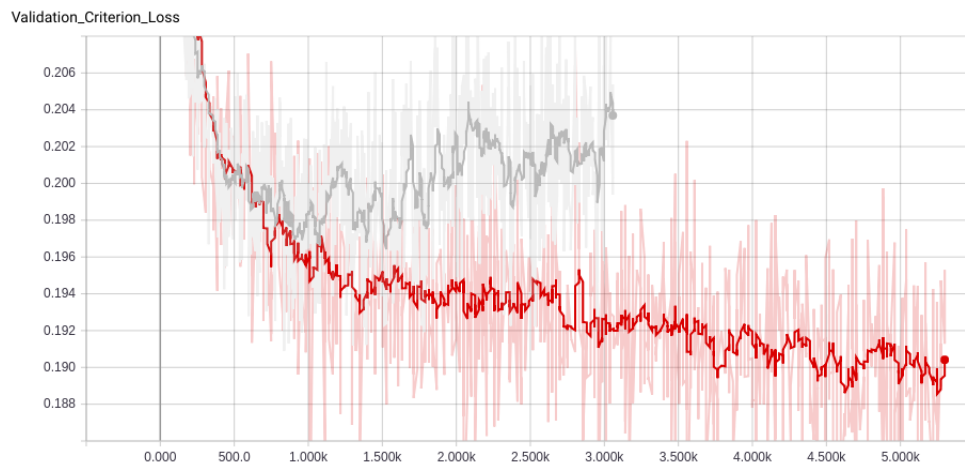


Figure 4.5: Validation Loss vs Timesteps for dataset of planes

Chapter 5

Conclusion and Future Work

At present, the CNN network estimates the light model fairly close to the original light model. The experimental results show that this method can render high quality photorealistic images. The proposed approach can cater various poses with a single CNN instead of a set of CNNs. However, the proposed method suffers from the constraint having a known object in environment.

In a future work, the results from learning different numbers of SH bands and variations of patterns for selecting values for the SH coefficients can be studied. Furthermore, to estimate light model from objects present in environment more experiments can be done. The network can be trained on a dataset of sphere as it has extremely good distribution of normals. A combination of planes from the environment can be assembled to create a sphere and use it effectively as a plausible lightprobe. Finally, this model can be integrated with an AR rendering application to render in real environment.

References

- [1] P. Debevec., 1998 Rendering synthetic objects into real scenes: Bridging traditional and image-based graphics with global illumination and high dynamic range photography.
- [2] M. Knecht, C. Traxler, O. Mattausch, W. Purgathofer, and M. Wimmer., 2010 Differential instant radiosity for mixed reality.
- [3] Gruber, L., Richter-Trummer, T., and Schmalstieg, D., 2012, Real-time photometric registration from arbitrary geometry
- [4] L. Gruber, T. Langlotz, P. Sen, T. Hoherer, and D. Schmalstieg., 2014, Efficient and robust radiance transfer for probeless photorealistic augmented reality.
- [5] L. Gruber, Jonathan Ventura, and D. Schmalstieg., 2015, Image-Space Illumination for Augmented Reality in Dynamic Environments
- [6] Hold-Geoffroy, Y., Sunkavalli, K., Hadap, S., Gambaretto, E., and Lalonde, J.-F., 2017, Deep outdoor illumination estimation
- [7] Mandl, D., Yi, K. M., Mohr, P., Roth, P. M., Fua, P., Lepetit, V., Schmalstieg, D., and Kalkofen, D., 2017, Learning lightprobes for mixed reality illumination
- [8] Ilya Loshchilov, Frank Hutter, 2017, SGDR: Stochastic Gradient Descent with Warm Restarts
- [9] Diederik P. Kingma, Jimmy Lei Ba, 2017, Adam: A Method for Stochastic Optimization
- [10] R. Green, 2003, Spherical harmonic lighting: The gritty details
- [11] Sergey Ioffe, Christian Szegedy, 2015, Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift
- [12] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, Ruslan Salakhutdinov, 2014 Dropout: A Simple Way to Prevent Neural Networks from Overfitting