

---

# Monocular Depth Estimation using Deep learning

---

Sai Karthik Vyas Akondi  
University at Buffalo  
50488251  
sakondi@buffalo.edu

## 1 Overview

Monocular depth estimation is a task that uses image data and estimates the depth or distance of the image from the camera. My idea is to use deep learning techniques like Convolution Neural Networks, Neural Style Transfer and Generative Adversarial Networks and compare them with most widely used state of the art MIDAS model. Monocular Depth Estimation is widely used in Autonomous Vehicles, UAV's and other robotic applications. Analysing different types of Deep Learning techniques to get depth maps would be very useful in these applications.

The main idea of this project is to analyze different types of deep learning techniques for Monocular Depth Estimation task. The main challenge with Monocular depth estimation is use of sequential data from a single camera which makes it challenging to abstract depth maps using simple computer vision techniques. Whereas in the case of Stereo depth estimation, where left and right sequences and images are given is much simpler.

Inputs : This Deep learning application takes RGB images from sequential data of indoor scenes Outputs : Depth maps, the outputs will be Depth matrix which will have depth matrix is a matrix with depth of each pixel from the camera.

### 1.1 Dataset

The data-set selected for this project is the NYU V2 data-set, consisting of 1449 indoor scenes with corresponding RGB images and ground truth depth maps. The data-set provides various indoor environments essential for the depth estimation task. Some of the most significant scenes are listed below:

- Scene Type: Playroom, Number of Images: 1330
- Scene Type: Living Room, Number of Images: 12242
- Scene Type: Dining Room, Number of Images: 7778
- Scene Type: Bedroom, Number of Images: 15466
- Scene Type: Office, Number of Images: 3270

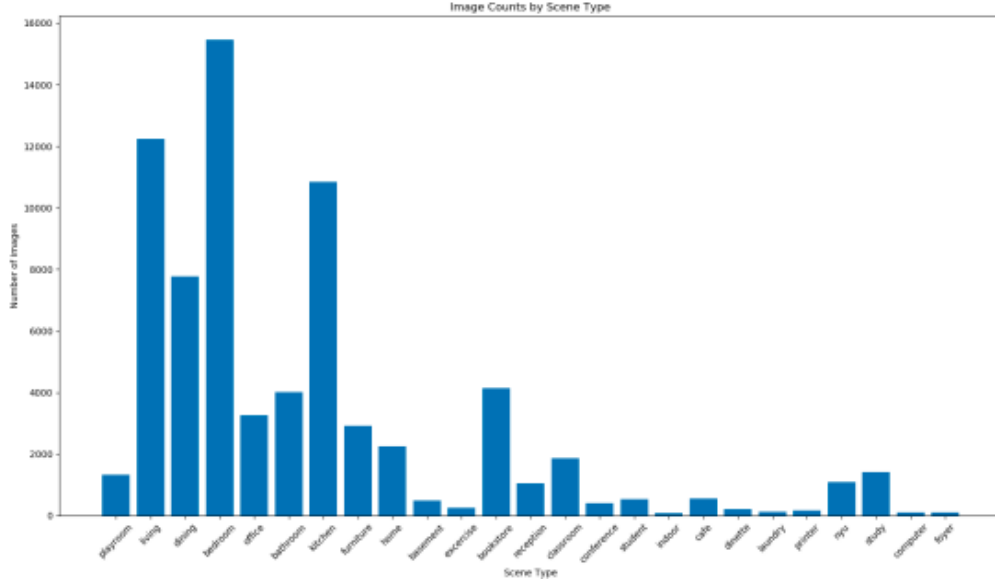


Figure 1: NYU V2 Scene distribution

28 Figure-1 shows the distribution of all scenes

## 29 1.2 State of the Art

30 Monocular depth estimation has seen significant advances in recent years, with  
 31 deep learning techniques at the forefront of these developments. Current state-of-  
 32 the-art solutions employ a variety of approaches, each offering unique insights and  
 33 advancements in depth estimation from single images. Key contributions include:

- 34 • The **MIDAS model** represents a breakthrough in depth estimation, offering  
 35 robust and accurate predictions. This model's strength lies in its ability to  
 36 generalize across different scenes and lighting conditions, making it highly  
 37 versatile for various applications.
- 38 • The work by Eigen et al. [1] entitled "*Depth Map Prediction from a Single*  
 39 *Image using a Multi-Scale Deep Network*" introduced a novel approach  
 40 that leverages multi-scale neural networks. This methodology significantly  
 41 improved the accuracy of depth predictions from single images, setting a  
 42 new benchmark for subsequent research.
- 43 • In "*Indoor Segmentation and Support Inference from RGBD Images*" [11], a  
 44 combination of image segmentation and the RANSAC algorithm is employed.  
 45 This study enhanced understanding of the indoor scenes' structural elements,  
 46 providing a foundational technique for depth estimation.
- 47 • Gatys et al. [2] in their work on "*Image Style Transfer Using Convolu-*  
 48 *tional Neural Networks*" provided a unique perspective on the use of neural  
 49 networks, influencing approaches in depth estimation.

## 50 **2 Approach**

### 51 **2.1 Deep Learning Algorithms**

52 Deep learning techniques have always been pretty useful for images tasks. Currently  
53 the following four techniques have been explored in Dept estimation task :

- 54 • **CNN with Batch normalization**
- 55 • **Pre-trained MIDAS model**
- 56 • **Neural Style transfer using VGG-19 weights**
- 57 • **Generative Adversarial Networks**

58 In the above mentioned algorithms: the CNN model with Batch normalization  
59 implemented using Pytorch was implemented without using any existing applications.  
60 The Neural Style transfer and GAN haven't been used for depth estimation till now,  
61 hence making this project novel.

### 62 **2.2 CNN with Batch normalization**

63 Contribution : Dataset class, Model, Training loop was written and implemented  
64 with minimal use online repositories and pre-trained models.

#### 65 **2.2.1 Model**

66 The CNN model uses the following layers:

- 67 • **Conv2d Layers** : Conv2D layers are the most important part of any CNN  
68 model, they take the input and perform Convolution operation on given input.  
69 When the convolution operation is performed, important features of input  
70 data is extracted into a generic way.
- 71 • **Batch normalization layers** : Batch normalization layers normalize (standard-  
72 ize) the data, and in turn the model's performance increases.
- 73 • **Transposed Convolution Layers** : These layers take the convolution outputs  
74 and then bring back the extracted features in higher dimensions. This is  
75 helpful in getting back the features that we require for getting depth map  
76 outputs.
- 77 • **Activation** : RELU activation function removes negative values and only  
78 keeps the positive one's. Acting as a filtering mechanism
- 79 • **Max-pooling** : Max pooling is used to take the max value in a given filter  
80 window, this skips rest of the unnecessary values and helps in making the  
81 training faster.

#### 82 **2.2.2 Training**

83 Hardware used for training : NVIDIA RTX 3050 , 4GB graphics card.

84 The model is trained on the whole dataset. The dataset is loaded using Pytorch  
85 Dataset and Dataloader classes. The shuffle parameter is set to false, due to a CUDA

86 error that was observed during initial training. If shuffle was set to true, the model  
87 would have performed a bit better.

88 Training is done on epoch sizes of 3 and 5, loss function Mean Squared Error and  
89 Adam optimizer. During initial runs, the model was outputting a blank image and  
90 the loss was very high. After changing the loss function to L1 loss, the model  
91 performance improved and the depths of some features was captured in the predicted  
92 image. The model performance might improve by increasing the number of epochs,  
93 but computationally it's not that easy with a 4GB graphics card. Model is saved as a  
94 pth file, which can be useful for further training and fine tuning.

## 95 **2.3 Pre-trained MIDAS model**

96 Contribution : Code for loading, transforming, and passing image through pretrained  
97 model was written. Pre-trained MIDAS model was imported from Torch hub.

98 Hardware used for training : NVIDIA RTX 3050 , 4GB graphics card.

99 MIDAS model is a part of PyTorch, a tool for building AI models, and is very good at  
100 estimating depths. This part of the project involved loading the MIDAS DPT\_Large  
101 model using torch load. The process involved loading an image using OpenCV,  
102 converting it from BGR to RGB format, and aligning it with the MiDaS model's  
103 format. After applying the necessary transformations to the image, it was fed into the  
104 MiDaS model to predict depth. To accurately represent the depth information, the  
105 model's output was resized to match the original image dimensions.

106 **Result :** MIDAS model had the best depth map outputs, which were close to the  
107 ground truth images.

## 108 **2.4 Neural Style transfer using VGG-19 weights**

109 Contribution: This was directly taken from Torch implementation of Neural Style  
110 transfer for a different application. Modifications were done to fit our application  
111 into the model, i.e changing the input features, output features, and conversion of  
112 output RGB to Grayscale to showcase depth images.

113 Hardware used for training : NVIDIA RTX 3050 , 4GB graphics card.

114 In this section, the VGG19 model was utilized with its pre-trained weights. The  
115 process began with loading and configuring the VGG19 model from PyTorch's  
116 library of pre-trained models. For image processing, content and style images were  
117 prepared using PyTorch's transform function.

118 Style Image: Existing depth images Content Image: RGB image

119 The neural style transfer implementation involved defining specific layers within the  
120 VGG19 model to track content and style losses. Content loss was computed using  
121 the Mean Squared Error between the feature maps from the content and generated  
122 images, while the style loss tracked with the Gram matrix.

123 **Result :** Neural Style transfer with VGG-19 model weights, did not output expected  
124 depth images. This might be due to the fact that, the main aim for neural style  
125 transfer is aimed for RGB to RGB conversion.

## 126 2.5 Generative Adversarial Networks

127 Hardware used for training : NVIDIA RTX 3050 , 4GB graphics card.

128 GAN for depth images is also a new research area, that hasn't been done yet.  
129 Currently this part of the application is still in progress. I'm going through some  
130 papers and some code online on GAN, to apply it to Depth estimation tasks.

## 131 3 Experimental Protocol

### 132 3.1 Dataset

#### 133 3.1.1 NYU V2 Dataset

134 The NYU Depth Dataset V2 was used for depth estimation, which is a comprehensive  
135 collection of RGB and depth images from indoor scenes.

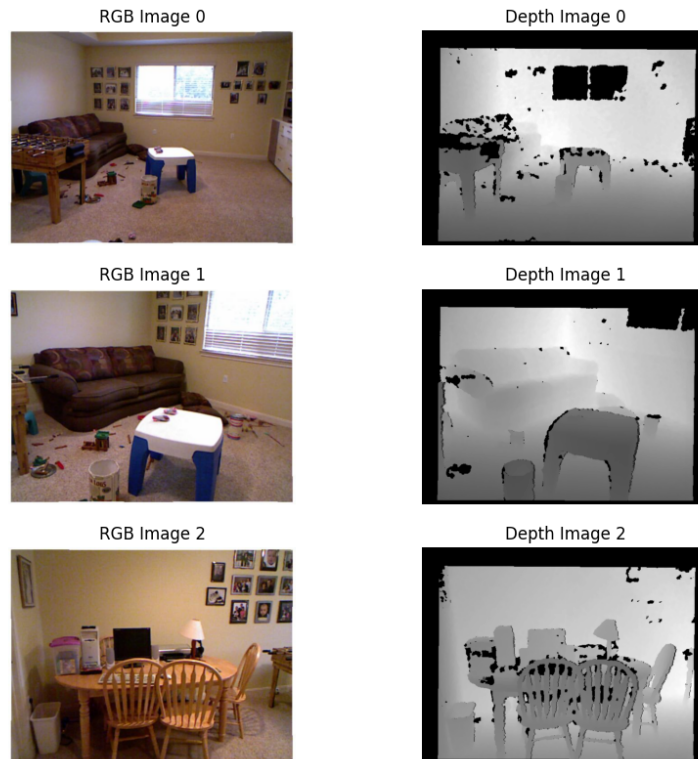


Figure 2: NYU V2 Dataset

### 136 3.1.2 Relevance to the Depth Estimation task

137 This dataset is particularly selected due to its diverse environments and scenarios,  
138 providing a better chances of diversity for training and evaluating the depth estima-  
139 tion model. For the neural style transfer, custom images were selected from NYU  
140 V2 dataset, with RGB images serving as content and depth images as style sources.

### 141 3.2 Evaluation of Success

142 The most simplest qualitative result for any depth estimation task is obtained by just  
143 comparing the ground truth depth images with that of the predicted images.

### 144 3.3 Computational Resources

#### 145 3.3.1 Hardware :

146 NVIDIA RTX 3050, 4GB graphics card was used to train the models for Depth  
147 Estimation task.

#### 148 3.3.2 Frameworks and Libraries:

149 PyTorch's library for pre-trained models (e.g., VGG19, MIDAS)

## 150 4 Results

### 151 4.1 CNN with Batch Normalization

152 The CNN model's performance was not as effective as state-of-the-art CNN models,  
153 primarily due to its less complex architecture. This is evident from the loss graph  
154 shown below:

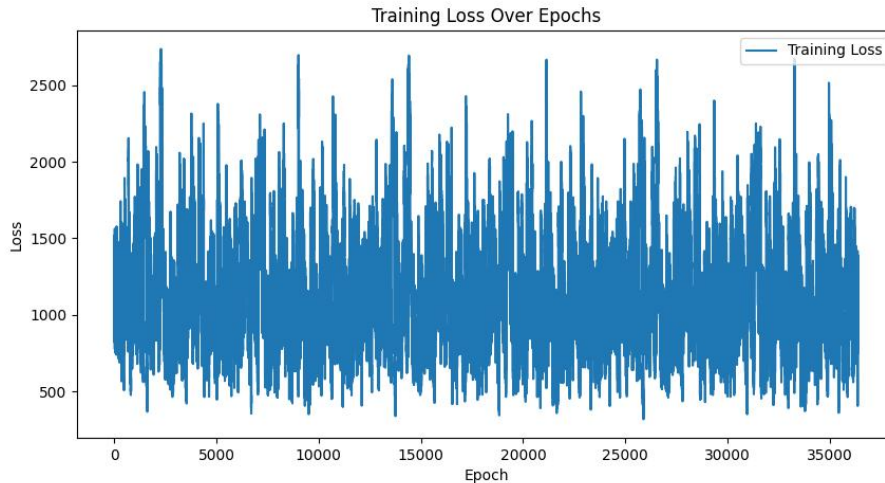


Figure 3: Loss graph of the CNN model



Figure 4: RGB Image



Figure 5: Ground truth depth image

155 Although the output images from the CNN model contained a significant amount  
156 of noise, some aspects of the depth maps were recognizable and corresponded  
157 reasonably well to the expected output.

158 Figure 2 shows the fluctuation in loss during the training. In the graph it's clearly  
159 observed that the training cross is in the range of 300 and 2300, which suggests that

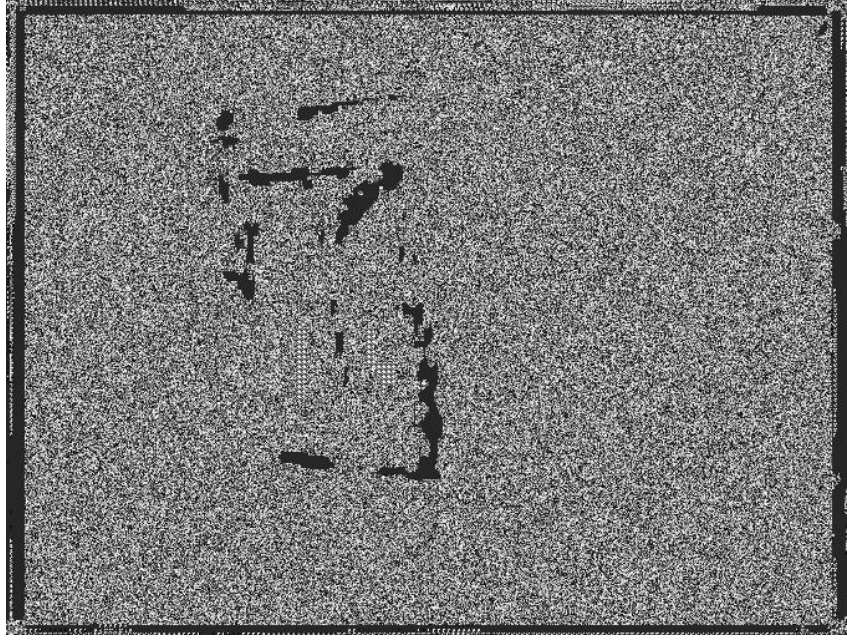


Figure 6: Predicted Image - CNN with Batch normalization

the Deep learning model should have more layers and complexity to fit the task of Depth estimation.

The output of the Deep learning model is shown in Figure-3:

Comparing the RGB (Figure 3), Ground truth depth (Figure 4) and CNN predicted depth images(Figure 5)

## 4.2 Pre-trained MIDAS Model

The Pre-trained MIDAS model showed a marked improvement in depth estimation accuracy. The depth maps generated by this model were closer to the ground truth, as illustrated in the image below:

The superior performance of the MIDAS model is attributed to its advanced architecture and the extensive training it has undergone, which is reflected in the quality of the depth maps it produces.

The pretrained MIDAS model has the best depth outputs out of all three

## 4.3 Neural Style transfer using VGG-19 weights

Neural style transfer involves fusing two images: a content image (such as a photograph) and a style image (usually an artwork), to produce a result that maintains the content of the first image but is stylistically similar to the second. VGG-19 weights were used, as the model is a bench mark in Neural Style transfer.

The novel idea was to implement Neural style transfer for our Depth estimation task. This method sought to integrate the depth information with the RGB content. However, the results indicated that this novel approach did not give fruitful results.





Figure 7: Depth map generated by the Pre-trained MIDAS model



Figure 8: Content Image

181 Considering Style Image as the Ground Truth Depth and Content image as it's  
 182 corresponding RGB image. Neural style transfer was implemented, the result was an  
 183 empty image, monochromatic, lacking any meaningful synthesis of the depth data  
 184 with the RGB content.

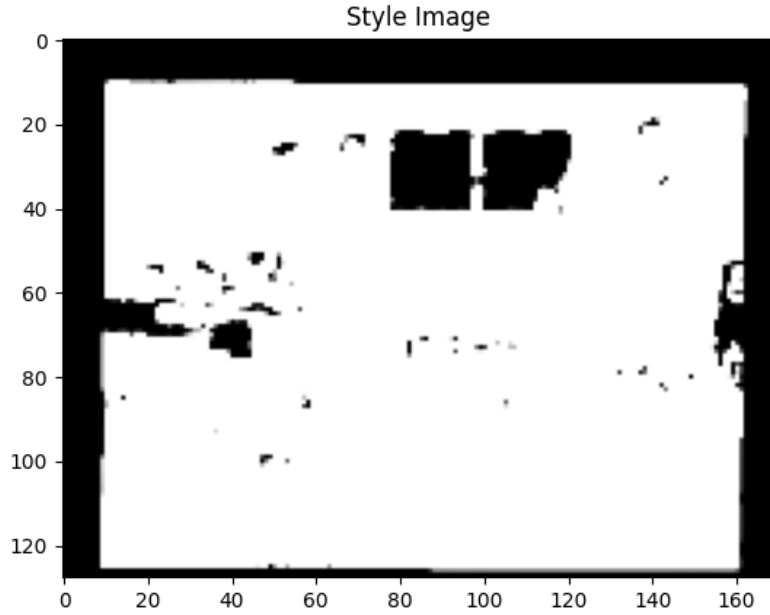


Figure 9: Style Image

#### 185 **4.4 GAN for Depth Estimation task**

186 In this experiment to use Generative Adversarial Networks (GANs) for the task  
 187 of estimating depth in images. The complexities involved in adapting GANs for  
 188 monocular depth estimation task required more time and couldn't be completed in  
 189 this time frame. Hence this section of the project was not completed as expected and  
 190 therefore no results.

### 191 **5 Analysis**

#### 192 **5.1 CNN with Batch normalization**

193 Improving the complexity of the

##### 194 **5.1.1 Advantages**

- 195 • The advantage of using a CNN-based approach with batch normalization is  
 196 its simplicity and ease of implementation.
- 197 • It is computationally less intensive compared to more complex models,  
 198 making it feasible for training on hardware with limited resources.

##### 199 **5.1.2 Limitations**

- 200 • The main limitation of this approach is its limited capacity to capture complex  
 201 depth information. The simple architecture may not be sufficient to learn  
 202 depth patterns.

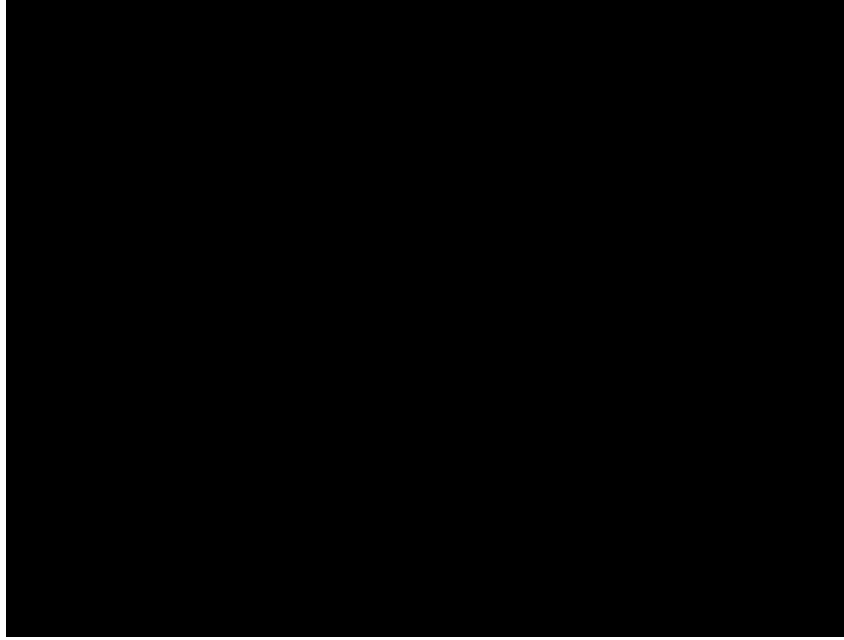


Figure 10: Output of Neural Style transfer - blank image

### 203 **5.1.3 Improvements**

- 204 • Although CNN with Batch normalization did not give expected results, this  
205 technique can be improved by augmenting data and increasing the number  
206 of epochs for training.

## 207 **5.2 Pre-Trained MIDAS model**

### 208 **5.2.1 Advantages**

- 209 • The depth maps generated by the MIDAS model closely resemble ground  
210 truth depth images, indicating high-quality results.
- 211 • Pre-trained MIDAS was the best model for Depth estimation task, it gave out  
212 best results both in terms of accuracy and the capturing of depth features.

### 213 **5.2.2 Limitations**

- 214 • One of the most important limitation to consider is that the Pre-trained  
215 MIDAS model was trained on the same NYU-V2 dataset. So it might be  
216 possible that the features might have over-fitted.

### 217 **5.2.3 Improvements**

- 218 • Training on different types of data will be the most important improvement  
219 for this model.

## 220 5.3 Neural Style Transfer

### 221 5.3.1 Advantages

- 222 • The main advantage about this technique is it's novelty. Neural Style Transfer  
223 is a creative approach to depth estimation, attempting to combine style (depth)  
224 with content (RGB) images to generate unique results.
- 225 • It uses pre-trained VGG-19 weights, which are known for their effectiveness  
226 in image processing tasks.

### 227 5.3.2 Limitations

- 228 • Neural Style transfer was mainly developed for converting images to paint-  
229 ings, here both the content and style images are RGB.
- 230 • In our application the style image is gray-scale and content image is RGB,  
231 research has to be done in this area for effectively applying Neural Style  
232 transfer for Depth task.

### 233 5.3.3 Improvements

- 234 • Instead of using VGG-19 pre-trained model, other models specifically de-  
235 signed for depth task can be used.

## 236 6 Discussion

237 In summary, this project aimed to explore different deep learning methods for  
238 monocular depth estimation, an important in computer vision with applications in  
239 robotics, autonomous vehicles, Four techniques were researched and three of them  
240 were implemented :

- 241 • CNN with Batch normalization
- 242 • MIDAS pre-trained model
- 243 • Neural Style transfer using VGG-19 weights
- 244 • GANs -> this was not implemented

245 The pre-trained MIDAS model produced the most promising results, almost resem-  
246 bling ground truth depth maps. However, the custom CNN model had limitations due  
247 to its simplicity. Adapting neural style transfer for depth estimation faced challenges  
248 related to image style differences between RGB and depth images. The application  
249 of GANs for depth estimation is still an ongoing research. To improve accuracy,  
250 future work should explore more complex CNN architectures, diverse data sets, and  
251 innovative techniques bridging RGB and depth information

## References

- [1] D. Eigen, C. Puhrsch, and R. Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Advances in Neural Information Processing Systems*, 2014.
- [2] L. A. Gatys, A. S. Ecker, and M. Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [3] C. Godard, O. Mac Aodha, and G. J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [4] C. Godard, O. Mac Aodha, M. Firman, and G. J. Brostow. Monodepth2: Unsupervised monocular depth estimation with left-right consistency. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2019.
- [5] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, 2014.
- [6] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab. Deeper depth prediction with fully convolutional residual networks. In *Proceedings of the 4th International Conference on 3D Vision (3DV)*, 2016.
- [7] F. Liu, C. Shen, and G. Lin. Learning depth from single monocular images using deep convolutional neural fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(10):2024–2039, 2015.
- [8] F. Liu, C. Shen, G. Lin, and I. Reid. Multi-scale continuous crfs as sequential deep networks for monocular depth estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [9] R. Ranftl, K. Lasinger, D. Hafner, K. Schindler, and J.-M. Frahm. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [10] A. Saxena, M. Sun, and A. Y. Ng. Make3d: Learning 3d scene structure from a single still image. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(5):824–840, 2009.
- [11] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus. Indoor segmentation and support inference from rgb-d images. In *European Conference on Computer Vision (ECCV)*, 2012.

## Online Resources

- a. PyTorch. (n.d.). *Neural Style Transfer Tutorial*. Retrieved from [https://pytorch.org/tutorials/advanced/neural\\_style\\_tutorial.html](https://pytorch.org/tutorials/advanced/neural_style_tutorial.html)

- 290 b. PyTorch. (n.d.). *MiDaS Model for PyTorch*. Retrieved from [https://](https://pytorch.org/hub/intelisl_midas_v2/)  
291 [pytorch.org/hub/intelisl\\_midas\\_v2/](https://pytorch.org/hub/intelisl_midas_v2/)
- 292 c. The Computer Vision Foundation. (n.d.). *Open Access to Computer Vision*  
293 *Research*. Provides access to a wide range of research papers and resources  
294 in computer vision. Retrieved from <https://www.thecvf.com/>
- 295 d. NYU Depth Dataset V2. (n.d.). *NYU Depth Dataset V2*. A dataset of  
296 indoor scenes for depth estimation. Retrieved from [https://cs.nyu.edu/](https://cs.nyu.edu/~silberman/datasets/nyu_depth_v2.html)  
297 [~silberman/datasets/nyu\\_depth\\_v2.html](https://cs.nyu.edu/~silberman/datasets/nyu_depth_v2.html)