

SmallDepthMask: Large DataSet Generation for Monocular Depth Estimation and Foreground Segmentation from few Internet Images

Anonymous CVPR 2021 submission

Paper ID ****

Abstract

Segmentation of the desired object along with depth estimation is useful in various applications like robotics and autonomous navigation. Any deep learning workflow to segment the desired foreground object in a scene require significant training data. The data generation process usually involve expensive hardware like RGB-D sensors, Laser Scanners or significant manual involvement. This paper presents a novel way to utilize only a small number of readily available png images with transparency for the foreground object, and representative background images from the internet and combine them to generate a huge dataset for deep learning utilizing current state of the art monocular depth estimation and segmentation techniques. Few example applications show the efficacy of the training data on detecting cattle on road for autonomous driving application, etc. The baseline models exhibit strong generalization to real scenarios.

1. Introduction

Expand the abstract with appropriate references. We need references for: 1. Monocular Depth 2. Image segmentation 3. Deep learning requiring large datasets 4. Generating data sets is cumbersome 5. Ours is novel work that uses few images to generate huge dataset that generalizes well Depth estimation and semantic segmentation are often used together in many vision tasks [10] like autonomous navigation of agents, augmented reality, self driving cars and other robotics applications. In all these application, identification of desired objects precisely in the scene and its depth estimate from the camera are crucial for safe and effective navigation. Modern RGB-D sensors like OAK-D are capable of simultaneously running advanced neural networks while providing depth from two stereo cameras and color information from a single 4K camera in the center. Deep learning based techniques using convolution neural nets have effective solution in both depth estimation and

semantic segmentation. In general for high accuracy outcomes, a deep learning network is dependent on large training dataset availability. To gather such data itself incur high cost and time. For specific applications requiring several foreground objects against variety of backgrounds become even more challenging in terms of simulating those scenarios. Synthetic datasets using Virtual Reality have been proposed to that end.

Recent research indicates effective use of readily available images on the internet to curate training data. This paper introduces SmallDepthMask, a way to curate custom dataset containing millions of images by multiplexing desired foreground objects over representative background scenes, while also generating corresponding depth and foreground mask images. This significantly reduces the cost and time overheads. The authors also experiment by creating baseline models for several application contexts and show that the generated data successfully generalizes to detect relevant objects in real scenes. Multiplexing, combined with random cropping, scaling and translation, make the data-generation fast and effective. With only 100 pair of background and foreground images, the authors generate 4 million image triplets and effectively leverage existing SOTA models for depth estimation and semantic segmentation.

The main contributions of this paper are the following:

1. A novel effort to mix and match foreground and background images reducing the need for complex scene generation for data curation.
2. Curate large dataset to effectively train models for custom applications of detecting depth and mask for specific foreground objects over any target background, from a limited input of ready available internet images.
3. Combine image, depthmap and foreground mask in a single dataset using current SOTA models for depth estimation and semantic segmentation.
4. To release curated dataset and the trained models making them publicly available. Researchers can use this

single dataset to do segmentation, train models to predict depth, or to predict both depth and mask.

Figure 1 represents a few representative examples of a dataset to help detect cattle on road, which is very common on Indian roads, leading to several accidents involving loss of life and property. The generated datasets and the trained models are publicly available.

2. Related Work

A depth image is an image channel in which each pixel relates to a distance between the image plane and the corresponding object in the RGB image. Monocular depth gives information about depth and distance and the Monocular Depth Estimation is the task of estimating scene depth using a single image[1]. Image Segmentation is the process of partitioning an image into multiple segments and it can be used for locating objects, boundaries [3]. RGBD image is a combination of a RGB image and its corresponding depth image[19].

Depth information is integral to many problems in robotics including mapping, localization and obstacle avoidance for terrestrial and aerial vehicles, autonomous navigation, and in computer vision, including augmented and virtual reality[11]. RGBD datasets usually collected using depth sensors, monocular cameras and LiDAR scanners are expensive and data collection is a time consuming job. The wellknown datasets for monocular 3D object detection are Context-Aware MixEd ReAlity (CAMERA), Objectron, Kitty3D, Cityscape3D, Synthia, etc. [] and these datasets have limitations like indoor only images, small number of training examples and sparse sampling.

To address the issues usage of expensive devices and small number of training examples, this paper proposes a technique to come up with a custom dataset by using existing accurate depth predictor models, like High Quality Monocular Depth Estimation via Transfer Learning(nyu.h5) [2].

A variety of RGBD datasets in which images are paired with corresponding depth maps(D) have been proposed through the years. Some of the frequently used RGBD datasets are the Kitti dataset [7], the Synthia dataset [14], Make3D dataset [15], NYU dataset [16]

The dataset Kiiti [7] is the well known RGB-D dataset collected using a vehicle equipped with a sparse Velodyne VLP-64 LiDAR scanner and RGB cameras, and features street scenes in and around the German city of Karlsruhe. The Primary application of this dataset involves perception tasks in the context of self-driving. Synthia [14] is a street scene dataset with depth maps of synthetic data, requiring domain adaptation to apply to real world settings. Cityscapes [5] provides a dataset of street scenes, albeit with more diversity than KITTI. Sintel [12] is another syn-

thetic dataset which mainly comprises of outdoors scenes.

Megadepth [9] is a large-scale dataset of outdoor images collected from internet, with depth maps reconstructed using structure-from-motion techniques, but this dataset lacks in ground truth depth and scale. The RedWeb [18] dataset provide depth maps generated from stereo images which are freely available in large-scale data platforms such as Flickr. The datasets MegaDepth and RedWeb can be easily computed with the existing MVS methods.

Make3D [15] provides RGB and depth information for outdoor scenes. The NYUv2 dataset [16] is widely used for monocular depth estimation in indoor environments. The data was collected with a Kinect RGBD camera, which provides sparse and noisy depth returns. These returns are generally in-painted and smoothed before they are used for monocular depth estimation tasks. As a result, while the dataset includes sufficient samples to train modern machine learning pipelines, the “ground-truth” depth does not necessarily correspond to true scene depth.

Most of the existing datasets consists of indoor images, or outdoor images of city streets. For every specific application, like detecting animals roaming on roads for self driving or assisted driving cars, or people inside a room for autonomus room cleaners etc. researchers need to curate specific dataset to train relevant deep learning models. The present work makes the task of curating dtaset extremely simple and cost effective.

3. Dataset Generation Method

The curated dataset must have following objectives:

1. dataset which dedicatedly includes foreground object.
2. dataset should drive deep learning models and generalize.
3. dataset should provide accurate dense depth maps.
4. dataset should provide foreground mask
5. dataset can be stored offline or generated online during training phase dynamically

4. Method

4.1. Data Acquisition

The first step to curate data is to determine a target application scenario and thus determine the foreground object(s) and the representative background context. At the same time the dataset must have sufficient variability to include majority of the types and views that the trained deep network may see when deployed.

We propose to download or take RGB image of n foreground object(s) and m background images (we used $n =$



Figure 1. Sample Record which contains the background image, a cow overlaid on top of background, its mask and depth images

$m = 100$) balancing the types and views. For example, for cattle on roads dataset, we chose several cow, bull and calf types, individual or in group, sitting, standing or walking, and from various angles. Similarly for background, we chose backgrounds of streets, storefronts, main roads, highways, markets, railway tracks, landscapes, garbage piles etc. PNG images with transparency are readily available on the internet for almost any desired foreground object. Such images will easily allow to generate foreground mask from non-transparent pixels. If not, tools like GIMP [8], combined with deep learning foreground extractors¹ that uses a combination of Image based techniques and DNN to separate foreground from background) can help generate the required PNG foreground images.

Fig. 2 shows few of the sample scene and foreground images used for the creation of this dataset.

4.2. Multiplexing and Depth Generation

This step is to place each foreground object on to several background images generating a fg-bg image and the corresponding mask corresponding to the foreground placement and scale. Depth is computed from the fg-bg image via the model proposed by Ibraheem Alhashim et al. in their paper titled "High Quality Monocular Depth Estimation via Transfer Learning" [2]². This model takes 448×448 size images as input, hence we resize all background images to this size while maintaining their aspect ratio.

The data generation process is completely online and produces one batch of images for training a deep model. By repeating one foreground object k times for each background image and repeating another k times with horizontally flipped version of the same foreground, one can generate $2kmn$ fg-bg images. For $k = 20$ and $n = m = 100$ this becomes 400,000 fg-bg images. Algorithm 1 describes the data generation process

¹<https://www.remove.bg/>

²source code for depth estimation model: <https://github.com/ialhashim/DenseDepth/blob/master/DenseDepth.ipynb>

4.3. Data Curation and Processing

Fig. 3 shows the three outcomes of the above algorithm for a set of images.

4.3.1 fg_bg images

To generate fg_bg images, the foreground image is overlaid on background images randomly for 20 times. To place a foreground image on the background, a center point (x, y) in the range of 0 to 447 is randomly picked, and a scale between 0.3 to 0.6, which identifies the area overlapped by foreground image on the background image is also randomly picked. Next the foreground image is scaled and placed on top of background image centered at (x, y) . Save this overlaid image with 224×224 resolution. As the number of foreground images are 200, the total number of overlaid images per background becomes 4000. By repeating the same procedure for all the 100 background images, the total number of fg on bg images becomes 400000. A set of sample images after overlaying foreground on background are shown in Fig.

4.3.2 masks of fg_bg images

The mask is calculated for every fg_bg image by setting a binary image to transparency channel of foreground image. These 400000 images are also stored in the particular folder of the dataset and a set of sample masks of fg on bg images are shown in Fig.

4.3.3 depth maps of fg_bg images

We used nyu.h5 model for depth calculation from Depth estimation proposed by [2]. This model requires input images to be of 448×448 resolution and produces 224×224 size depth image. A set of sample depth images are shown in fig. 3.



Figure 2. Scene and foreground object images

Algorithm 1: Generate Dataset($[bgimages]$, $[fgimages]$, k , b)

input : m Background Image paths, $2n$ Foreground Image paths, multiplexing factor k , batch size b must be multiple of k

output: Yield $2kmn$ fg-bg, mask and depth images in batches of size b

for offline use it creates 3 folders with fg_bg, mask and depth each having $2kmn$ images;

```

for bg ← 1 to m do
  for fg ← 1 to 2n do
    for i ← 1 to k do
      croppedbg ← take maximal random crop of 448 × 448 from bg without affecting the aspect ratio;
      randomly pick a center point (x, y) in range [0, 447];
      randomly pick a scale in range [0.3, 0.6] (ratio of area fg covers bg);
      create fg - bg image by resizing the fg to scale and place it on top of croppedbg centered at x, y
      calculated;
      calculate binary mask from current placement of fg by thresholding transparency channel. save
      fg - bg image and mask add fg - bg image to a batch;
    if b new fg-bg images generated then
      run depth model on batch and save corresponding depth images
  
```

4.3.4 Directory Structure**5. Example Applications****5.1. Detecting cattle on roads**

Describe importance and if possible related work and how deep learning is not yet applied. Give example images

We can include adaptive placement of foreground to be more realistic and also occlusion handling. However note that this is only preliminary. The segmentation I used gave us road, sky, tree etc. If I remember correctly, I used the bottommost row in image that have a threshold number of sky pixels. And based on the max depth and span of the non sky part determined a formula to scale. The cattle was placed on ground area only. Finally, any objects that come on top of cow region in segmentation are put on top in the generated image.

I can formally write the algo tomorrow.

5.1.1 Baseline Model

Discuss the model with loss function used and results along with generalized results

5.2. Example 2**5.3. Example 3****5.4. Data Statistics**

Every image (fgbg, mask, depth) in the dataset are of size 224 X 224. The distribution of fgbg, mask and depth values for XYZ dataset are shown in fig. The dataset has 400000 records. A train-test-split of (70-30) gives a training set size of 280000 and 120000. A sample record in the dataset contains paths to all the images as shown below.

```

('./data/bgimages/bgimg099.jpg',
 './data/out2/images/fgbg392483.jpg',
 './data/out2/masks/mask392483.jpg',
 './data/out2/depth/fgbg392483.jpg')

```

6. Experiments

In this section, we provide a baseline for monocular depth estimation and segmentation on the XYZ dataset. The state-of-the-art models for image segmentation are variants of U-Net and fully convolutional networks (FCN)[6]. long skip connections are used to skip features from the contracting path to the expanding path in order to recover spatial in-



Figure 3. Three images resulted from the algorithm used. (top) A scene image on which a foreground object is positioned at random location with random scale, (middle) respective mask for the scene image, and (bottom) calculated depth by using a model.

formation lost during downsampling [20]. Short skip connections can be used to build deep FCNs. By using both long and short skip connections we proposed a model following U-Net architecture. This model has one encoder and two decoders, each meant for mask and depth prediction.

6.1. Model

We have designed a model with one encoder and two decoders with skip connections. The architecture is shown in Fig 4. The total no. of parameters of this model are 5,525,568. We have trained this model on the entire XYZ dataset from scratch. During training the network is trained with the batch size of 64 for 10 epochs using SGD optimizer [4]. We have used OneCycleLR scheduler [17] with a maximum Learning rate of 0.1. This made the initial learning rate as 0.0099. The Deep Convolutional Neural Networks encoder is fed with a image (224 X 224) and the first decoder outputs a mask image and second decoder outputs a depth image. To reduce overfitting[13], this work employed Random Rotation, Random Grayscale, Color Jitter, random horizontal flips and random channel swaps for data augmentation.

The Loss is calculated with the help of L1 loss and Structural Similarity (SSIM) at both the decoders. We have also employed regularization for weight penalty.

For training our network with two decoders, we defined the same loss function L between y and \hat{y} as the weighted sum of two loss function values.

$$L(y, \hat{y}) = \lambda L_{term1}(y, \hat{y}) + (1 - \lambda) L_{term2}(y, \hat{y}) \quad (1)$$

The first loss term $L_{term1}(y, \hat{y})$ is the point-wise L1 loss defined on mask values at the first decoder and on depth values at the second decoder.

$$L_{term1}(y, \hat{y}) = \frac{1}{n} \sum_{x=1}^n |y_i - \hat{y}_i| \quad (2)$$

we have also used weight decay... do we need to add it in the form of equation??????????

The second loss term $L_{term2}(y, \hat{y})$ uses a commonly used metric for image reconstruction task i.e., SSIM. Many recent too depth prediction CNNs employed this metric. The loss term is redefined as shown in equation as SSIM has an upper bound of one.

$$L_{term1}(y, \hat{y}) = \frac{1 - SSIM(y, \hat{y})}{2} \quad (3)$$

Different weight parameters λ were tried and we have ended with a value $\lambda = 0.84$. The final loss function is as follows.

$$L(y, \hat{y}) = 0.84 * L_{term1}(y, \hat{y}) + 0.16 * L_{term2}(y, \hat{y}) \quad (4)$$

6.2. Evaluation

6.3. Analysis

7. Conclusion

8. Acknowledgement

This paper and the research behind it would not have been possible without the exceptional support and computing facilities of my Institution, Vishnu Institute of Technology. Mention TSAI

References

- [1] Abdullah Abuolaim and Michael S Brown. Defocus deblurring using dual-pixel data. In *European Conference on Computer Vision*, pages 111–126. Springer, 2020. 2
- [2] Ibraheem Alhashim and Peter Wonka. High quality monocular depth estimation via transfer learning. *arXiv preprint arXiv:1812.11941*, 2018. 2, 3
- [3] Catalin Amza. A review on neural network-based image segmentation techniques. *De Montfort University, Mechanical and Manufacturing Engg., The Gateway Leicester, LE1 9BH, United Kingdom*, pages 1–23, 2012. 2

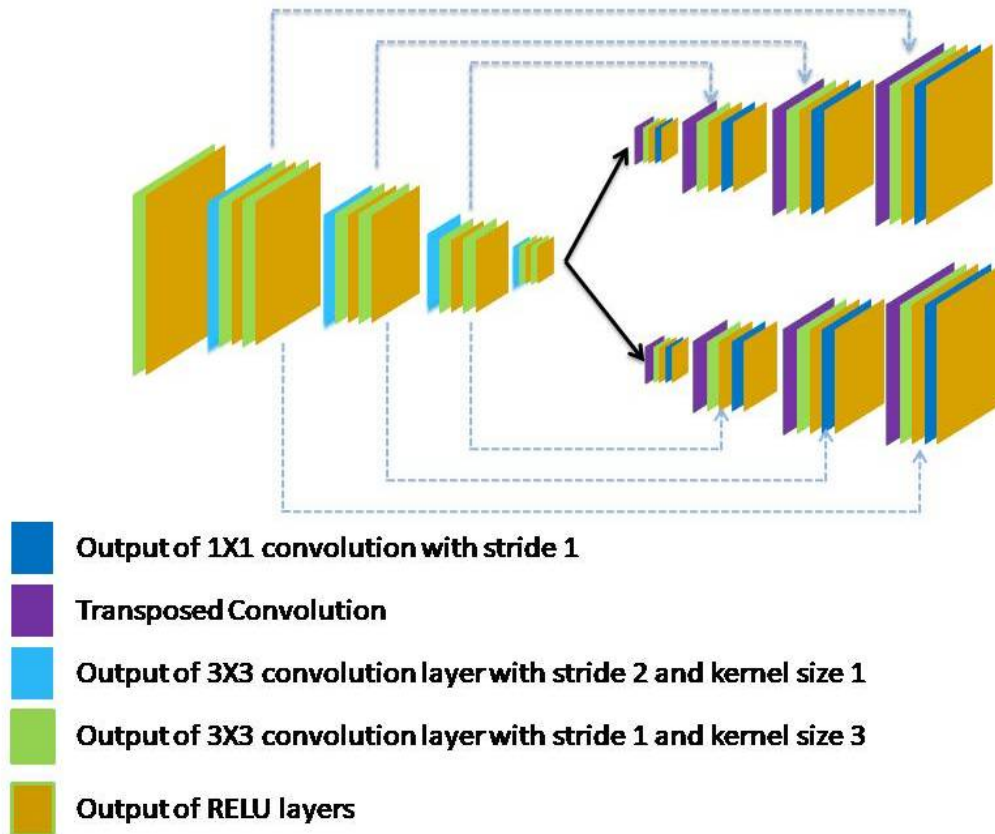


Figure 4. Network Architecture

- [4] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010*, pages 177–186. Springer, 2010. 5
- [5] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. 2
- [6] Michal Drozdal, Eugene Vorontsov, Gabriel Chartrand, Samuel Kadoury, and Chris Pal. The importance of skip connections in biomedical image segmentation. In *Deep learning and data labeling for medical applications*, pages 179–187. Springer, 2016. 4
- [7] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013. 2
- [8] Ian M Howat, A Negrete, and Benjamin E Smith. The greenland ice mapping project (gimp) land classification and surface elevation data sets. *The Cryosphere*, 8(4):1509–1518, 2014. 3
- [9] Zhengqi Li and Noah Snavely. Megadepth: Learning single-view depth prediction from internet photos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2041–2050, 2018. 2
- [10] Zhe Lin, Scott D Cohen, Peng Wang, SHEN Xiaohui, and Brian L Price. Joint depth estimation and semantic segmentation from a single image, July 10 2018. US Patent 10,019,657. 1
- [11] Eric Marchand, Hideaki Uchiyama, and Fabien Spindler. Pose estimation for augmented reality: a hands-on survey. *IEEE transactions on visualization and computer graphics*, 22(12):2633–2651, 2015. 2
- [12] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4040–4048, 2016. 2
- [13] Luis Perez and Jason Wang. The effectiveness of data augmentation in image classification using deep learning. *arXiv preprint arXiv:1712.04621*, 2017. 5
- [14] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio M Lopez. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3234–3243, 2016. 2
- [15] Ashutosh Saxena, Min Sun, and Andrew Y Ng. Make3d: Depth perception from a single still image. In *AAAI*, volume 3, pages 1571–1576, 2008. 2

- [16] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *European conference on computer vision*, pages 746–760. Springer, 2012. 2
- [17] Leslie N Smith. A disciplined approach to neural network hyper-parameters: Part 1—learning rate, batch size, momentum, and weight decay. *arXiv preprint arXiv:1803.09820*, 2018. 5
- [18] Ke Xian, Chunhua Shen, Zhiguo Cao, Hao Lu, Yang Xiao, Ruibo Li, and Zhenbo Luo. Monocular relative depth perception with web stereo data supervision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 311–320, 2018. 2
- [19] Yinda Zhang and Thomas Funkhouser. Deep depth completion of a single rgb-d image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 175–185, 2018. 2
- [20] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++: Redesigning skip connections to exploit multiscale features in image segmentation. *IEEE transactions on medical imaging*, 39(6):1856–1867, 2019. 5

702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755