

# Image Segmentation and Dimension Analysis

Naman Goyal  
2015csb1021@iitrpr.ac.in  
Koustav Das  
2015csb1017@iitrpr.ac.in

Department of Computer Science and Engineering,  
Indian Institute of Technology Ropar  
Department of Computer Science and Engineering,  
Indian Institute of Technology Ropar

## Abstract

Automated dimension analysis from an image is of immense importance in today's world. Current projects present in this domain do not provide a robust and easy to use UI. In the following paper we would like to present a project that creates a robust software that works for most of the pictures while keeping in mind the ease of access. At the present stage this provides all the basic dimension that are important.

## 1 Introduction

Around the world we find an rapid automation in all sectors of the industry. But still there lies one important industry where most of the work is done manually. Welcome to the packaging industry. Here an user needs to measure the height, width,length and various other dimensions manually. We wanted to tackle this problem with an easy to use software.

Currently there is only a few apps that are currently present in this domain. But there scope in practical sense is extremely limited. Most of them can determine the height of the phone or its distance from an object.

In our present work we propose a working model that lets us take a photo of the desired object with a reference object. The desired object is selected and further the reference objects help us to extract the real life dimensions.

## 2 Related Work

The major literature is based on methods for segmentation. The supervised techniques to semantically segment learn a mapping from pixel to class. The most of the techniques here are employ neural networks like Fully Convolution Networks, Segmentation networks which employ encoder-decoder networks, Dilated Convolution [7].

There are also unsupervised techniques using super-pixels, k-means and other clustering techniques.

While the original focus was on semantic segmentation of the image. Different supervised methods which models pixel to class mapping were tested such as DeepLab, a state-of-art deep learning system for semantic image segmentation built on top of Caffe. It combines densely-computed deep convolutional neural network (CNN) responses with densely connected conditional random fields (CRF) on MS COCO dataset.

But soon it was realized that semantic segmentation is not fit for task since it can only segment objects on which present with corresponding labels in the training dataset. Hence whenever a new object is given; any pretrained model needs to be trained again and cannot work on the fly.

The focus was then shifted on segmenting the object using initial seeds. While k-means and superpixels segment the whole image into a number of cluster or super-pixels; we required actually to extract the foreground from background.

Finally Lazy Snapping was used an interactive user input based technique based on Graph cut that can be formulated in terms of energy minimization can be approximated by solving a maximum flow problem in a graph. The idea was presented in ACM SIGGRAPH 2004 [3].

## 3 Literature Review on Segmentation Methods

This subsection presents a review of literature on pixel level semantic segmentation of images. Various techniques are compared including earlier works like unsupervised learning and decision forest. Then the recent works such as fully convolution networks, SegNet and dilated convolution are presented. The important data-sets and metrics are reported. Then dimension analysis is discussed and a pipeline model is proposed.

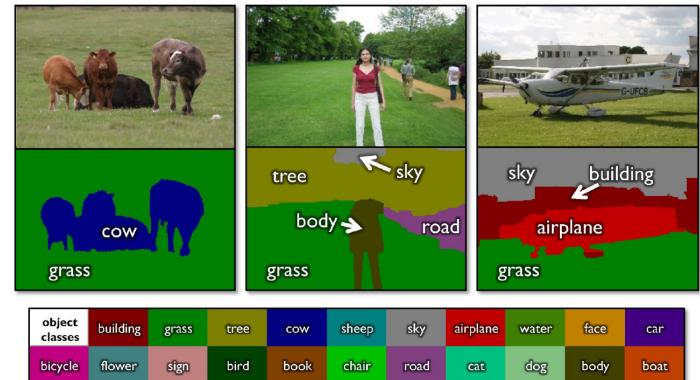


Figure 1: Semantic segmentation

### 3.1 Earlier Work on Segmentation

The earlier work can be mainly divided into

#### 3.1.1 Unsupervised Segmentation

These are non semantic approach using clustering algorithms and Graph based image segmentation.

Clustering algorithms can directly be applied on the pixels, when one gives a feature vector per pixel. Two clustering algorithms are k-means and the mean-shift algorithm.

Graph-based image segmentation algorithms typically interpret pixels as vertices and an edge weight is a measure of dissimilarity such as the difference in color.

**Pros** Few parameters to prune. Faster since no training phase required.

**Cons** Accuracy is hard to improve above a threshold. Semantics information is lost.

#### 3.1.2 Random Decision Forests

This type of classifier applies techniques called ensemble learning, where multiple classifiers are trained.

There are two techniques either the feature or training data bagging. It is observed that an ensemble/ Random Forest from random sampling of features works very well, where the classifiers. are decision trees. A decision tree is a tree where each inner node uses one or more features to decide in which branch to descend. Each leaf is a class.

**Pros** One strength of Random Decision Forests compared to many other classifiers like SVMs and neural networks is that the scale of measure of the features can be arbitrary.

**Cons** Training takes significant time.

### 3.2 Recent Works

#### 3.2.1 Fully Convolution Networks

The image is convoluted with a kernel which covers entire image. It works similar to patch model but is often much faster.

It does not have any of the fully-connected layers at the end, which are typically use for classification. Instead, it uses convolution layers to classify each pixel in the image. [4]

Upsampling is done either through "de-convolution" layers or short-cut connections.

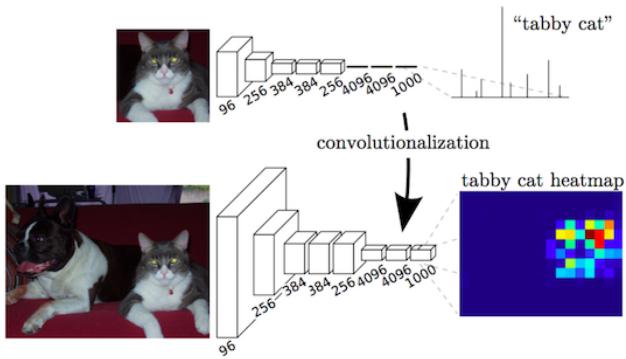


Figure 2: Fully Convolution Networks

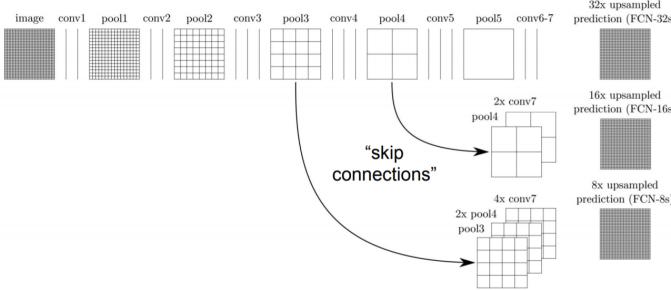


Figure 3: Upsampling using Shortcut Connections between convolution layers

**Pros** End to end fully convolution layer are easier to train because of parameter sharing.

**Cons** It produces coarse segmentation as dimension of image is reduced at each step due to pooling operation.

### 3.2.2 SegNet

A SegNet rather than upsampling image uses encoder and decoder network. The max-pooling indices are copied to decoder network.

The core trainable segmentation engine consists of an encoder network, a corresponding decoder network followed by a pixel-wise classification layer. The role of the decoder network is to map the low resolution encoder feature maps to full input resolution feature maps for pixel-wise classification. The novelty of SegNet lies in the manner in which the decoder upsamples its lower resolution input feature map(s). Specifically, the decoder uses pooling indices computed in the max-pooling step of the corresponding encoder to perform non-linear upsampling. This eliminates the need for learning to upsample. [1]

**Pros** Efficient both in terms of memory and computational time during inference. It is also significantly smaller in the number of trainable parameters than other competing architectures

**Cons** Poor benchmark performance over data-set.

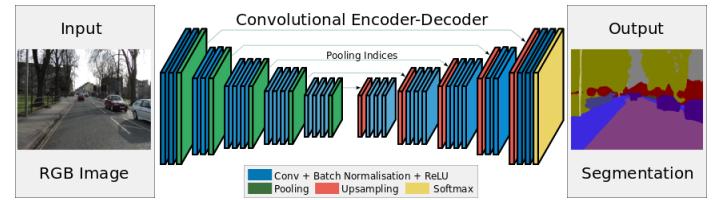


Figure 5: Segmentation Net Architecture

### Accumulated dataset importance

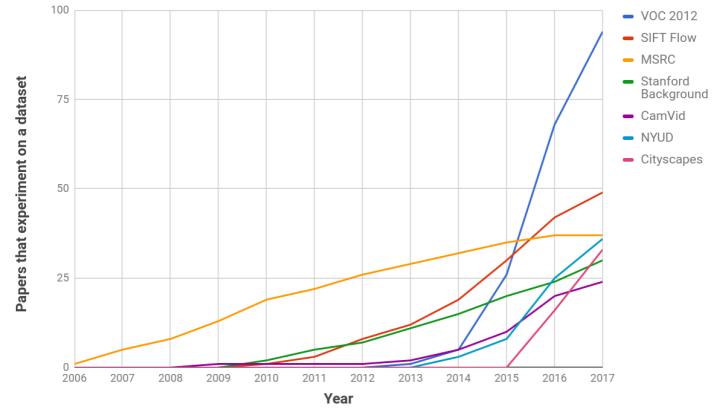


Figure 6: Accumulated data set importance

### 3.3 Dilated Convolutions

A new convolution network module that is specifically designed for dense prediction. The module uses dilated convolutions to systematically aggregate multi-scale contextual information without losing resolution. The architecture is based on the fact that dilated convolutions support exponential expansion of the receptive field without loss of resolution or coverage. [9]

### 3.4 Data Set

The computer vision community produced a couple of different datasets which are publicly available.

The most important datasets are VOC2012 and MSCOCO (MSRC) which are large-scale object detection, segmentation, and captioning dataset.

### 3.5 Experimental Protocol

A typical segmentation pipeline gets raw pixel data, applies preprocessing techniques like scaling. For training, data augmentation techniques such as image rotation can be applied. For every single image, patches of the image called windows are extracted and those windows are classified. The resulting semantic segmentation can be refined by simple morphologic operations. The output is compared and benchmarked. [7]

### 3.6 Measuring the size of objects in an image

There is a whole industry that is working particularly on this field known as Photogrammetry. Photogrammetry has been defined by the American Society for Photogrammetry and Remote Sensing (ASPRS) as the art, science, and technology of obtaining reliable information about physical objects and the environment through processes of recording, measur-

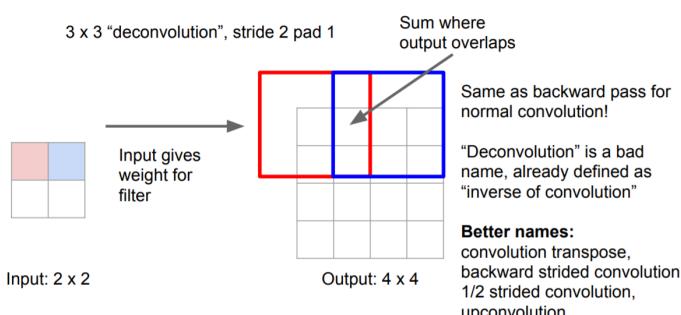


Figure 4: Upsampling using "Deconvolution" layers

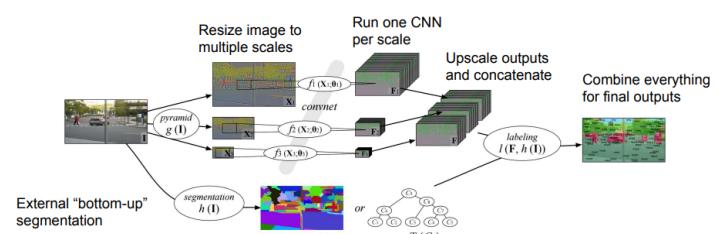


Figure 7: Multi-Scale Context Aggregation by Dilated Convolutions

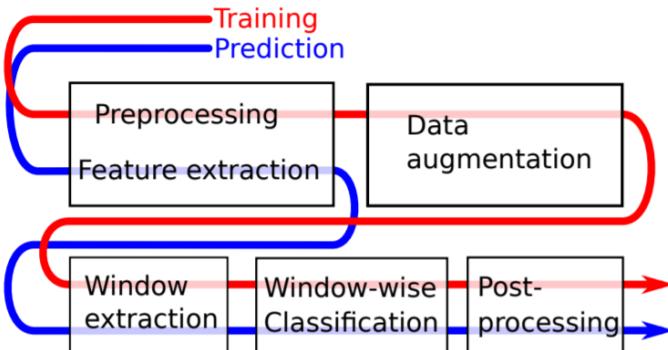


Figure 8: Experimental Protocol



Figure 9: Dimension computation of other objects from a reference object

ing and interpreting photographic images and patterns of recorded radiant electromagnetic energy and other phenomena. A special case, called stereophotogrammetry, involves estimating the three-dimensional coordinates of points on an object employing measurements made in two or more photographic images taken from different positions (see stereoscopy). Common points are identified on each image. A line of sight (or ray) can be constructed from the camera location to the point on the object. It is the subintersection of these rays (triangulation) that determines the three-dimensional location of the point. [8]

There exists a simple approach called "pixels per metric" ratio. This method requires a reference object that is quite distinctive in the segmented image and this helps us to extract pixels per metric ratio for a particular image. The reference object should also have a standard size for example a coin. After we have this ratio all the other dimension of the image is computed from the ratio.

It has been claimed that photogrammetry systems are able to measure smooth three-dimensional objects with surface height deviations less than  $1\mu m$ . [6]

### 3.7 Eariler Proposed Pipeline

The idea is to make an app based model. The user would take an snapshot of the object with their smart phone. The user needs to take the snapshot with a reference object kept beside the target object. The working model then runs an image segmentation algorithm using convolution neural network with learnable upsampling. The architecture of the same would resemble FCNs and based on GoogleNet and AlexNet. On this segmented image we apply dimension determination using "pixels per metric" ratio. The rest of the dimension computed from this ratio.

## 4 Method

1. Foreground extraction of the desired Object
2. Get the image to pixel ratio

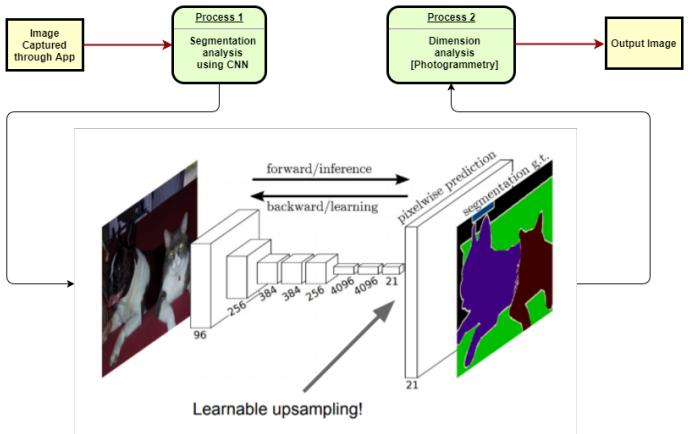


Figure 10: Earlier Proposed Pipeline of Model

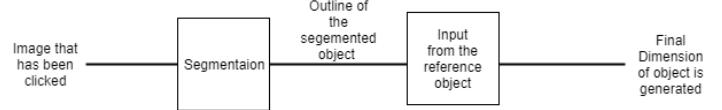


Figure 11: Overall Pipeline

3. Estimation of dimension the perimeter (boundary) and bounding box.

### Listing 1: Workflow

#### 4.1 Foreground extraction of the desired Object

There are three main segmentation algorithms that have been deployed in the first segment of the code.

**Active contours** to separate the image into foreground and the background. The whole algorithm works in an iterative method. So it was extremely essential to properly tune the number of iterations so that foreground and the background is effectively separated without too much loss in the essential information.

After the separation is done we apply morphological operators on the image. The information that are partially present in only the boundaries are removed. Eroding has been done in order to extract only the desired information. Till some extent dilation(hole filling) had also been used to get a proper segmented image. We then extract the boundaries of this foreground object.

**Lazy Snapping** an interactive image cutout tool. Lazy Snapping separates coarse and fine scale processing, making object specification and detailed adjustment easy. It works on max-flow min-cut algorithm to minimize energy based on graph cut [2].

Was used in *LAB color space* as an region of interest(ROI) based model. This also tries to segment the image into foreground and the background. But in this the user needs to select some portion of what he considers like the foreground and what portion it considers as the background. Based on this selection the code separates out the foreground part of the image. This foreground image is then converted to binary image. In the last part we extract the boundary of the segmented image.

**GrabCut** is an image segmentation method based on graph cuts. It works both on the ROI model and the foreground and the background selection. [5] The algorithm estimates the color distribution of the target object and that of the background using a Gaussian mixture model. This is used to construct a Markov random field over the pixel labels, with an energy function that prefers connected regions having the same

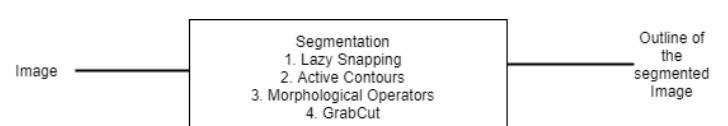


Figure 12: Segmentation Pipeline

label, and running a graph cut based optimization to infer their values. As this estimate is likely to be more accurate than the original, taken from the bounding box, this two-step procedure is repeated until convergence. Estimates can be further corrected by the user by pointing out misclassified regions and rerunning the optimization. The method also corrects the results to preserve edges.

**Morphological operators** using sobel edge detector and simple morphological operators to segment the image into foreground and the background. The algorithm detects edge over a certain threshold value. The next we implement dilation and hole filling in on this edge map. This gives us complete object such that there are no holes or void left in the image. Next we apply erosion to smoothen the boundaries of the foreground. On this segmented foreground we extract the boundary.

## 4.2 Get the image to pixel ratio

To make implementation generic the user is asked choose the dimension of any object according to their convenience of which they know the real world dimension.

A scaling factor is then computed based on

$$\text{factor} = \frac{\text{Actual Distance}}{\text{Euclidean Distance in image}}$$

This factor can be multiplied by the euclidean distance in the image to compute the real world dimension of the objects.

## 4.3 Estimation of dimension the boundary and bounding box

A close bounding box is fitted around the segmented object. The length and width of is obtained by multiply by scaling factor.

The perimeter of the object that has been segmented is also estimated. A Depth First search is performed from starting a source and each neighbour is added with euclidean distance with previous neighbours; to get cumulative sum of distances between adjacent points.

This data is helpful in for calculating the length and width of the packing that might be required for such an object.

## 4.4 Results

The results are compared using the methods for active contour, lazy snapping, grab cut, morphological operators.

Results using Lazy Snapping

Image	Width	Height	Boundary	Time
Shoe	7.785	12.947 cm	38.931 cm	0.132 sec
Phone 1	16.426 cm	9.230 cm	59.559 cm	5.573 sec
Phone 2	17.069 cm	13.036 cm	56.122 cm	5.606 sec
Mouse	10.073 cm	6.352 cm	35.097 cm	5.592 sec

Results using Active Contour

Image	Width	Height	Boundary	Time
Shoe	8.409	14.137 cm	50.36 cm	5.512 sec
Phone 1	15.634 cm	8.626 cm	51.111 cm	54.756 sec
Phone 2	17.306 cm	13.405 cm	51.955 cm	53.003 sec
Mouse	10.001 cm	6.211 cm	33.144 cm	45.268 sec

Results using Morphological

Image	Width	Height	Boundary	Time
Shoe	12.039	14.760 cm	47.476 cm	0.027 sec

Results using inbuilt Matlab

Image	Width	Height	Boundary
Pen	15.076	4.830 cm	40.191 sec

## 4.5 Observations

**Observation** The lazy snapping is most generic with decent output; and low running time.



Figure 13: Shoe Lazy Snapping



Figure 14: Shoe - Segment - Lazy Snapping



Figure 15: Phone 1 Lazy Snapping

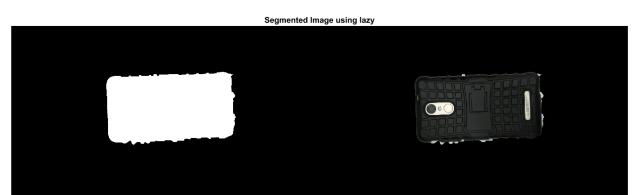


Figure 16: Phone 1 - Segment - Lazy Snapping



Figure 17: Phone 2 Lazy Snapping

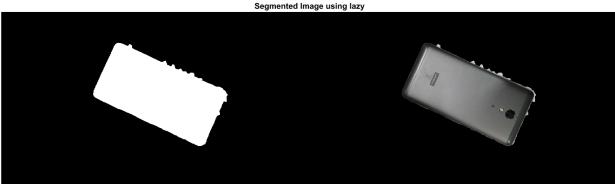


Figure 18: Phone 2 - Segment - Lazy Snapping



Figure 19: Mouse Active Contour



Figure 20: Phone 1 Active Contour

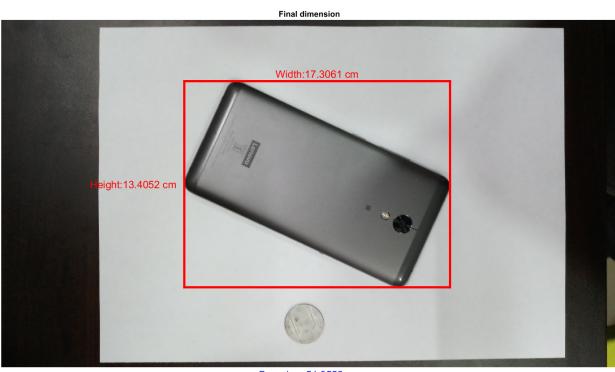


Figure 21: Phone 2 Active Contour



Figure 22: Shoe Morphological

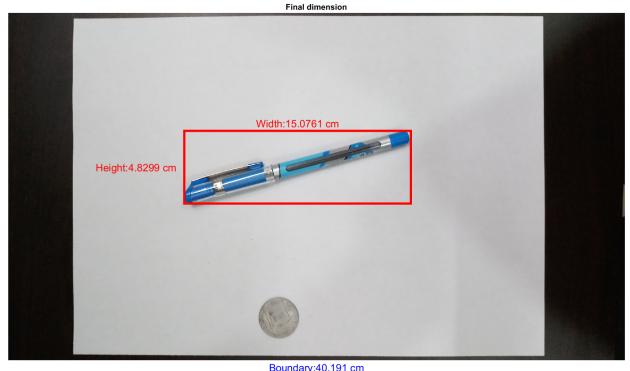


Figure 23: Pen Inbuilt Matlab

**Explanation** Instant results is made possible by a novel image segmentation algorithm which combines graph cut with pre-computed oversegmentation.

**Observation** Morphological are least generic for segmentation

**Explanation** Morphological operators have very large number of parameter like size of kernel, kernel itself and hence a single model cannot work for all instances.

**Observation** Active contour takes highest running time and gives best segmentation

**Explanation** Active contour has large number of iterations and removes background pixels to maximal extend.

#### 4.6 Future Works

The model has variety of areas for improvement

1. The model currently neglects the depth information; which can lead large error since the scaling factor varies across the depth. To accommodate the same; a depth map needs to be obtained using the multi-view geometry using 2 or more images.
2. The reference object can be detected using a deep neural network, cascade object detection or any other method to minimize requirements of inputs by user.

#### References

- [1] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *arXiv preprint arXiv:1511.00561*, 2015.
- [2] Vivek Kwatra, Arno Schödl, Irfan Essa, Greg Turk, and Aaron Bobick. Graphcut textures: image and video synthesis using graph cuts. In *ACM Transactions on Graphics (ToG)*, volume 22, pages 277–286. ACM, 2003.
- [3] Yin Li, Jian Sun, Chi-Keung Tang, and Heung-Yeung Shum. Lazy snapping. In *ACM Transactions on Graphics (ToG)*, volume 23, pages 303–308. ACM, 2004.
- [4] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015.
- [5] Carsten Rother, Vladimir Kolmogorov, and Andrew Blake. Grabcut: Interactive foreground extraction using iterated graph cuts. In *ACM transactions on graphics (TOG)*, volume 23, pages 309–314. ACM, 2004.
- [6] Danny Sims-Waterhouse, Samanta Piano, and Richard Leach. Verification of micro-scale photogrammetry for smooth three-dimensional object measurement. *Measurement Science and Technology*, 28(5): 055010, 2017.
- [7] Martin Thoma. A survey of semantic segmentation. *arXiv preprint arXiv:1602.06541*, 2016.

- [8] Wikipedia. Photogrammetry — Wikipedia, the free encyclopedia. <http://en.wikipedia.org/w/index.php?title=Photogrammetry&oldid=796993379>, 2017. [Online; accessed 24-September-2017].
- [9] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015.