

# DL Course Project

*Joint Depth Estimation and Object Detection*

MOHAMMAD AMIN ALAMALHODA  
AMIRREZA HATAMIPOUR  
MOHAMMADREZA ALIMOHAMMADI

# Contents

|       |  |    |
|-------|--|----|
| 1     | Git and Project Dependencies . . . . .   | 1  |
| 1.1   | Git . . . . .                            | 1  |
| 1.2   | Project Dependencies . . . . .           | 1  |
| 2     | Project Structure . . . . .              | 2  |
| 2.1   | train.ipynb . . . . .                    | 2  |
| 2.2   | coco_eval.py and coco_utils.py . . . . . | 2  |
| 2.3   | datas.py . . . . .                       | 2  |
| 2.4   | engine.py . . . . .                      | 2  |
| 2.5   | models.py . . . . .                      | 2  |
| 2.6   | transforms.py . . . . .                  | 2  |
| 2.7   | utils.py . . . . .                       | 2  |
| 2.8   | GUI Directory . . . . .                  | 2  |
| 3     | Datas . . . . .                          | 3  |
| 3.1   | Converting Datas . . . . .               | 3  |
| 3.2   | Loading Datas . . . . .                  | 3  |
| 3.3   | Data Augmentation . . . . .              | 4  |
| 3.4   | Data Loader . . . . .                    | 4  |
| 3.5   | Loading .mat Dataset . . . . .           | 4  |
| 4     | Models . . . . .                         | 5  |
| 4.1   | Object Detector . . . . .                | 5  |
| 4.1.1 | Architecture . . . . .                   | 5  |
| 4.1.2 | Concept . . . . .                        | 5  |
| 4.1.3 | Computation Pricedure . . . . .          | 6  |
| 4.1.4 | Fine-Tuning . . . . .                    | 7  |
| 4.1.5 | Evaluation . . . . .                     | 7  |
| 4.1.6 | Test Images . . . . .                    | 8  |
| 4.2   | Depth Estimator . . . . .                | 9  |
| 4.2.1 | Architecture . . . . .                   | 9  |
| 4.2.2 | Computation Pricedure . . . . .          | 9  |
| 4.2.3 | Fine-Tuning . . . . .                    | 10 |
| 4.2.4 | Evaluation . . . . .                     | 10 |
| 4.2.5 | Test Images . . . . .                    | 11 |
| 5     | Contacenator . . . . .                   | 12 |
| 5.1   | First Method . . . . .                   | 12 |
| 5.1.1 | Evaluation . . . . .                     | 12 |
| 5.2   | Second Method . . . . .                  | 13 |
| 5.2.1 | Evaluation . . . . .                     | 13 |

# List of Figures

|    |  |    |
|----|--|----|
| 1  | Some Sample Images from The Dataset . . . . .                    | 3  |
| 2  | An Image and its Augmented Version . . . . .                     | 4  |
| 3  | Bounding Box Regression in the Mask-RCNN . . . . .               | 5  |
| 4  | A Schematic of the Mask-RCNN . . . . .                           | 6  |
| 5  | Training Loss of the Object Detector Model . . . . .             | 7  |
| 6  | Validation Loss of the Object Detector Model . . . . .           | 7  |
| 7  | Output of Object Detector Network for Some Test Images . . . . . | 8  |
| 8  | Training Loss of the Depth Estimator Model . . . . .             | 10 |
| 9  | Validation Loss of the Depth Estimator Model . . . . .           | 10 |
| 10 | Output of Depth Estimator Network for Some Test Images . . . . . | 11 |
| 11 | Schematic of the AE used for Concatenating . . . . .             | 12 |
| 12 | Perplexity of Training and Validation . . . . .                  | 12 |
| 13 | BELU Metrics . . . . .   | 13 |

# List of Tables

|   |                          |    |
|---|--------------------------|----|
| 1 | Object Detection Metrics | 7  |
| 2 | Depth Estimator Metrics  | 10 |





# 1 Git and Project Dependencies

## 1.1 Git

This Project is open source and is published on Github. You can watch it using [this link](#).

You can use the following bash command for cloning this project:

```
$ git clone https://github.com/MohammadAminAlamhoda/Deep-Object
```

If you don't have `git` installed on your device, you can use the following bash command:

- Linux

```
$ sudo apt-get install git
```

- MacOS

MacOS already have `git` installed, check its version using bash command below:

```
$ git --version
```

If you uninstalled it, you can install it using `brew`:

```
$ brew install git
```

- Windows

You can download source code of `git` and makeinstall it using [this link](#).

## 1.2 Project Dependencies

This project needs the following stuff in order to be compiled successfully.

- PyTorch - Python Lib
- Torchvision - Python Lib
- OpenCV - Python Lib
- Matplotlib - Python Lib
- mat73 - Python Lib
- Os - Python Lib
- Sys - Python Lib
- **Mask-RCNN**- Torch Hub Model
- **MiDaS** - Torch Hub Model



## 2 Project Structure

This project contains different files. We will explain them in the following.

### 2.1 train.ipynb

This jupyter notebook file contains the configs and essential properties for training the networks.

### 2.2 coco\_eval.py and coco\_utils.py

This files contain coco tools for image cropping and preparing images for object detection.

### 2.3 datas.py

This files contains dataloader classes and handles the loading and augmentation of the datas during training.

### 2.4 engine.py

This file contains the functions for training and evaluation of the networks.

### 2.5 models.py

All the model classes (Object Detector, Depth Estimator, and Concatener) are defined in this file.

### 2.6 transforms.py

Torchvision transfromation for augmenting the data set are defined in this file.

### 2.7 utils.py

This file contain some useful function such as model saver-loader, loggers, and ... .

### 2.8 GUI Directory

This directory contains files related to our application for loading and joint object detection and depth estimation of the datas.



### 3 Datas

#### 3.1 Converting Datas

We converted the datas which were in the .mat format to the .png for better RAM management. This way it is possible to load the images directly from the hard drive in each iteration.

#### 3.2 Loading Datas

Figure 1 shows a sample from the dataset.

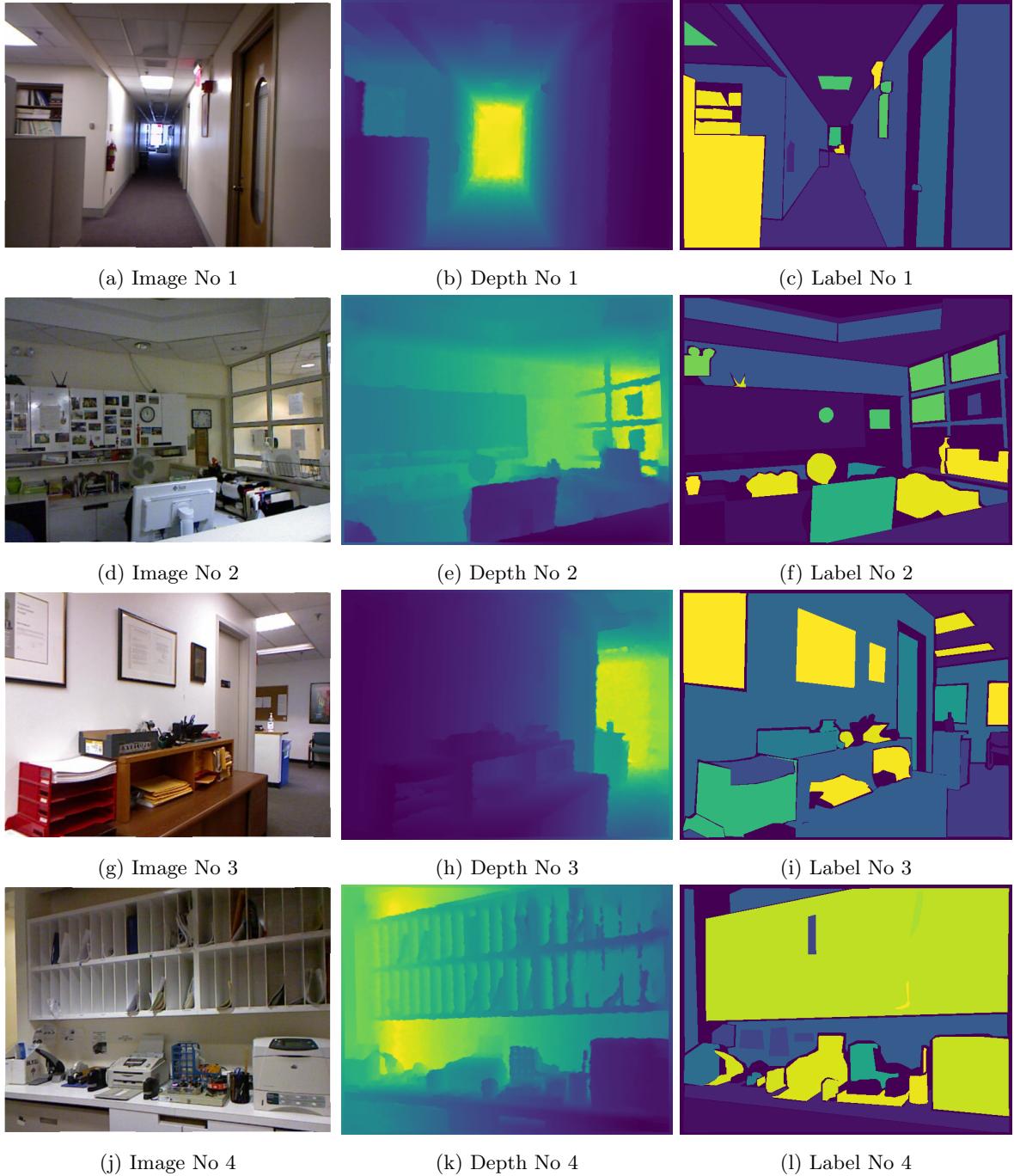


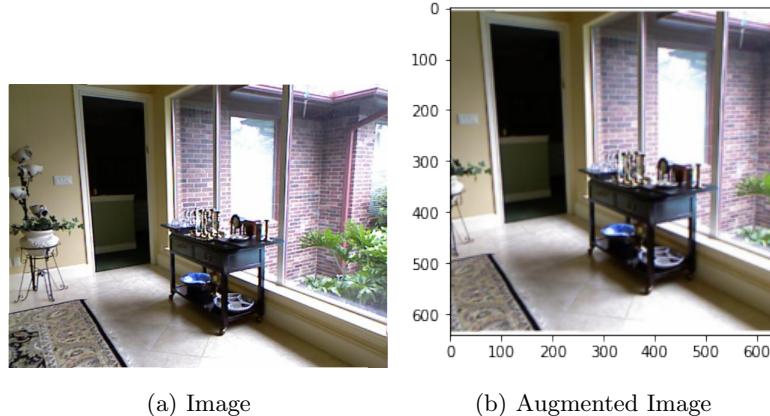
Figure 1: Some Sample Images from The Dataset

As can be seen in Figure 1, dataset contains image, label, and depth.



### 3.3 Data Augmentation

We augmented the datas by resizing to 640, random cropping an  $640 \times 640$  square, and random horizontal flip. Some of the Augmented Images are plotted in Figure 2. It is noteworthy to mention that all the datas(images, labels, and depths) which have belong to a specific scene, should be augmented by same properties.



(a) Image

(b) Augmented Image

Figure 2: An Image and its Augmented Version

Transforms are done using `torchvision.datasets.transforms`.

### 3.4 Data Loader

We used `torch data loaders` for managing datas and augmentation. You can find this files in `datas.py` file.

### 3.5 Loading .mat Dataset

Due to RAM shortage, we didn't load the datas on the RAM. All the images, labels and depths in .mat file were converted to .png images and the loaded directly from the hard drive during training. This was done using function `mat2png` which can be find in `utils.py`



## 4 Models

### 4.1 Object Detector

We used pretrained *Mask-RCNN* as the object detector network and fine tuned it using our dataset. You can find the paper of *Mask-RCNN* using [this link](#).

#### 4.1.1 Architecture

This network develops, for instance segmentation. Instance segmentation is challenging because it requires the correct detection of all objects in an image while precisely segmenting each instance. It therefore combines elements from the classical computer vision tasks of object detection, where the goal is to classify individual objects and localize each using a bounding box, and semantic segmentation, where the goal is to classify each pixel into a fixed set of categories without differentiating object instances. Mask R-CNN, extends Faster R-CNN by adding a branch for predicting segmentation masks on each Region of Interest (RoI), in parallel with the existing branch for classification and bounding box regression. You can see an example of box regression in the Figure 3.

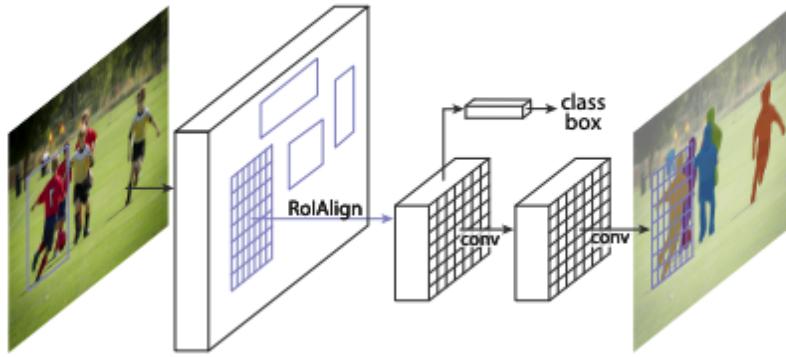


Figure 3: Bounding Box Regression in the Mask-RCNN

The mask branch is a small FCN applied to each RoI, predicting a segmentation mask in a pixel-to-pixel manner. Mask R-CNN is simple to implement and train given the Faster R-CNN framework, which facilitates a wide range of flexible architecture designs. Additionally, the mask branch only adds a small computational overhead, enabling a fast system and rapid experimentation. Faster RCNN was not designed for pixel-to-pixel alignment between network inputs and output. To fix the misalignment, proposed a simple, quantization-free layer, called RoIAlign, that faithfully preserves exact spatial locations. Despite being a seemingly minor change, RoIAlign has a large impact: it improves mask accuracy by relative 10% to 50%, showing bigger gains under stricter localization metrics. Second, we found it essential to decouple mask and class prediction: we predict a binary mask for each class independently, without competition among classes, and rely on the network's RoI classification branch to predict the category.

#### 4.1.2 Concept

Mask R-CNN is conceptually simple: Faster R-CNN has two outputs for each candidate object, a class label, and a bounding-box offset; to this, we add a third branch that outputs the object mask. Mask R-CNN is thus a natural and intuitive idea. But the additional mask output is distinct from the class and box outputs, requiring extraction of an object's much finer spatial layout. Next, we introduce the key elements of Mask R-CNN, including pixel-to-pixel alignment, which is the main missing piece of Fast/Faster R-CNN.



#### 4.1.3 Computation Procedure

Mask R-CNN adopts the same two-stage procedure, with an identical first stage (RPN). In parallel to predicting the class and box offset in the second stage, Mask R-CNN also outputs a binary mask for each ROI. This contrasts with most recent systems, where classification depends on mask predictions. Our approach follows the spirit of Fast R-CNN that applies bounding-box classification and regression in parallel (which turned out to largely simplify the multi-stage pipeline of original R-CNN). Formally, during training, we define a multi-task loss on each sampled ROI as:

$$L = L_{cs} + L_{box} + L_{mask}$$

The classification loss  $L_{cs}$  and bounding-box  $L_{bo}$ . The mask branch has a  $Km^2$ -dimensional output for each ROI, which encodes  $K$  binary masks of resolution  $m \times m$ , one for each of the  $K$  classes. To this, apply a per-pixel sigmoid and define  $L_{mask}$  as the average binary cross-entropy loss, for an ROI associated with ground-truth class  $k$ ,  $L_{mask}$  is only defined on the  $k^{th}$  mask (other mask outputs do not contribute to the loss).

$L_{mask}$  allows the network to generate masks for every class without competition among classes; They relied on the dedicated classification branch to predict the class label used to select the output mask. This decouples mask and class prediction. This is different from common practice when applying FCNs to semantic segmentation, which typically uses a per-pixel softmax and a multinomial cross-entropy loss. In that case, masks across classes compete; they do not have a per-pixel sigmoid and a binary loss.

You can see an overview of *Mask-RCNN* in the Figure 4

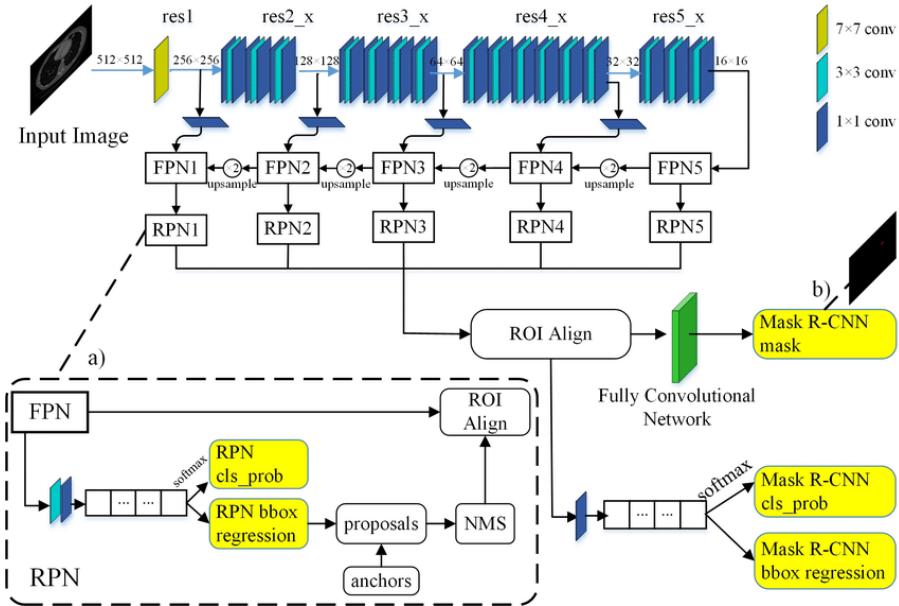


Figure 4: A Schematic of the Mask-RCNN



#### 4.1.4 Fine-Tuning

As mentioned earlier, we fine tuned the model by using our datas. You can see the training loss plot in Figure 5.

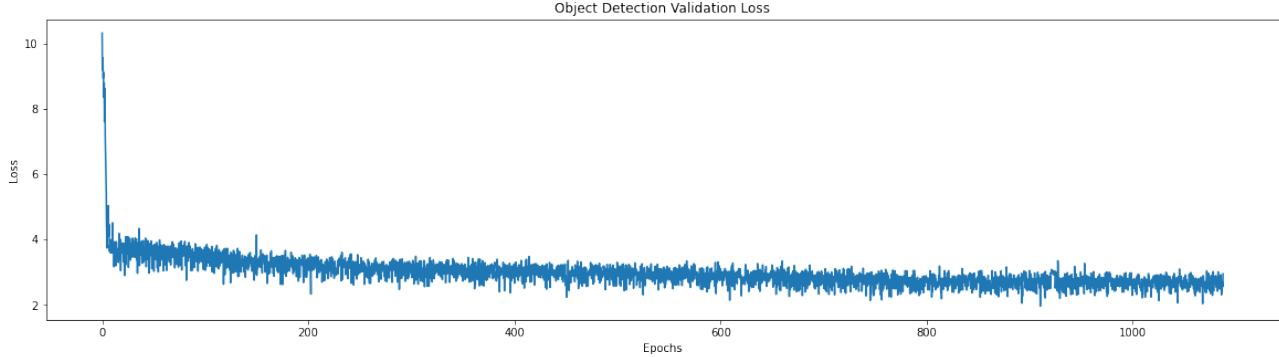


Figure 5: Training Loss of the Object Detector Model

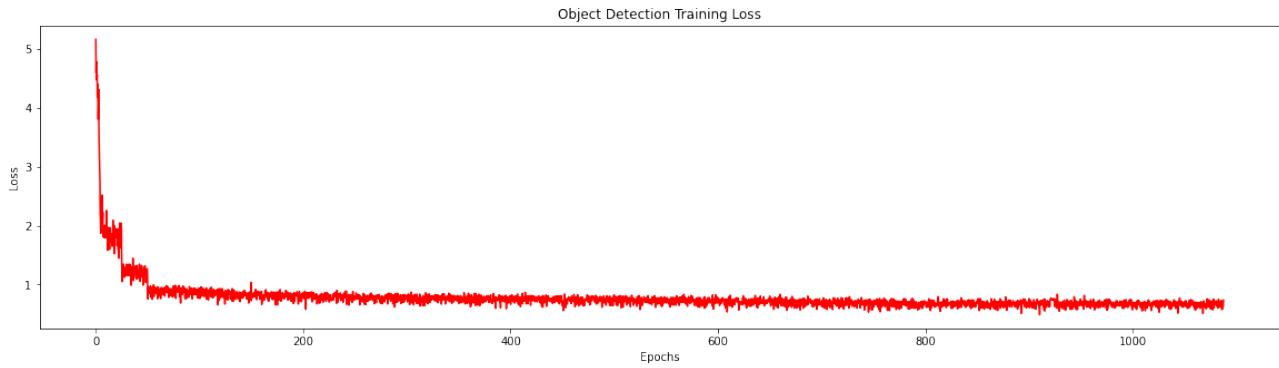


Figure 6: Validation Loss of the Object Detector Model

#### 4.1.5 Evaluation

Evaluating object detection models is not straightforward because each image can have many objects and each object can belong to different classes. This means that we need to measure if the model found all the objects and also a way to verify if the found objects belong to the correct class. This means that an object detection model needs to accomplish two things

- Find all the objects in an image
- Check if the found objects belong to the correct class

Therefore, we used different metrics for evaluating the object detector model which are written in Table 1:

Table 1: Object Detection Metrics

| Metric          | Initial Value-Training | Final Value-Training | Initial Value-Validation | Final Value-Validation |
|-----------------|------------------------|----------------------|--------------------------|------------------------|
| Total loss      | 8.892                  | 2.7060               | 4.323                    | 1.413                  |
| Classifier Loss | 3.902                  | 1.4148               | 1.032                    | 0.242                  |
| Box Reg Loss    | 2.023                  | 0.6253               | 1.231                    | 0.321                  |
| Mask Loss       | 1.834                  | 0.4961               | 0.891                    | 0.214                  |
| Objectness Loss | 0.673                  | 0.0435               | 0.413                    | 0.001                  |



#### 4.1.6 Test Images

You can see some test images in the Figure 7.

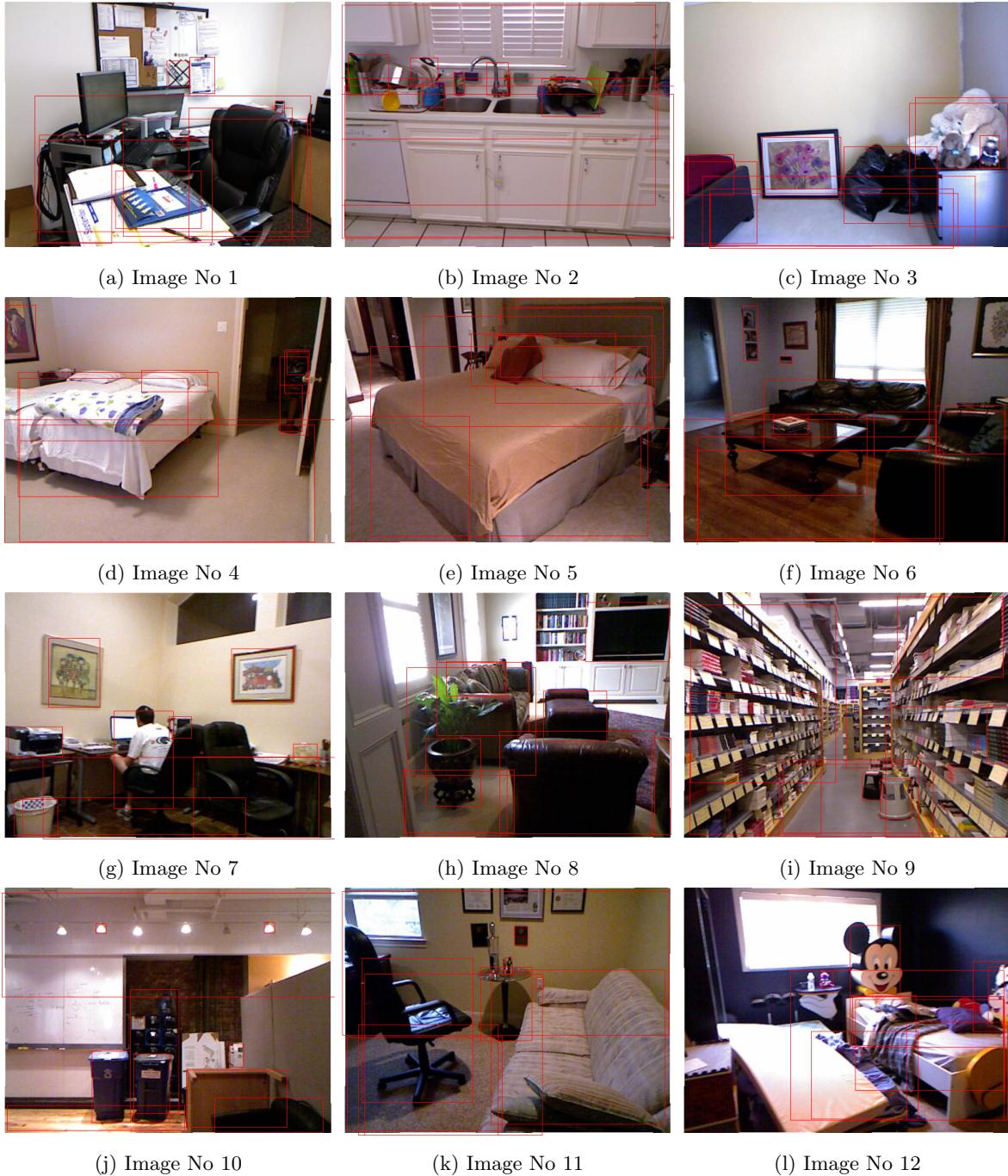


Figure 7: Output of Object Detector Network for Some Test Images



## 4.2 Depth Estimator

We used pretrained *Intelisl MiDaS* as the depth estimation network and fine tuned it using our dataset. You can find the paper of *MiDaS* using [this link](#).

### 4.2.1 Architecture

### 4.2.2 Computation Pricedure



#### 4.2.3 Fine-Tuning

As mentioned earlier, we fine tuned the model by using our datas. You can see the training loss plot in Figure 8.

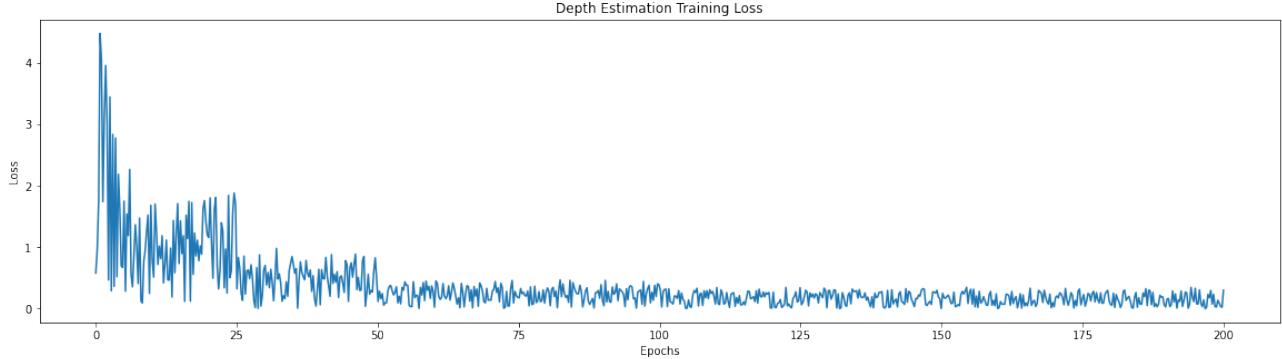


Figure 8: Training Loss of the Depth Estimator Model

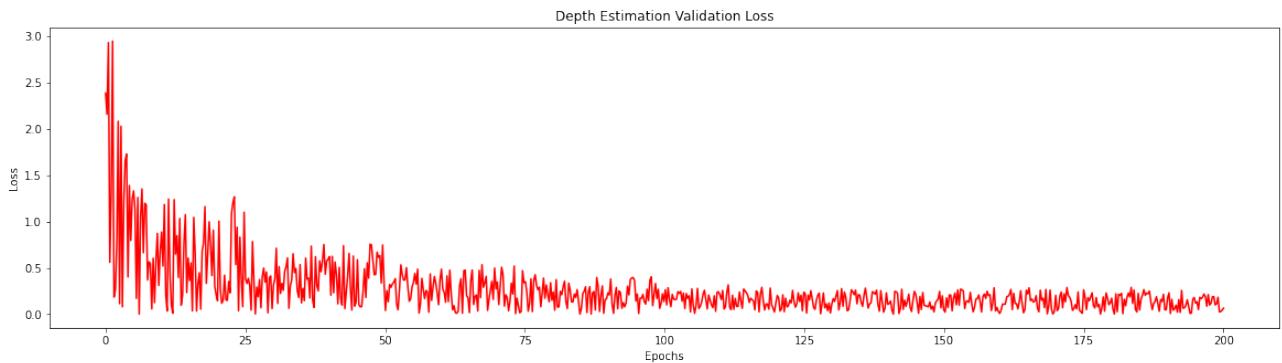


Figure 9: Validation Loss of the Depth Estimator Model

#### 4.2.4 Evaluation

Estimating depth of different pixels in an image is a regression problem, so we used MSE(mean squad error) as loss. You can see this loss values for training and validation in the Table 2.

Table 2: Depth Estimator Metrics

| Metric   | Initial Value-Training | Final Value-Training | Initial Value-Validation | Final Value-Validation |
|----------|------------------------|----------------------|--------------------------|------------------------|
| MSE loss | 4.592                  | 0.2060               | 6.234                    | 0.124                  |

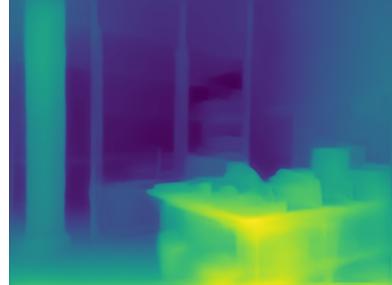


#### 4.2.5 Test Images

You can see some test images in the Figure 10.



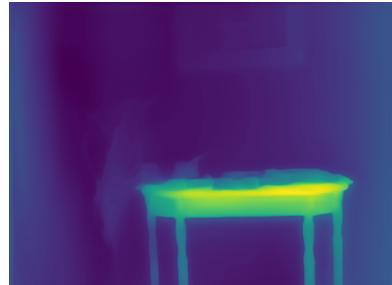
(a) Image No 1



(b) Depth of Image No 1



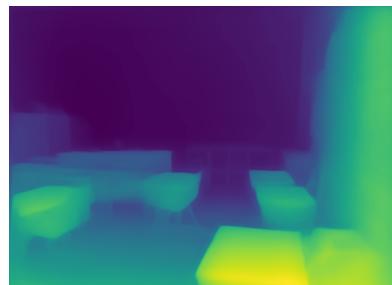
(c) Image No 2



(d) Depth of Image No 2



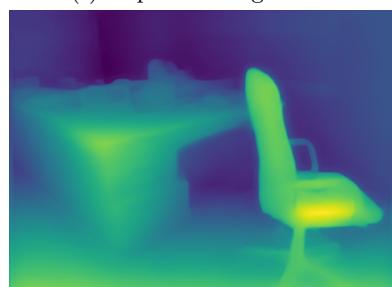
(e) Image No 3



(f) Depth of Image No 3



(g) Image No 4



(h) Depth of Image No 4

Figure 10: Output of Depth Estimator Network for Some Test Images



## 5 Contacenator

We used two different methods for concatenating the object detector and depth estimator:

### 5.1 First Method

In the first method, we used an auto-encoder network which was placed after the depth estimator and before the object detector. This AE's input was the estimated depth and was the detected objects and masks. This network was learned when both the other networks were frozen. You can see a schematic of this AE in Figure 11.

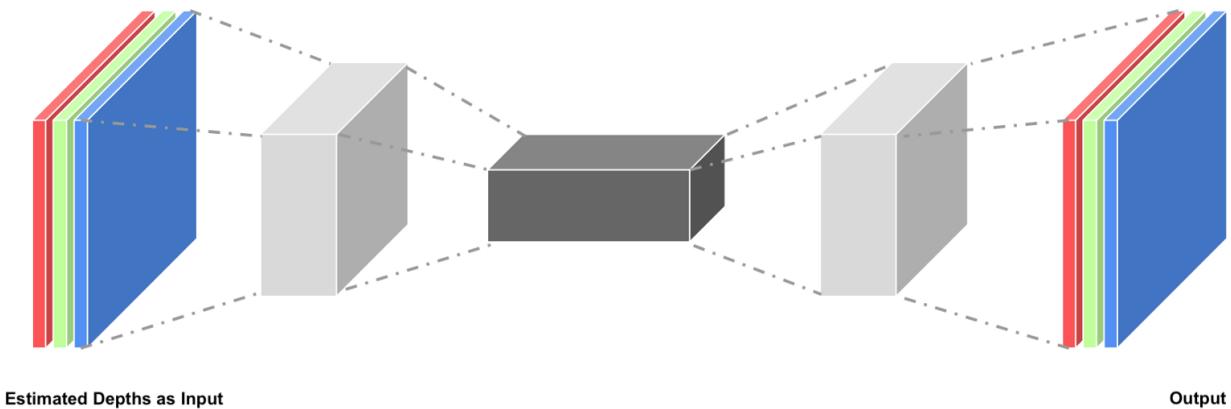


Figure 11: Schematic of the AE used for Concatenating

This method can perform well if we have large dataset and very fast computing hardware, but because of the small dataset and slow computing hardware, we couldn't fully train this model.

#### 5.1.1 Evaluation

We used BLEU and perplexity metrics for evaluating functionality of the concatenator. You can see the perplexity of the training and validation in the Figure 12 and definition of the perplexity of a probability distribution in 1.

$$PP(p) = 2^{H(p)} = 2^{-\sum p(x)} \quad (1)$$

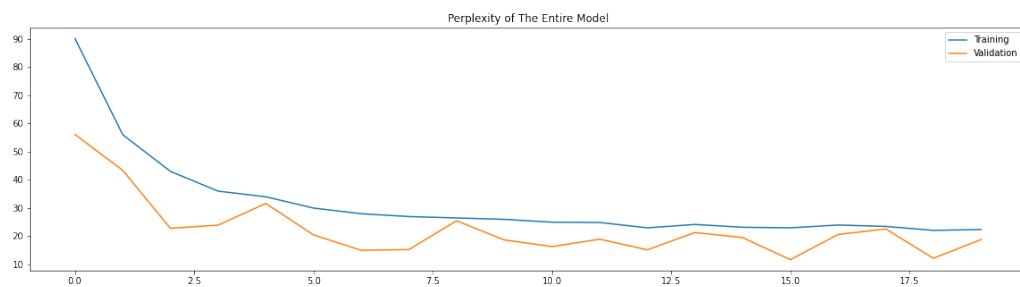
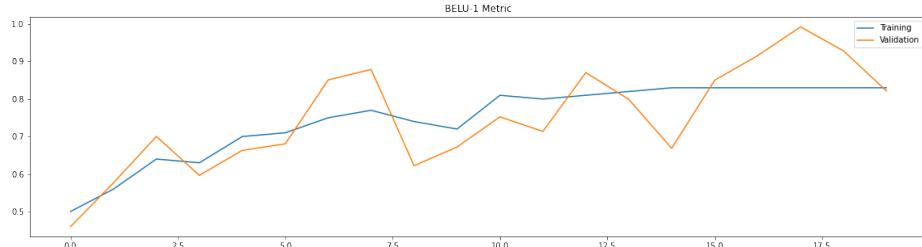


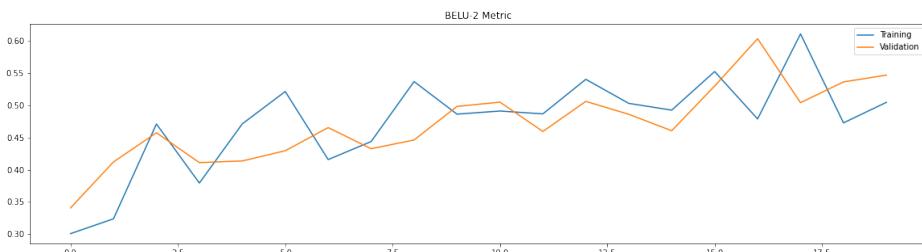
Figure 12: Perplexity of Training and Validation



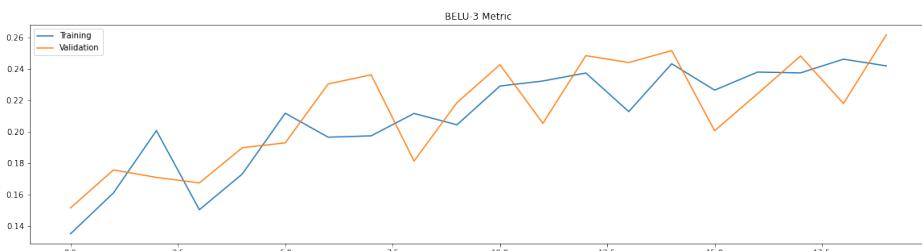
You can see BELU-1, 2, 3, and 4 metrics in Figure 13.



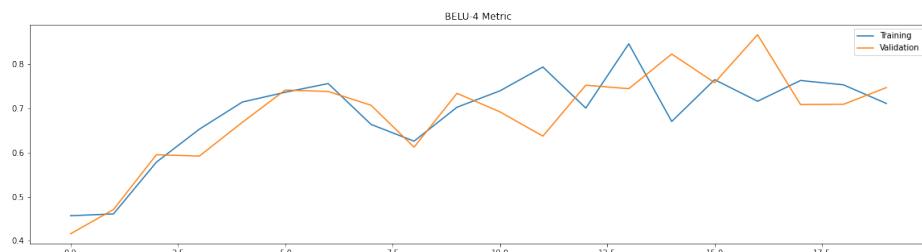
(a) BELU-1 Metric



(b) BELU-2 Metric



(c) BELU-3 Metric



(d) BELU-4 Metric

Figure 13: BELU Metrics

## 5.2 Second Method

In this method, there is just object detector and depth estimator. by thresholding on the depth estimator output, we split the detected objects to thee classes: Close, Middle, and Far.

### 5.2.1 Evaluation

Because in this method we perform concatenation by thresholding on the depths, there is no numerical metric for measuring performance, but you can see the final outputs of the network in the next pages.