

# MidTerm Report: Bully Images Recognition

Haotian Deng and Qingbo Lai

<sup>1</sup> School of Computing, Clemson University

<sup>2</sup> [hdeng@clemson.edu](mailto:hdeng@clemson.edu) and [qingbol@clemson.edu](mailto:qingbol@clemson.edu)

## 1 Introduction

Detecting In the mid-term project, we classified the bully pictures and nonbully pictures which are totally 10 categories through the modified VGG16 model[4], and subdivided into 9 categories in the bully picture. The project achieved a satisfactory accuracy. In the final project we will inherit the success of the mid-term project and implement object detection on the basis of the mid-term project. In the final project, we implement object detection based on Single Shot MultiBox Detector (SSD)[2].

In the report, section 2 talks about our project setting and models, section 3 section 4 illustrate our experiments, section 5 describe the future work and plans of final project, section 6 is a conclusion.

## 2 Project Description

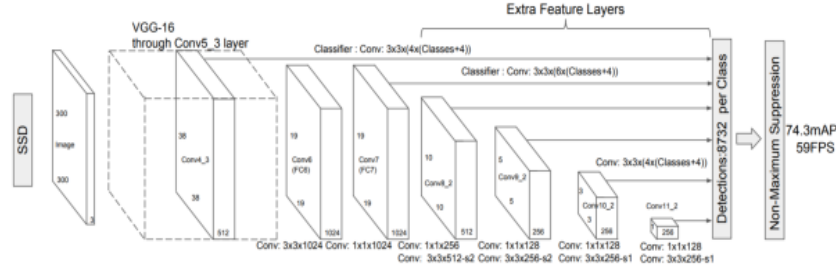
The project's goal is develop deep learning models to detect images with bullying actionsModel able to identify perpetrators which do bully action and victims. In our whole project design, We first analyzed the well-known object detection model and found that the Single Shot MultiBox Detector(SSD)[2] model is the most suitable for us. we use the VGG16 model for image classification in our mid-term project. We know more about the VGG16 model and SSD is based on a modified VGG- 16 network. In the project, We don't have a ready data set for this projectthe data set provided by the CPSC 8810 class which do not have the information and labels about objects like bully and victim, so preprocessing data and make own dataset is also a important procedure. SSD model combines the regression idea in YOLO [5] and the Anchor mechanism in Faster-RCNN, and uses the multi-scale regional features of each position of the whole graph to perform regression, which not only maintains the fast YOLO characteristics, but also guarantees the window prediction and the Faster-RCNN.

The core of the SSD is to use a convolution kernel on the feature map to predict the category and coordinate offset of a series of Default Bounding Boxes. In order to improve the detection accuracy, SSD predicts on different scales of feature maps. More details about model like the specific accuracy and performance, parameters in Section 3.

## 2.1 Model Architecture

The idea of SSD network main body design is feature layer extraction, and the frame regression and classification are performed in turn. Because different levels of feature maps can represent different levels of semantic information, low-level feature maps can represent low-level semantic information (with more details), can improve the quality of semantic segmentation, and are suitable for small-scale target learning. High-level feature maps can represent high-level semantic information, smooth segmentation results, and are suitable for in-depth learning of large-scale goals. Therefore, the network proposed by the author of SSD can be suitable for target detection at different scales.

The object detection algorithm is divided into one-stage and two-stage. SSD is a one-stage algorithm, which has many advantages. SSD model combines the regression idea in YOLO [5] and the Anchor mechanism in Faster-RCNN, and uses the multi-scale regional features of each position of the whole graph to perform regression, which not only maintains the fast YOLO characteristics, but also guarantees the window prediction and the Faster-RCNN. SSD model use VGG16 network to get features, On the basis of the VGG model, remove the last fully connected layer, add 6 convolution layers (0, 1, 2, 3, 4, 5) and take two convolution layers Conv4-3 and Conv7 in the VGG network. (Fc7), take 4 (2, 3, 4, 5) in the next 6 convolutional layers. If the resolution of the input image is 300X300, then the size of the 6 selected features is (38X38, 19X19, 10X10, 5X5, 3X3, 1X1), compared to the original image, the reduced steps are (8, 16, 32, 64, 128, 256). We draw fixed-scale frames on different feature maps. Small frames on large resolutions are helpful for detecting small targets. Boxes on small resolutions are useful for detecting large targets. Therefore, it is possible to obtain predicted values of multiple scales.



**Fig. 1.** Single Shot MultiBox Detector(SSD) architecture[2]

So the SSD network is divided into 6 stages, each stage can learn a feature map, and then carry out border regression and classification. The SSD network uses the first five-layer convolutional network of VGG16 as the first stage, and

then converts the two fully connected layers of fc6 and fc7 in VGG16 into two convolutional layers, Conv6 and Conv7, as the second and third stages of the network. . Then on this basis, the SSD network continues to add Conv8, Conv9, Conv10 and Conv11 four-layer networks to extract higher-level semantic information. The network structure of the SSD is shown in Figure 3.1 below. In each stage operation, the network contains multiple convolutional layer operations, each of which is essentially a small convolution.

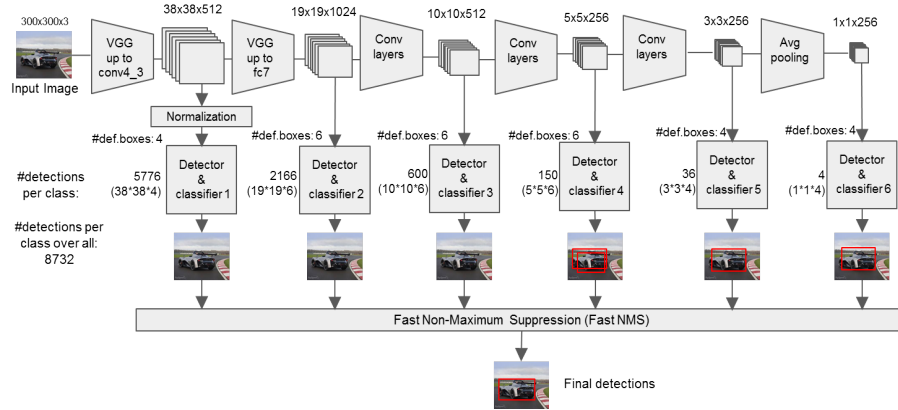


Fig. 2. SSD model architecture[2]

### 3 Dataset

We use the dataset provided in the classroom as the basis, but for the model training, the dataset in the classroom can not meet the requirements, and further processing, such as labeling the object, we use the "labelling" [6] tool to mark the file. we can get "xml" files about images which contains information of labels and regions. And then annotations' formula same as the PASCAL VOC data set. Our data mainly contain two objects, one is bully, the other is victim.

```

1 <annotation>
2   <folder>JPEGImages</folder>
3   <filename>000050.jpg</filename>
4   <path>/home/naruto/Documents/19Spring/Deep Learning/Project/JPEGImages/000050.jpg</path>
5   <source>
6     <database>Unknown</database>
7   </source>
8   <size>
9     <width>360</width>
10    <height>240</height>
11    <depth>3</depth>
12  </size>
13  <segmented>0</segmented>
14  <object>
15    <name>bully</name>
16    <pose>Unspecified</pose>
17    <truncated>0</truncated>
18    <difficult>0</difficult>
19    <bndbox>
20      <xmin>32</xmin>
21      <ymin>26</ymin>
22      <xmax>180</xmax>
23      <ymax>174</ymax>
24    </bndbox>
25  </object>
26  <object>
27    <name>victim</name>
28    <pose>Unspecified</pose>
29    <truncated>0</truncated>
30    <difficult>0</difficult>
31    <bndbox>
32      <xmin>186</xmin>
33      <ymin>69</ymin>
34      <xmax>312</xmax>
35      <ymax>201</ymax>
36    </bndbox>
37  </object>
38 </annotation>
39

```

Fig. 3. Annotations about Images[2]

## 4 Experimental setup

Our enviroments are Python3.6 and Tensorflow basic framework 1.12all the codes implemented on Palmetto Cluster[3] ”<https://palmetto.clemson.edu/palmetto/>” which provided by Clemson University. Our bully picture dataset provided by class and nonbully dataset[1], we download from website. The totally dataset image amount are 2387.

### 4.1 Parameters setting

We test our model so many times and fix the parameters, the following tables are our model parameter setting include Hyperparameters and parameter of networks.

Table 1. Hyperparameters

Parameters	Values
Learning rate	0.001
iteration times	10000 s
batch size	16
Image size	300

## 4.2 Train Process

The image is about train process.

```

Step 901---Loss:[5.3316.25] ** Location:0.966 ** Class:5.29 ** pred_class:[8.89e+04|3.91| 0.0] ** pred_location:[70018.297|5.131|0.000]
Step 902---Loss:[5.3316.36] ** Location:1.07 ** Class:5.29 ** pred_class:[8.89e+04| 3.9| 0.0] ** pred_location:[70036.383|5.105|0.000]
Step 903---Loss:[5.3316.98] ** Location:0.859 ** Class:6.12 ** pred_class:[8.9e+04|3.88| 0.0] ** pred_location:[70047.992|5.135|0.000]
Step 904---Loss:[5.3317.17] ** Location:1.61 ** Class:5.56 ** pred_class:[8.91e+04|3.88| 0.0] ** pred_location:[70072.461|5.103|0.000]
Step 905---Loss:[5.3316.33] ** Location:0.654 ** Class:5.67 ** pred_class:[8.91e+04|3.88| 0.0] ** pred_location:[70078.070|5.102|0.000]
Step 906---Loss:[5.3316.05] ** Location:1.03 ** Class:5.02 ** pred_class:[8.92e+04|3.91| 0.0] ** pred_location:[70101.422|5.092|0.000]
Step 907---Loss:[5.3317.14] ** Location:1.29 ** Class:5.83 ** pred_class:[8.92e+04| 3.9| 0.0] ** pred_location:[70102.195|5.090|0.000]
Step 908---Loss:[5.3316.89] ** Location:0.892 ** Class:6.0 ** pred_class:[8.93e+04| 3.9| 0.0] ** pred_location:[70137.530|5.100|0.000]
Step 909---Loss:[5.3316.53] ** Location:0.802 ** Class:5.73 ** pred_class:[8.93e+04|3.92| 0.0] ** pred_location:[70129.398|5.111|0.000]
Step 910---Loss:[5.3216.46] ** Location:1.6 ** Class:4.86 ** pred_class:[8.94e+04|3.89| 0.0] ** pred_location:[70145.961|5.071|0.000]
Step 911---Loss:[5.3216.34] ** Location:0.802 ** Class:5.54 ** pred_class:[8.94e+04|3.91| 0.0] ** pred_location:[70151.350|5.091|0.000]
Step 912---Loss:[5.3216.23] ** Location:0.91 ** Class:5.32 ** pred_class:[8.95e+04|3.92| 0.0] ** pred_location:[70261.336|5.087|0.000]
Step 913---Loss:[5.3216.54] ** Location:1.03 ** Class:5.51 ** pred_class:[8.96e+04|3.95| 0.0] ** pred_location:[70356.844|5.085|0.000]
Step 914---Loss:[5.321 6.3] ** Location:0.797 ** Class:5.51 ** pred_class:[8.96e+04|3.97| 0.0] ** pred_location:[70410.945|5.093|0.000]
Step 915---Loss:[5.321 6.8] ** Location:0.781 ** Class:6.02 ** pred_class:[8.97e+04|3.98| 0.0] ** pred_location:[70462.531|5.099|0.000]
Step 916---Loss:[5.3216.91] ** Location:1.35 ** Class:5.56 ** pred_class:[8.97e+04|4.02| 0.0] ** pred_location:[70594.391|5.081|0.000]
Step 917---Loss:[5.3216.81] ** Location:1.1 ** Class:5.71 ** pred_class:[8.98e+04|4.05| 0.0] ** pred_location:[70559.633|5.088|0.000]
Step 918---Loss:[5.3216.33] ** Location:0.744 ** Class:5.58 ** pred_class:[8.98e+04|4.05| 0.0] ** pred_location:[70612.961|5.083|0.000]
Step 919---Loss:[5.3216.18] ** Location:0.566 ** Class:5.61 ** pred_class:[8.98e+04|4.06| 0.0] ** pred_location:[70673.055|5.110|0.000]
Step 920---Loss:[5.3216.79] ** Location:1.91 ** Class:5.78 ** pred_class:[8.99e+04|4.05| 0.0] ** pred_location:[70662.242|5.113|0.000]
Step 921---Loss:[5.3216.31] ** Location:0.839 ** Class:5.47 ** pred_class:[8.99e+04|4.08| 0.0] ** pred_location:[70679.727|5.097|0.000]
Step 922---Loss:[5.3216.71] ** Location:1.36 ** Class:5.34 ** pred_class:[9e+04|4.07| 0.0] ** pred_location:[70684.781|5.108|0.000]
Step 923---Loss:[5.3216.75] ** Location:1.41 ** Class:5.33 ** pred_class:[9e+04|4.12| 0.0] ** pred_location:[70731.398|5.112|0.000]
Step 924---Loss:[5.3216.28] ** Location:0.96 ** Class:5.32 ** pred_class:[9.01e+04|4.15| 0.0] ** pred_location:[70739.062|5.114|0.000]
Step 925---Loss:[5.3216.52] ** Location:1.04 ** Class:5.49 ** pred_class:[9.01e+04|4.14| 0.0] ** pred_location:[70718.109|5.104|0.000]
Step 926---Loss:[5.3216.38] ** Location:0.893 ** Class:5.49 ** pred_class:[9.02e+04|4.09| 0.0] ** pred_location:[70710.648|5.099|0.000]
Step 927---Loss:[5.321 6.4] ** Location:0.902 ** Class:5.5 ** pred_class:[9.02e+04|4.05| 0.0] ** pred_location:[70701.664|5.111|0.000]
Step 928---Loss:[5.3216.37] ** Location:0.862 ** Class:5.51 ** pred_class:[9.02e+04|4.04| 0.0] ** pred_location:[70654.383|5.102|0.000]
Step 929---Loss:[5.3216.06] ** Location:0.777 ** Class:5.28 ** pred_class:[9.03e+04|4.02| 0.0] ** pred_location:[70668.977|5.145|0.000]
Step 930---Loss:[5.3216.47] ** Location:0.982 ** Class:5.49 ** pred_class:[9.03e+04|4.07| 0.0] ** pred_location:[70620.070|5.115|0.000]
Step 931---Loss:[5.3216.72] ** Location:1.42 ** Class:5.3 ** pred_class:[9.04e+04|4.04| 0.0] ** pred_location:[70574.531|5.086|0.000]
Step 932---Loss:[5.3217.51] ** Location:1.6 ** Class:5.91 ** pred_class:[9.04e+04|4.06| 0.0] ** pred_location:[70545.617|5.089|0.000]
Step 933---Loss:[5.3216.27] ** Location:0.785 ** Class:5.48 ** pred_class:[9.04e+04|4.05| 0.0] ** pred_location:[70498.530|5.091|0.000]
Step 934---Loss:[5.3216.29] ** Location:1.08 ** Class:5.21 ** pred_class:[9.05e+04|4.08| 0.0] ** pred_location:[70426.891|5.098|0.000]

```

Fig. 4. Process

## 4.3 Results

We design two Api (script) to develop two functions of our models, we test on the dataset which provided by class.

Our model can predict single image or a group images, we can see the ground truth and we are predicted class indicated the object (Arrange by possibility), predicted class valuepossibility and predicted location (region of object).

We can find our models successfully predict the image and it contains objects.

```

[img-25 actual]:[0.2125, 0.5541666666666667, 0.325, 0.8833333333333333, 2]
[img-25 actual]:[0.7652777777777777, 0.5229166666666667, 0.4694444444444444,
0.8458333333333333, 2]
predicted class is [2, 2, 1, 2]
predicted class value:[0.99591124, 0.99006337, 0.8602344, 0.61778736]
predicted location:[array([0.58308107, 0.5073111, 0.8015503, 0.8796249 ], dtype=float32),
array([0.50063336, 0.51950413, 0.8553538 ], dtype=float32), array([0.70523787,
0.5702946, 0.5486614, 0.804874 ], dtype=float32), array([0.29035494, 0.5391314,
0.31666473, 0.7510599 ], dtype=float32)]
[img-26 actual]:[0.3097222222222222, 0.4854166666666666, 0.4805555555555557, 0.5625, 1]
[img-26 actual]:[0.7694444444444445, 0.5166666666666667, 0.4222222222222222,
0.6166666666666667, 2]

```

Fig. 5. Predict image

```

dtype=float32)]
[img-25 actual]:[0.2125, 0.5541666666666667, 0.325, 0.8833333333333333, 2]
[img-25 actual]:[0.7652777777777777, 0.5229166666666667, 0.4694444444444444,
0.8458333333333333, 2]
predicted class is [2, 2, 1, 2]
predicted class value:[0.99591124, 0.99006337, 0.8602344, 0.61778736]
predicted location:[array([0.58308107, 0.5073111, 0.8015503, 0.8796249 ], dtype=float32),
array([0.50063336, 0.51950413, 0.8553538 ], dtype=float32), array([0.70523787,
0.5702946, 0.5486614, 0.804874 ], dtype=float32), array([0.29035494, 0.5391314,
0.31666473, 0.7510599 ], dtype=float32)]
[img-26 actual]:[0.3097222222222222, 0.4854166666666666, 0.4805555555555557, 0.5625, 1]
[img-26 actual]:[0.7694444444444445, 0.5166666666666667, 0.4222222222222222,
0.6166666666666667, 2]

```

Fig. 6. Predict image

## 5 Conclusion

We have put a lot of effort into this project. From the very beginning, we have blurred the concept of deep learning. Now we can write our own model, adjust the hyperparameters, and complete the project completely. In this, we have learned a lot and improved the ability to solve problems.

## Bibliography

- [1] Actions, S. . (2015). Stanford 40 actions. In *Stanford 40 Actions*.
- [2] Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.-Y., and Berg, A. C. (2016). SSD: Single shot multibox detector. In *ECCV*.
- [3] palmetto (2015). palmetto. In <https://www.palmetto.clemson.edu/palmetto/>.
- [4] Qassim, H., Verma, A., and Feinzimer, D. (2018). Compressed residual-vgg16 cnn model for big data places image recognition. In *2018 IEEE 8th Annual Computing and Communication Workshop and Conference (CCWC)*, pages 169–175.
- [5] Redmon, J. and Farhadi, A. (2018). Yolov3: An incremental improvement. *arXiv*.
- [6] Tzutalin (2015). Labelimg. In <https://github.com/tzutalin/labelImg>.