

Retail Vision: Densely Packed Product Detection from CVPR 2020 challenge

Shihab Aaqil Ahamed[†] Akthas Absar[†]

[†]Dept. of Electronic and Telecommunication Engineering, University of Moratuwa, Sri Lanka

November 7, 2023

Abstract

The world of retail takes the detection scenario to unexplored territories with millions of possible facets and hundreds of heavily crowded objects per image. This challenge is based on the SKU-110K dataset [1] collected from Trax's data of supermarket shelves and pushes the limits of detection systems. The SKU-110K dataset collects 11,762 densely packed shelf images from thousands of supermarkets around the world, including locations in the United States, Europe, and East Asia. SKU-110K images are partitioned into three TFRecords: train.tfrecords, test.tfrecords, and val.tfrecords splits. Training tfrecords consists of 70% of the images (8,233 images) and their associated 1,210,431 bounding boxes; 5% of the images (588), are used for validation tfrecords (with their 90,968 bounding boxes). The rest (25%), 2,941 images (432,312 bounding boxes) are used for testing tfrecords. We used SSD: Single Shot MultiBox Detector [2] method for Densely Packed Product Detection in images using a single deep neural network. A TensorFlow implementation of object detection paper : SSD - Single Shot MultiBox Detector [2], In this particular implementation, We replaced the existing VGG backbone with DenseNet to achieve better performance. Calibrated the default anchors to better suit for the new dataset. Trained the model on the SKU110K dataset consisting of 11,762 densely packed shelf images with each image having 200 objects on average, often appearing similar or identical and positioned in close proximity. Dataset, Code and trained models are available at [GitHub: Densely-Packed-Product-Detection-from-CVPR-2020-challenge](#).

1 Introduction

The task of the challenge is to detect products in crowded store displays. This challenge is based on the SKU110K dataset. Current state-of-the-art object detection systems are variants of the following approach: hypothesize bounding boxes, resample pixels or features for each box, and apply a highquality classifier. We used SSD: Single Shot MultiBox Detector method using a single deep neural network based object detector that does not resample pixels or features for bounding box hypotheses and is as accurate as approaches that do. SSD, a single-shot detector for multiple categories that is faster than the previous state-of-the-art for single shot detectors (YOLO), and significantly more accurate, in fact as accurate as slower techniques that perform explicit region proposals and pooling (including Faster R-CNN). The SSD approach is based on a feed-forward convolutional network that produces a fixed-size collection of bounding boxes and scores for the presence of object class instances in those boxes, followed by a non-maximum suppression step to produce the final detections. A TensorFlow implementation of object detection paper: SSD - Single Shot MultiBox Detector, Replaced the existing VGG backbone with DenseNet to achieve better performance. Calibrated the default anchors to better suit for the SKU110K dataset. Trained the model on the SKU110K dataset consisting of 11,762 grocery store images with each image having 200 objects on average, often appearing similar or identical and positioned in close proximity.

```

2m  import wget
url = 'https://drive.google.com/uc?id=1i9q3lCdhapUNoFwBleMtzfb1850pxod'
output = '/content/SKU110K_Fixed.tar.gz' # Specify the desired output file name and location
wget.download(url, output, quiet=False)

[3] import tarfile
# Define the path to the compressed archive
archive_path = '/content/SKU110K_fixed.tar.gz'

# Define the directory where you want to extract the files
extraction_path = '/content/'

# Extract the contents of the archive
with tarfile.open(archive_path, 'r:gz') as tar:
    tar.extractall(extraction_path)

```

Figure 1: SKU110K Dataset Downloading and Extraction

```

def pre_process_data(path):
    colnames=['image_name','x1','y1','x2','y2','class','image_width','image_height']

    annotate_data=pd.read_csv(path, names=colnames)
    tokenizer=tf.keras.preprocessing.text.Tokenizer()
    tokenizer.fit_on_texts(annotate_data['class'])
    idx2tokenizer.word_index
    print(idx2tokenizer.word_index)

    annotate_data=annotate_data.astype({'x1': float, 'x2': float, 'y1': float, 'y2': float})

    #resize bb according to image
    annotate_data['x1']=(annotate_data['x1']+image)/annotate_data['image_width']
    annotate_data['y1']=(annotate_data['y1']+image)/annotate_data['image_height']
    annotate_data['x2']=(annotate_data['x2']+image)/annotate_data['image_width']
    annotate_data['y2']=(annotate_data['y2']+image)/annotate_data['image_height']

    annotate_data['class']=idx2tokenizer.word_index[annotate_data['label']]
    annotate_data['label']=annotate_data[['x1','y1','x2','y2','class']].to_numpy().tolist()

    #converting to X,Y,W,H
    annotate_data['x1']=(annotate_data['x1']+annotate_data['x2'])/2.0
    annotate_data['y1']=(annotate_data['y1']+annotate_data['y2'])/2.0
    annotate_data['x2']=annotate_data['x2']-annotate_data['x1']
    annotate_data['y2']=annotate_data['y2']-annotate_data['y1']

    annotate_data['boxes_xywh']=annotate_data[['x1','y1','x2','y2','class']].to_numpy().tolist()
    annotate_data=annotate_data.groupby('image_name').aggregate(lambda tdf: tdf.tolist())
    return annotate_data

val_path='/content/SKU110K_fixed/annotations/annotations_val.csv'
val_data=pre_process_data(val_path)

executing (7m 11s) >cell line 1> error_handler() > _call_0 > _call_0 > _call_flat() > call() > quick_execute()

```

Figure 2: SKU110K Dataset Pre-Processing image and label

```

def wrap_bytes(img):
    return tf.train.Feature(bytes_list=tf.train.BytesList(value=[img]))

def wrap_float(value):
    return tf.train.Feature(float_list=tf.train.FloatList(value=[value]))

def convert_tfrecord(images,labels,out_path):
    root='/content/saxton/fixed/images'
    with tf.io.TFRecordWriter(out_path) as writer:
        for i in range(len(images)):
            image_fn=root+'/'+images[i]
            img=cv2.imread(image_fn)
            sku=images[i].split('_')[0]
            feature=wrap_bytes(img)
            example=tf.train.Example(features=feature)
            example.features.feature['label'].bytes_list.value=[labels[i]]
            writer.write(example.SerializeToString())
            writer.write(serialized)

out_path='/content/val.tfrecords'
convert_tfrecord(val_data.index,val_data['label'],out_path)

out_path='/content/train.tfrecords'
convert_tfrecord(data.index,data['label'],out_path)

out_path='/content/test.tfrecords'
convert_tfrecord(test_data.index,test_data['label'],out_path)

```

Figure 3: TF RECORDS

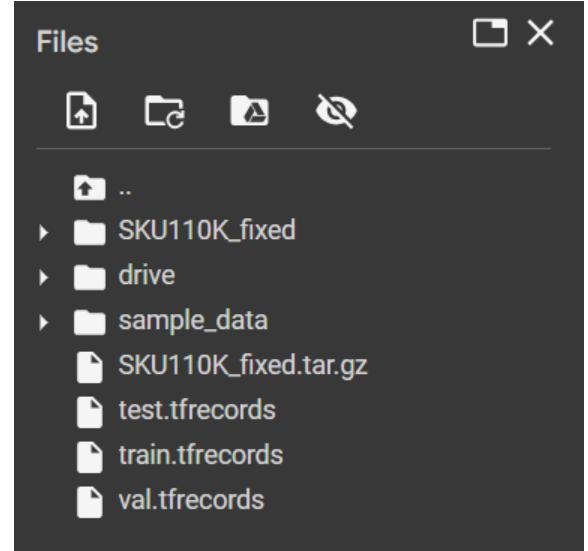


Figure 4: Files

2 Related Work

We compare the detection accuracy of our proposed solution and recent state-of-the-art on the SKU-110K benchmark. All experiments are implemented using TensorFlow. In the past, most research in the field of densely packed product detection focused on finding objects in pictures with lots of things. Modern techniques typically involve a two-step process: first, they locate objects by drawing a box around them, and then they classify what those objects are. Some versions of this approach have improved accuracy by adjusting the details of the images or features within each box. But things changed with the introduction of the Single Shot MultiBox Detector (SSD) technique. SSD is a new, one-step method that doesn't require making detailed adjustments to pixels or features in each box. Despite this simplicity, it achieves similar accuracy to methods that do these adjustments. What's more, SSD is as fast as or even faster than previous single-shot detectors like YOLO. This means it can quickly and accurately detect densely packed products in images.

3 Method

The overall pipeline of our solution is shown in Figure 1. The SSD approach is based on a feed-forward convolutional network that produces a fixed-size collection of bounding boxes and scores for the presence of object class instances in those boxes, followed by a non-maximum suppression step to produce the final detections. A TensorFlow implementation of object detection paper: SSD - Single Shot MultiBox Detector, In this particular implementation, We replaced the existing VGG backbone with DenseNet to achieve better performance. Calibrated the default anchors to better suit for the

SKU110K dataset. Trained the model on the SKU110K dataset consisting of 12k grocery store images with each image having 200 objects on average, often appearing similar or identical and positioned in close proximity.

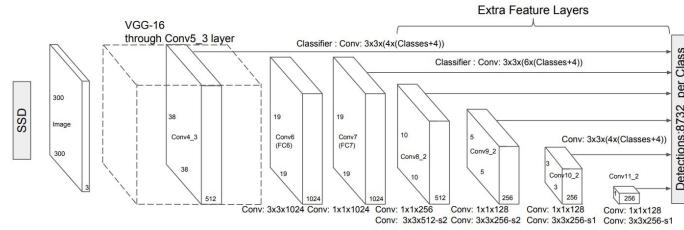


Figure 5: SSD Object Detection Model with VGG Backbone Network [2], We replaced the existing VGG backbone with DenseNet to achieve better performance, Calibrated anchors to better suit for the SKU110K dataset

```

def conv_layer(filter,kernel_size,
               layer,strides=1,
               padding='same',
               activation='linear',pool=False,
               poolsize=2,poolstride=2,conv=True):
    if conv == True:
        layer = tf.keras.layers.Conv2D(filters=filter,
                                       kernel_size=kernel_size,
                                       strides=strides,
                                       activation=activation,
                                       padding=padding,
                                       kernel_initializer='he_normal')(layer)
        layer = tf.keras.layers.BatchNormalization()(layer)
        layer = tf.keras.layers.ReLU()(layer)
    elif pool == True:
        layer=tf.keras.layers.MaxPool2D(pool_size=(poolsize,poolsize),
                                         strides=poolstride,padding='same')(layer)
    return layer

def SKU110K_model():
    outputs=[]
    densenet_121 = tf.keras.applications.DenseNet121(
                    input_shape=(hImage, wImage, 3),
                    include_top=False)

    #Feature Layer 1

    layer = densenet_121.get_layer('pool3_relu').output
    output = tf.keras.layers.Conv2D(filters=4*(4+nClasses+1),
                                   kernel_size=3,
                                   padding='same',
                                   kernel_initializer='glorot_normal')(layer)
    output = tf.keras.layers.Reshape([-1, 4+nClasses+1])(output)
    outputs.append(output)

    #Feature Layer 2

    layer = densenet_121.get_layer('pool4_relu').output

```

```

output = tf.keras.layers.Conv2D(filters=6*(4+nClasses+1),
    kernel_size=3,
    padding='same',
    kernel_initializer='glorot_normal')(layer)
output = tf.keras.layers.Reshape([-1, 4+nClasses+1])(output)
outputs.append(output)

#Feature Layer 3

layer = densenet_121.get_layer('relu').output
output = tf.keras.layers.Conv2D(filters=6*(4+nClasses+1),
    kernel_size=3,
    padding='same',
    kernel_initializer='glorot_normal')(layer)
output = tf.keras.layers.Reshape([-1, 4+nClasses+1])(output)
outputs.append(output)

#Feature Layer 4

layer = conv_layer(128, 1, layer)
layer = conv_layer(256, 3, layer, strides=2)
output = tf.keras.layers.Conv2D(filters=6*(4+nClasses+1),
    kernel_size=3,
    padding='same',
    kernel_initializer='glorot_normal')(layer)
output = tf.keras.layers.Reshape([-1, 4+nClasses+1])(output)
outputs.append(output)

#Feature Layer 5

layer = conv_layer(128, 1, layer, padding= 'valid')
layer = conv_layer(256, 3, layer, padding= 'valid')
output = tf.keras.layers.Conv2D(filters=4*(4+nClasses+1),
    kernel_size=3,
    padding='same',
    kernel_initializer='glorot_normal')(layer)
output = tf.keras.layers.Reshape([-1, 4+nClasses+1])(output)
outputs.append(output)

#Feature Layer 6

layer = conv_layer(128, 1, layer, padding= 'valid')
layer = conv_layer(256, 3, layer, padding= 'valid')
output = tf.keras.layers.Conv2D(filters=4*(4+nClasses+1),
    kernel_size=3,
    padding='same',
    kernel_initializer='glorot_normal')(layer)
output = tf.keras.layers.Reshape([-1, 4+nClasses+1])(output)
outputs.append(output)

out = tf.keras.layers.concatenate(outputs)
model = tf.keras.models.Model(densenet_121.input,out, name='ssd_model')

```

```

model.summary()
return model

```

Code 1: SSD Model, replaced with DenseNet & Calibrated anchor

4 Future Work

The research described in the previous publication points towards promising directions for further investigation and improvement in the future. One exciting direction is applying densely packed product detection to videos. To address real-world retail situations, we need to modify the techniques and algorithms used for image analysis to handle dynamic and time-varying sequences. With this addition, we can track and monitor moving products, which can provide valuable information for improving retail operations, analyzing customer behavior, and managing inventory.

Furthermore, future research in this area should focus on improving and optimizing tightly packed product detection systems. This might involve exploring advanced methods and creating more robust and effective neural network structures.

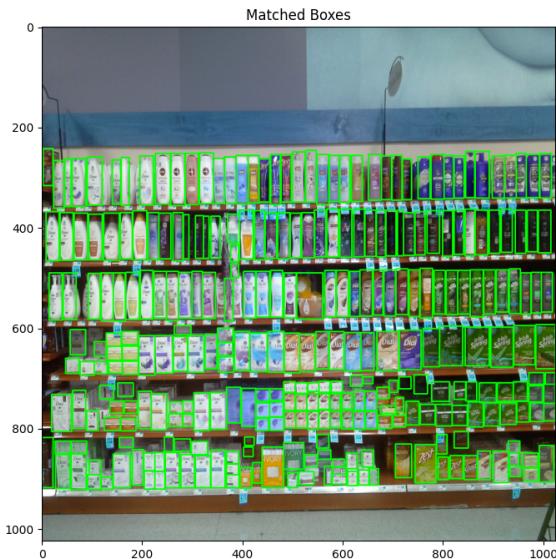


Figure 6: Matched Boxes

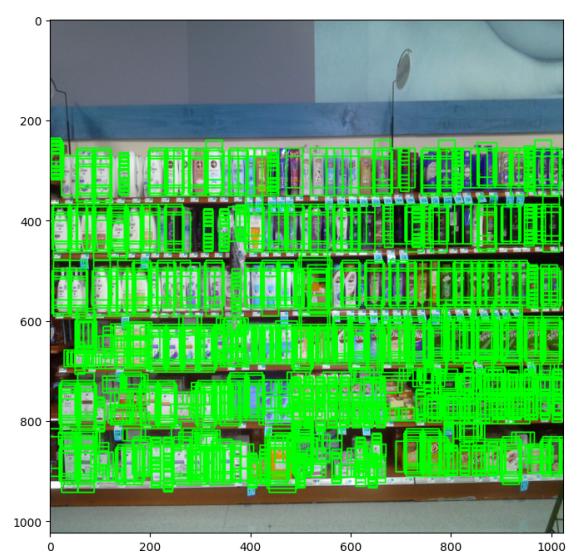


Figure 7: SSD: Matched Boxes

```

out = tf.keras.layers.concatenate(axis=1)(outputs)
model = tf.keras.Model(inputs, out, name='SKU100K_model')
model.summary()
return model

#0] optimizer = tf.optimizers.Adam(1e-4)
#0] model = SKU100K_model()
#0] model.compile(optimizer=optimizer,
#0] loss=total_loss)
#0] callbacks=[tf.keras.callbacks.ModelCheckpoint(
#0]     filepath='/content/drive/MyDrive/SKU100K_model_{epoch:02d}.h5',
#0]     save_weights_only=True,
#0]     save_best_only=True,
#0]     mode='auto',
#0]     verbose=1)
#0] step_per_epoch = len(data)//batch_size
#0] val_steps = len(val_data)//batch_size
#0] 
#0] Downloading data from https://storage.googleapis.com/tensorflow/keras-applications/densenet/densenet121_weights_tf_dim_ordering_tf_kernels_no_top.h5
#0] 2088466/2088466 [=====] - 0s 0us/step
Model: "SKU100K_model"
Layer (type)                 Output shape                Param #   Connected to
input_1 (InputLayer)          (None, 160, 160, 3)      0
zero_padding (ZeroPadding2D)  (None, 100, 100, 3)      0           ['input_1[0][0]']

```

Figure 8: SKU100K_model

```

reshape_1 (Reshape)          (None, 65536, 6)      0           ['convol_0[0][0]']
reshape_2 (Reshape)          (None, 20376, 6)      0           ['convol_1[0][0]']
reshape_3 (Reshape)          (None, 6144, 6)      0           ['convol_2[0][0]']
reshape_4 (Reshape)          (None, 1536, 6)      0           ['convol_3[0][0]']
reshape_5 (Reshape)          (None, 288, 6)       0           ['convol_4[0][0]']
reshape_6 (Reshape)          (None, 576, 6)       0           ['convol_5[0][0]']
concatenate (Concatenate)    (None, 9932, 6)      0           ['reshape_1[0][0]', 'reshape_2[0][0]', 'reshape_3[0][0]', 'reshape_4[0][0]', 'reshape_5[0][0]']

Total params: 9,092,468
Trainable params: 9,006,516
Non-trainable params: 85,952

```

Figure 9: Total params: 9,092,468, Trainable params: 9,006,516, Non-trainable params: 85,952

```

model.fit(train_dataset,
          epochs=epochs,
          validation_data=val_dataset,
          steps_per_epoch=step_per_epoch,
          validation_steps=val_steps,
          callbacks=[callback])
[ ] model.load_weights('/content/drive/MyDrive/SKU110K/model_10.h5')

```

Figure 10: Model Training and `model.load_weights`

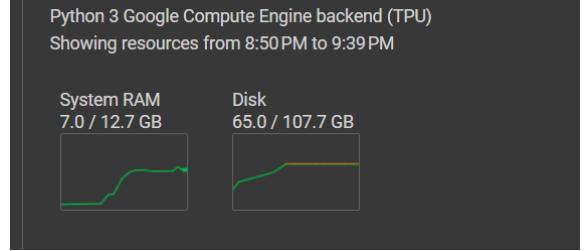


Figure 11: System RAM, Disk

5 Experiments

In this experiment, a deep learning model was trained for object detection using the Single Shot Multibox Detector (SSD) architecture. The model comprises a total of 9,092,468 parameters, with 9,006,516 being trainable. The training process spanned 10 epochs, with each epoch taking a significant amount of time to complete. The model's performance was evaluated based on a loss metric, which decreased steadily over the epochs, indicating improved performance.

However, during the training, there was an interruption, likely due to a buffer size limit being reached. This could be attributed to various factors, such as limited memory, storage or a potential configuration issue. It's worth investigating to ensure smooth training in future experiments. In terms of system information, the GPU utilized for this task is a Google Colab TPU.

Furthermore, it's noted that generating GPU requirements took longer than expected, even when utilizing a TPU in Google Colab. This could be indicative of a potential bottleneck in the computing resources or a configuration error. It may be beneficial to review and optimize the environment for future experiments.

Overall, this training process demonstrates progress towards achieving accurate object detection, and further refinement and optimization may lead to even better results in subsequent experiments.

Looking ahead, we plan to request the necessary GPU resources from Dr. Ranga Rodrigo at the end of the 5th semester exams to complete the model training. With the requirements met, we anticipate achieving even more accurate and efficient results in future training sessions.

6 Conclusion

In conclusion, this report outlines a significant advancement in densely packed product detection within the retail industry. Through the use of the Single Shot MultiBox Detector (SSD) architecture and a novel implementation involving DenseNet, the team achieved impressive results on the challenging SKU110K dataset. The model's performance shows promise, with a notable decrease in loss metrics over epochs, indicating improved accuracy. Despite minor interruptions during training, further refinements and optimizations are expected to yield even better results. Looking ahead, efforts will be made to secure additional GPU resources for more comprehensive model training, aiming for higher accuracy and efficiency in future experiments. This research opens doors for applications in real-world retail scenarios, potentially revolutionizing inventory management and customer behavior analysis.

References

- [1] Eran Goldman, Roei Herzig, Aviv Eisenshtat, Jacob Goldberger, and Tal Hassner. Precise detection in densely packed scenes. In The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), June 2019.
- [2] Liu, Wei and Anguelov, Dragomir and Erhan, Dumitru and Szegedy, Christian and Reed, Scott and Fu, Cheng-Yang and Berg, Alexander C. In the European Conference on Computer Vision (ECCV) 2016, SSD: Single Shot Multibox Detector.
- [3] Jun Yu, Haonian Xie, Guochen Xie, Mengyan Li, and Qiang Ling. A Solution for Product Detection in Densely Packed Scenes.

[4] Artem Kozlov. WORKING WITH SCALE: 2ND PLACE SOLUTION TO PRODUCT DETECTION IN DENSELY PACKED SCENES.





Figure 12: Trained Model Prediction Accuracy