
Place of Attention Matters!

An End-to-End Object Detection with Vision Transformation

Saeed Firouzi
Sharif University of Technology
saeedmr881@gmail.com

Abstract

In the object detection task, the purpose is to find the class of object and a bounding box around it. Most works have focused on just finding the class of object without considering bounding box features properly. We present a new method that focuses on relationships between patches of the image as a feature for bounding box detector.

Also, we combine convolutional neural network as a local feature detector and Transformer network as a long-distance feature detector. We were also inspired by the method that has been used in Transformer as a relationship between patches in the image.

Our implementation can perform in real-time and improve the accuracy of previous works. The training code and pre-training model are available at

https://github.com/saeed5959/object_detection_transformer

to help open source community.

1 Introduction

Most works has shown that convolutional networks are basically considered as local feature detectors. It is quite easy to see if we are familiar with CNN architecture. Using a filter with small dimensions like : 3×3 , 5×5 , ... , 9×9 can just get a local feature around a pixel value. The strength of using CNN is that we produce a high number of these filters to capture features with different styles. Some attempts like max-pooling, deformable convolution, dilated casual convolution have been applied to increase receptive fields of a pixel, to look further and to see more around it. But still they have not overcome this issue.

On the other hand, after introducing Transformer, and specifically vision transformer, many works have been tried to apply transformer to image input to achieve a long-term dependency between pixels or patches of image. The main idea behind the transformer is to get a weighted sum for a vector by point-wise multiplying of vectors as a similarity term.

That means by giving a vector to a transformer, we will get a vector with the same dimension as the input but some values related to a feature have been increased or decreased based on whether that feature exists in other vectors or not.

We were inspired by this point from transform, to use it for bounding box feature to be our major contribution to this work. We call it place of attention matrix (POA matrix) to be a feature guide for bounding box detection.

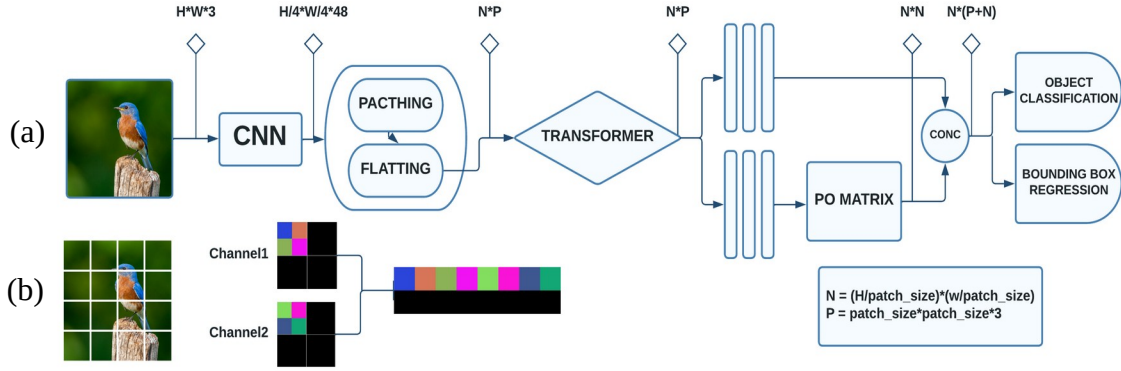


Figure 1. **An illustration of model architecture (PoA-ViT).** In (a) we pass the image to the CNN network. After patching and flattening, we give it to the transformer block. Then concatenate it with POA matrix and using 2 linear layers, we can get the class and bounding box of objects. In (b) we show how we patch and flat output of the CNN network to make it one dimensional vector for the transformer.

Mostly there are two kinds of network: two-stage and one-stage networks. In two-stage networks, they mostly use a region-proposal-network(RPN) for suggesting some box for objects. These boxes are predefined and have different ratios and aspects. Ratios are from 20,50,100,200 pixels in size and aspects are from 1:1, 1;2, 2;1, ... size. In most cases, the number of these boxes becomes more than 2000 boxes. These proposal boxes have 2 drawbacks; first, they need a lot of computation, which makes it impossible for real-time tasks, and second, some objects can not be fit in these predefined boxes.

In a one-stage network, the main idea is to predict the bounding box coordinate directly. By dividing the input image into patches and predicting a bounding box for these patches, it is possible to make the network work in real-time tasks.

In this method, we will use a one-stage network, and divide the input image into patches and pass patches to the network that consisted of CNN and Transformer to get a feature vector for every patch. Then we will calculate POA matrix by multiplying a feature vector of a patch on other vectors to get a similarity term as a factor of how much this patch is similar to another patch. And then concatenate this POA matrix with feature vector and then apply a linear layer to get a score for class and another linear layer that outputs 4 values: x of center, y of center, width and length as bounding box parameters.

2 Model Architecture

2.1 CNN Network

Instead of giving raw pixel value to network and just applying FC layer [Detr] that can't get local features properly, we use a convolutional network by inspiring from ResNet [] architecture. The main idea in ResNet is using skip-connection to pass data directly into deeper layers. Also, we just use two max pooling in our block, to keep spatial features for model to detect the location of features in image. because it can decrease the resolution of input 4x smaller.

As a regularization factor, we use both batch normalization and layer normalization due to a different variety of input images.

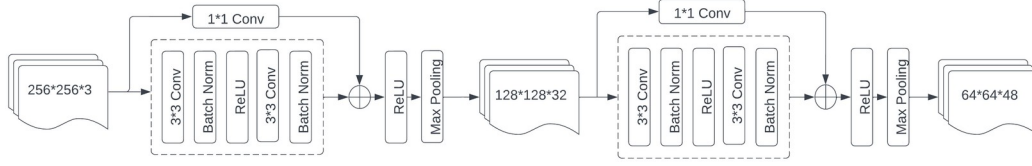


Figure 2. The CNN network

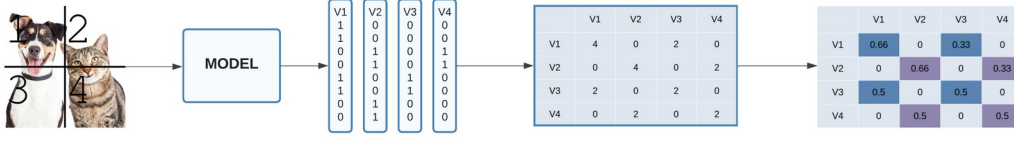


Figure 3. **POA matrix**. If we have an image with 4 patches and for every patch we have a feature vector, then the task of POA matrix would be to find similar patches by point-wise multiplying of feature vectors to get a similarity factor. This helps the model to know which patches belong to an object, then the model can predict the bounding box of object properly. In this case, patch 1,3 has similar features, then they belong to an object and also the 2,4 patches.

2.2 Patching and Flatting

After applying CNN network in input image with size $H \times W$ with 3 channels as RGB values, we will get data with size $H/4 \times W/4$ with 48 channels as a feature values.

Now, inspired by YOLO architecture [], we divide this into the same size patches and then flat them in raster order alongside their channels to get a feature vector for every patch. See Figure 1.b.

3.3 Transformer

Now we have N patches with size P . In Transformer we add positional encoding to feature vectors to give spatial information to the model for bounding box prediction. This positional encoding is an embedding layer, giving input of $(0, \text{number of patches})$ in a raster order to be learnable. If we consider a feature vector as a container of features of different objects, then Transformer's task is to increase values of some objects in this feature vector with respect to other patch feature vector, wheter there is a object in that patch or not.

We use multi-head attention with a sequence of linear layers [], and giving feature vector to query, key, value as inputs of Transformer block.

3.4 Place of Attention Matrix (POA Matrix)

By dividing the image into small patches, we can capture small objects. But for large objects, we should make a relationship between patches. By point-wise multiplying every patch with other patches, we can get a similarity and positional relationship between them. Then we have a matrix that we call the place of attention matrix (POA matrix). For every patch vector, we concatenate its corresponding row of POA matrix and by giving it to two linear layers, we can classify the patch and get the bounding box coordinate if it does not belong to the background class. This method is inspired by attention in transformer[].

4 TRAINING

4.1 Training Strategy

Despite YOLO [], which makes the center patch responsible for classifying the object and bounding box coordinate, and despite Faster-RCNN [], which makes all pixels in the box

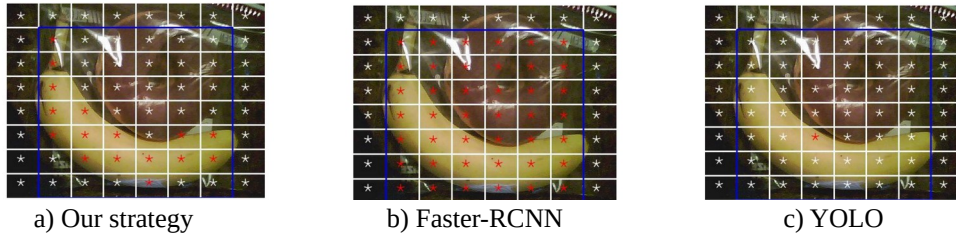


Figure 4. **Illustration of different training methods.** in (a) we just calculate the loss for red-star patches that contain objects. In (b) it uses all pixels of bounding box. As you can see, some pixels and patches don't contain the object and are just misinformation. In (c) it just uses one center patch and ignore some other patches that contain objects.

responsible for object detection, we use segmentation data beside the object box, and instead of making all patches inside of box responsible for object detection, just use patches that have intersection with segmentation data.

This approach helps to ignore misinformation that may come from patches inside the box that does not contain the object.

4.2 Augmentation

Instead of pre-training the model on the ImageNet dataset, we use the augmentation method to increase the number of images for better accuracy.

For augmentation, we use shift, scale, rotation and also the Cutout method[] to force model to learn different features of an object.

5 EXPERIMENT

We train the model on COCO dataset [] that includes 118,000 images with 90 categories. Every image contains a bounding box for objects and a segmentation image.

We evaluate our model based on two factors; accuracy and inference time. Then we will compare our model with Faster-RCNN and YOLO-V5 and Detr.