

DEEP LEARNING PROJECT

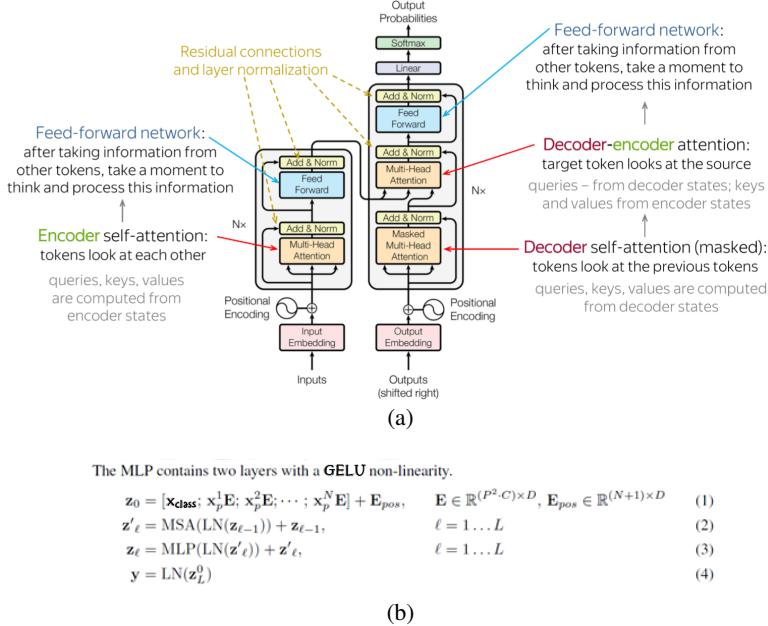
FINAL REPORT

ABSTRACT

We have focused on upgrading the Vision Transformer (ViT) model, which is commonly used in NLP, for image classification. In order to enhance the performance of the ViT model, we plan to implement an additional preprocessing step using the SSD model before the Input Embedding process.

1 BACKGROUND

1.1 ViT



As illustrated in the figure above, the ViT model is composed of an encoder and decoder structure. In this process, we can see that the input, which has been through the embedding process, is used as the input for the Multihead Attention (MSA).

When utilizing a large-sized input image, the self-Attention operation must be performed on all patches. This means the inference is based on the entire image, rather than focusing on a specific object. This approach can be seen as observing the overall context rather than focusing on a specific point, especially when dealing with images requiring high accuracy. However,

there are instances where we want to closely observe a specific point. For example, in a photo of a crowd, we might want to identify a 'Person-Name' rather than just confirming the presence of 'Person'. Similarly, when interpreting medical images (CT, MRI, etc.), we might want to focus on a particular abnormality to diagnose a disease. In these special scenarios, rather than inferring from the whole picture through the ViT, we would want to intensely analyze the content within a specific bounding box. For these reasons, we plan to use SSD as a preprocessing step to implement a higher-level classifier.

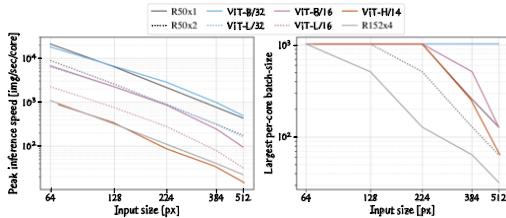
1.2 Attention Map



Through the Attention map, we can visually identify the areas being attended to in the photo. Let me explain using the above photo as an example.

Unlike in the bird photos where Attention functioned well, it can be visually observed that Attention is not properly functioning in the baby photos.

1.3 Comparison of Inference Speed Depending on Input Size in ViT Model

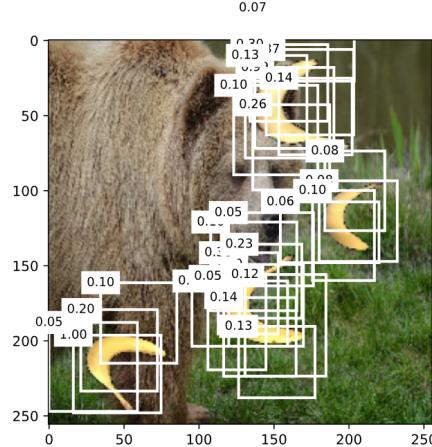


This graph represents the correlation between input size and inference speed in several models, including ViT. It shows a high dependency on the input size. Therefore, if smaller images are inputted through single, the inference speed of ViT will improve.

However, there exists a trade-off between the increased time consumption due to single and the reduced ViT inference time due to the decrease in input size.

1.4 Single Shot multibox Detection (SSD)

SSD differs from traditional predefined anchor boxes by generating anchor box centers only at reasonable locations. This approach prevents the creation of anchor boxes in unnecessary locations, thereby reducing unnecessary computations. The figure below illustrates an example of how anchor boxes are created in SSD.



We have also trained the model with objects other than bananas and confirmed that the prediction bounding box is accurately generated at the corresponding location.



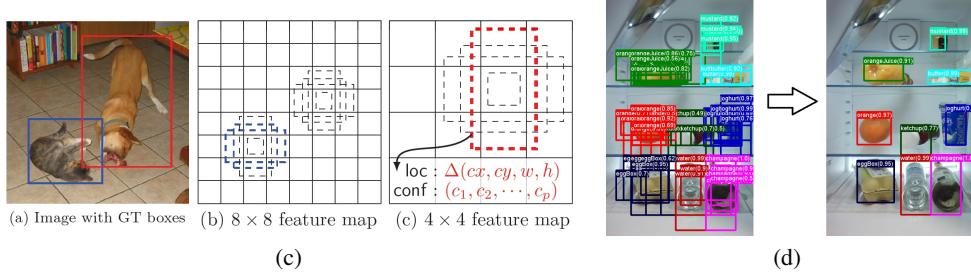


Figure 1: (c) Multiscale object (d) NMS

1.5 NMS(*Non-Maximum Suppression*) Algorithm

Before understanding the NMS algorithm, it's important to understand confidence scores.

When we say Confidence Score, we mean that we have n% confidence that the network has produced the correct answer.

means that when the network comes up with the correct answer, it has n% confidence in that answer.

Every bounding box carries a confidence score that indicates how well it captures the corresponding object. Bounding boxes with confidence scores below a certain threshold are eliminated, which is a preliminary filtering process for boxes with low-level certainty.

The remaining bounding boxes are then sorted in descending order based on their confidence scores.

Algorithm 1 Non-Max Suppression

```

1: procedure NMS( $B, c$ )
2:    $B_{nms} \leftarrow \emptyset$ 
3:   for  $b_i \in B$  do
4:      $discard \leftarrow \text{False}$ 
5:     for  $b_j \in B$  do
6:       if same( $b_i, b_j$ )  $> \lambda_{nms}$  then
7:         if score( $c, b_j$ )  $>$  score( $c, b_i$ ) then
8:            $discard \leftarrow \text{True}$ 
9:         if not  $discard$  then
10:           $B_{nms} \leftarrow B_{nms} \cup b_i$ 
11:   return  $B_{nms}$ 
```

Taking the bounding box at the front as a standard, the Intersection over Union (IoU) values with other bounding boxes are calculated. Boxes with an IoU exceeding a certain threshold are removed. This is because the higher the IoU between bounding boxes, the more likely they are detecting the same object.

This process is sequentially executed to compare and eliminate all bounding boxes. The higher the confidence threshold and the lower the IoU threshold, the more bounding boxes are removed.

1.6 Related Works

Attention Is All You Need

Ashish Vaswani^{1*}, Noam Shazeer², Niki Parmar³, Jakob Uszkoreit⁴, Llion Jones⁵, Aidan N. Gomez⁶, Łukasz Kaiser⁷, Illia Polosukhin⁸
Google Brain, Google Brain, Google Research, Google Research, Google Research, University of Toronto, aiida@cs.toronto.edu, lukasz Kaiser@google.com, illia.polosukhin@gmail.com

Abstract

The dominant sequence translation models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism that allows them to selectively attend to different parts of the input sentence. In this work, we show that this reliance on CNNs is not necessary and a pure transformer applied directly to sequences of image patches can perform very well. We demonstrate this by applying a multi-scale patch-based architecture of data and transferred to multiple mid-sized or small image recognition benchmarks (ImageNet, CIFAR-10, VTFB, etc.). Visual Transformer (ViT) attains excellent results compared to state-of-the-art convolutional networks while requiring significantly fewer computational resources to train.

(a) Attention Is All You Need¹

Published as a conference paper at ICLR 2021

AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

Alexey Dosovitskiy^{1*}, Lucas Beyer², Alexander Kolesnikov³, Dirk Weissenborn⁴, Xiaohua Zhai⁵, Thomas Unterthiner⁶, Mostafa Dehghani⁷, Matthias Minderer⁸, Georg Heigold⁹, Sylvain Gelly¹⁰, Jakob Uszkoreit¹¹, Neil Houlsby¹²
“equal technical contribution, equal advising
1Google Research, Brain Team
{adosovitskiy, neilhoulsby}@google.com

ABSTRACT

While the Transformer architecture has become the de-facto standard for natural language processing tasks, its applications to computer vision remain limited. In vision, attention is either applied in conjunction with convolutional networks, or used as a post-processing step to refine detections made by a model without an overall structure in place. We show that this reliance on CNNs is not necessary and a pure transformer applied directly to sequences of image patches can perform very well. We demonstrate this by applying a multi-scale patch-based architecture of data and transferred to multiple mid-sized or small image recognition benchmarks (ImageNet, CIFAR-10, VTFB, etc.). Visual Transformer (ViT) attains excellent results compared to state-of-the-art convolutional networks while requiring substantially fewer computational resources to train.

(b) Transformers For Image Recognition At Scale²

CvT-ASSD: Convolutional vision-Transformer Based Attentive Single Shot MultiBox Detector
1[†] Weiqiang Jin
School of Microelectronics and
Science, University
Shanghai, China
Postal Code: 200444
Email: weiqiangjin@sjtu.edu.cn
2[‡] Hang Yu
School of Microelectronics and
Science, University
Shanghai, China
Postal Code: 200444
Email: yuhang@sjtu.edu.cn
3[§] Xiuming Luo
School of Microelectronics and
Science, University
Shanghai, China
Postal Code: 200444
Email: yuhang@sjtu.edu.cn

Authors—Due to the success of Bidirectional Encoder Representations from Transformers (BERT) in natural language processing tasks, there has been a significant interest in applying similar approaches to computer vision tasks. However, most of these attempts have focused on image captioning and visual question answering tasks. In this work, we propose CvT-ASSD, a novel multi-scale patch-based vision-Transformer detector that achieves state-of-the-art performance on COCO and PASCAL3D+ datasets. CvT-ASSD consists of two main components: a multi-scale patch-based backbone and a multi-scale multi-head self-attention module. The backbone takes image patches as input and processes them sequentially. The multi-head self-attention module is used to capture global dependencies between patches. CvT-ASSD is able to achieve better performance than existing methods while being faster and more memory efficient. CvT-ASSD is also able to handle large-scale datasets like COCO and PASCAL3D+.

(c) Convolutional vision-Transformer Based Attentive Single Shot MultiBox Detector³

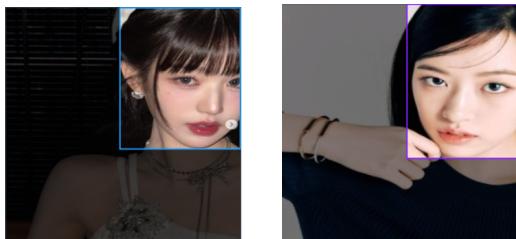
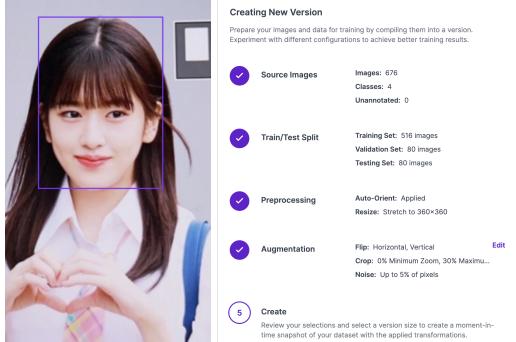
¹Attention Is All You Need

²Transformers For Image Recognition Scale

³Convolutional vision-Transformer Based Attentive Single Shot MultiBox Detector

2 Dataset and method

2.1 Dataset

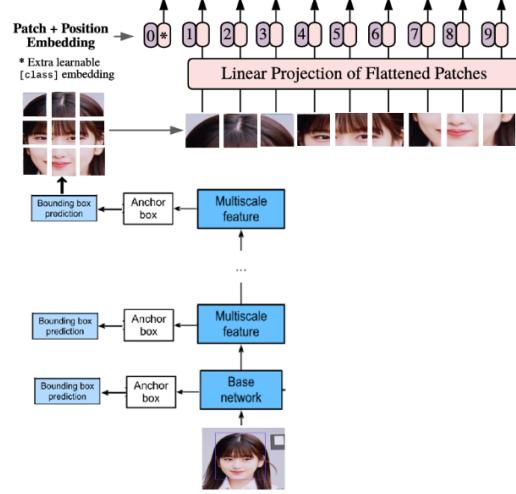


Jang Won-young, An Yu-jin, Kwon Eun-bi and other members picture. We used Roboflow to classify and set the Ground Truth for 671 original images.



After augmenting these images, the dataset was expanded to include 1548 training, 80 validation, and 80 testing images.

3 Method



1. Use SSD to crop the image area that will be input to ViT.
2. Insert the bounding box produced by SSD as an input to ViT.

Anticipating the location of the objects, the image is passed as input to the Vision Transformer. It is akin to providing ViT with images containing a single object.

This is expected to have a similar effect as noise reduction.

The aim of this process is to enhance the prediction accuracy of the ViT.

| ViT model Parameter | Parameter Value |
|---------------------|-----------------|
| n_patches | 16x16 (256) |
| model_dim | 768 |
| hidden_dim | 3072 |
| n_class | 4 |
| n_heads | 12 |
| n_layers | 12 |

4 Result

4.1 Environment

| | |
|--------------|-----------------------|
| Google Colab | Pro version |
| System RAM | 50GB |
| GPU | NVIDIA A100, V100, T4 |

4.2 SSD

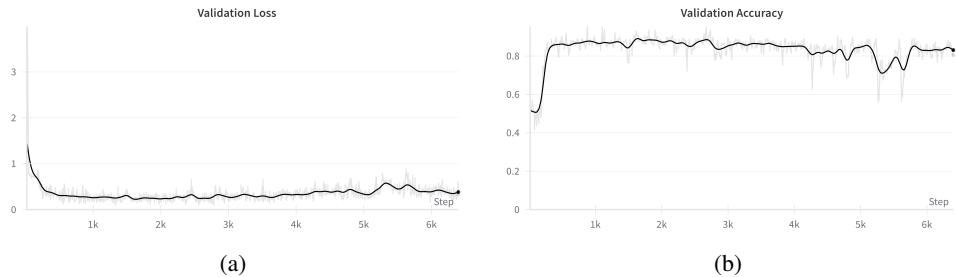


Figure 2: (a) Validation loss (b) Validation Accuracy

4.3 ViT

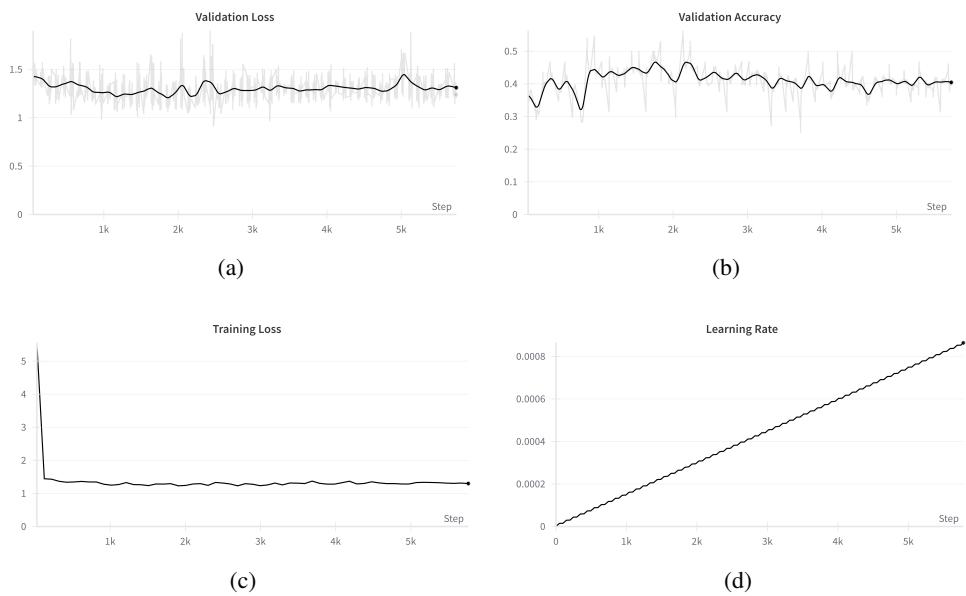


Figure 3: (a) Validation loss (b) Validation Accuracy
(c) Training loss (d) learning rate

5 *Discussion*

The anticipated benefit of using the ViT model and SSD's prediction bounding box together is an increase in ViT's inference speed due to a decrease in input size. However, as we are using the SSD's prediction bounding box, the speed is expected to slow down.

The primary goal is to accurately analyze a specific object, hence the design isn't focused on real-time inference.

Let's first consider why the result accuracy did not come out as well as expected. One possibility is that the dataset we manually classified did not become a good quality dataset because the criteria for setting the ground truth on the face were not consistent. Secondly, a large amount of training data is required to train the ViT model, but the dataset we created only contains just over 1800 images, which I think is a problem. Considering that it requires more data for training compared to other models, we even proceeded with augmentation, but it was insufficient. Lastly, I believe not being able to utilize a pre-trained model was a problem. As introduced in the paper, it is recommended to use a pre-trained model.

We've realized that the computing resources required to train the ViT are crucial. The long training time of up to 30 hours for ViT makes it hard to run experiments and check results, this is a regrettable aspect.

6 *Reference*

6.1 *Baseline*

*ViT baseline*⁴

*SSD baseline*⁵

⁴ViT baseline

⁵SSD baseline