

Transformer Architecture (Vision-Transformer)

• ABSTRACT

In recent years, Vision Transformers (ViTs) have emerged as a promising architecture in the field of computer vision, demonstrating remarkable performance on various image datasets. This project explores the application of Vision Transformers on the CIFAR dataset, known for its smaller image resolution compared to other benchmarks. The primary objective is to adapt and optimize the Transformer architecture for CIFAR images while maintaining high classification accuracy.

The project involves implementing a Vision Transformer model tailored specifically for CIFAR, leveraging self-attention mechanisms to capture intricate relationships between image patches. Extensive experimentation is conducted to fine-tune hyperparameters, including patch sizes, depth of the network, and attention mechanisms, to achieve optimal performance. Additionally, techniques such as data augmentation, regularization, and learning rate scheduling are employed to enhance the model's robustness and generalization.

Evaluation of the proposed CIFAR-specific Vision Transformer is performed against established baselines and state-of-the-art convolutional neural networks (CNNs). The model's performance is assessed based on accuracy, computational efficiency, and scalability. Comparative analysis sheds light on the strengths and limitations of Vision Transformers in handling smaller-sized image datasets like CIFAR.

The findings from this project contribute insights into the viability of Vision Transformers for low-resolution image classification tasks. The exploration of Vision Transformers on CIFAR dataset not only provides a deeper understanding of their adaptability but also offers valuable implications for leveraging transformer-based architectures in real-world applications with constrained computational resources and smaller image resolutions.

• **OBJECTIVE**

1. The primary objective of this project is to investigate and optimize the applicability of Vision Transformers (ViTs) specifically tailored for the CIFAR dataset.
2. Customize and refine the Vision Transformer architecture to effectively handle the unique characteristics of the CIFAR dataset.

• INTRODUCTION

In recent years, the Transformer architecture, initially popularized in natural language processing tasks, has garnered substantial attention in the realm of computer vision. Particularly, the emergence of Vision Transformers (ViTs) has sparked interest due to their remarkable success in handling image classification tasks on large-scale datasets like ImageNet. However, the adaptability and performance of ViTs on smaller image datasets, such as CIFAR, remain relatively unexplored.

The CIFAR dataset is renowned for its smaller image resolution compared to ImageNet, comprising 32x32 pixel RGB images across multiple classes. Traditional convolutional neural networks (CNNs) have historically dominated the landscape for CIFAR classification tasks. However, the inherent parallelization and self-attention mechanisms within Vision Transformers present an intriguing prospect for surpassing the capabilities of CNNs on such datasets.

This project delves into the adaptation and optimization of Vision Transformers explicitly designed for CIFAR. The primary focus is on tailoring the architecture to effectively handle lower-resolution images, exploiting the self-attention mechanism to capture intricate relationships between smaller image patches. By fine-tuning parameters and exploring various architectural modifications, this study aims to achieve competitive classification accuracy on the CIFAR dataset.

The investigation involves comprehensive experimentation, encompassing adjustments to patch sizes, network depth, attention mechanisms, regularization techniques, and other hyperparameters. Through rigorous evaluation and comparative analysis against established

CNN baselines and state-of-the-art models on CIFAR, this research seeks to ascertain the efficacy of Vision Transformers in this specific context.

The insights garnered from this endeavor are expected to shed light on the adaptability and limitations of Vision Transformers for low-resolution image classification tasks. Moreover, this exploration holds promise in providing valuable implications for leveraging transformer-based architectures in real-world applications where computational resources are constrained, and datasets exhibit smaller image resolutions, fostering a deeper understanding of their potential applicability beyond conventional high-resolution benchmarks.

• METHODOLOGY

1. Data Preparation

- Obtain the CIFAR dataset comprising 32x32 RGB images divided into predefined training and test sets.
- Preprocess the images by normalizing pixel values and applying suitable augmentation techniques (e.g., random crops, flips, rotations) to augment the dataset for better generalization.

2. Vision Transformer Architecture

- Implement the Vision Transformer model architecture, adapting it to suit the constraints and characteristics of the CIFAR dataset.
- Define the model's components, including embedding layers, positional encodings, multi-head self-attention mechanisms, feed-forward neural networks, and classification heads.

4. Training and Evaluation

- Train the CIFAR-specific Vision Transformer using the prepared dataset.
- Monitor the model's performance on the validation set, ensuring convergence and avoiding overfitting.
- Evaluate the trained model on the test set to measure classification accuracy and other relevant metrics.

7. Result and Conclusion

- Interpret the experimental results to draw conclusions regarding the efficacy and adaptability of Vision Transformers for CIFAR.

- Discuss the strengths, limitations, and potential avenues for further improvement or research in utilizing Vision Transformers for low-resolution image classification tasks.
- Provide insights and implications for leveraging transformer-based architectures in practical scenarios with smaller image resolutions and limited computational resources.

• CODE

```
!pip install tensorflow==2.8.0
```

```
!pip install keras==2.8.0
```

```
!pip install tensorflow_addons==0.20.0
```

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
import tensorflow as tf
```

```
from tensorflow import keras
```

```
tf_utils = tf.keras.utils
```

```
from tensorflow.keras import layers
```

```
import tensorflow_addons as tfa
```

```
num_class = 10
```

```
input_shape= (32,32,3)
```

```
(x_train, y_train), (x_test,y_test)= keras.datasets.cifar10.load_data()
```

```
print(f"x_train shape: {x_train.shape} - y_train shape: {y_train.shape}")
```

```
print(f"x_test shape: {x_test.shape} - y_test shape: {y_test.shape}")
```



```
learning_rate = 0.001

weight_decay = 0.001

batch_size= 256

num_epochs = 40

image_size = 72 #resizing input images

patch_size = 6

num_patches= (image_size//patch_size)**2

projection_dim= 64

num_heads = 4

transformer_units = [

    projection_dim*2,

    projection_dim

]          #size of transformer layers

transformer_layers = 8

mlp_head_units = [2048, 1024] #size of dense layers of the final classifier


data_augmentation = keras.Sequential(

    [

        layers.Normalization(),

        layers.Resizing(image_size, image_size),

        layers.RandomFlip("horizontal"),
```

```

        layers.RandomRotation(factor=0.02),

        layers.RandomZoom(height_factor=0.2, width_factor=0.2)

    ],

    name="data_augmentation"

)

data_augmentation.layers[0].adapt(x_train)    #compute mean and variance of x_train


def mlp(x, hidden_units, dropout_rate):    #multi layer perceptrons

    for units in hidden_units:

        x = layers.Dense(units, activation = tf.nn.gelu)(x)

        x= layers.Dropout(dropout_rate)(x)    #sets fraction of input to 0

    return x


class Patches(layers.Layer):

    def __init__(self,patch_size):    #constructor

        super(Patches, self).__init__()

        self.patch_size =patch_size


    def call(self, images):

        batch_size = tf.shape(images)[0]    #will dynamically determine the batch size at run
time

```

```

patches = tf.image.extract_patches(

    images=images,

    sizes= [1,self.patch_size, self.patch_size, 1],

    strides = [1,self.patch_size, self.patch_size, 1],  #step size between patches

    rates=[1,1,1,1],          #dilation rates.....1 means no dilation

    padding= "VALID",

)

patch_dims=patches.shape[-1]          #will flatten the patch so that each image is
presented as seq of patches

patches = tf.reshape(patches, [batch_size, -1,patch_dims])

return patches


plt.figure(figsize=(4,4))

image= x_train[np.random.choice(range(x_train.shape[0]))]

plt.imshow(image.astype("uint8"))

plt.axis("off")


resized_image = tf.image.resize(

    tf.convert_to_tensor([image]), size=(image_size, image_size)

)

```

```
patches= Patches(patch_size)(resized_image)

print(f"Image size: {image_size} X {image_size}")

print(f"Patch size: {patch_size} X {patch_size}")

print(f"Patches per image: {patches.shape[1]}")

print(f"Elements per patch: {patches.shape[-1]}")


n= int(np.sqrt(patches.shape[1]))

plt.figure(figsize= (4,4))

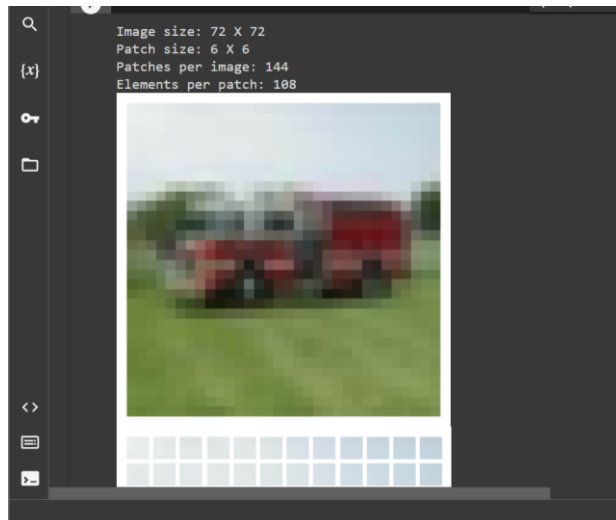
for i, patch in enumerate(patches[0]):

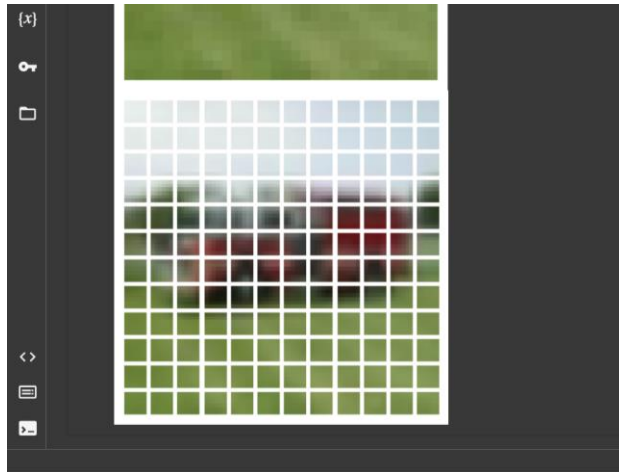
    ax = plt.subplot(n, n, i+1)

    patch_img = tf.reshape(patch, (patch_size, patch_size, 3))

    plt.imshow(patch_img.numpy().astype("uint8"))

    plt.axis("off")
```





```
class PatchEncoder(layers.Layer):
```

```
    def __init__(self, num_patches, projection_dim):
```

```
        super(PatchEncoder, self).__init__()
```

```
        self.num_patches= num_patches
```

```
        self.projection = layers.Dense(units= projection_dim)    #performs linear projection of
input patches to lower dimensions specified by proj dim
```

```
        self.position_embedding = layers.Embedding(              #instance of layer embedding
#embedding layer is a layer in Neural networks that maps the input info from from higher
dimension to o/p layer (lower dim) allowing network to learn more about the relation between
inputs to process data more efficiently
```

```
        input_dim = num_patches, output_dim = projection_dim
```

```
    )
```

```
    def call(self,patch):
```

```
        positions = tf.range(start=0, limit= self.num_patches, delta=1)    #to generate tensors of
seq num from 0 to range.....will represent the index of each patch
```

```
encoded = self.projection(patch) + self.position_embedding(positions)

return encoded
```

```
def create_vit_classifier():
```

```
    inputs = layers.Input(shape= input_shape)
```

```
    #Augment data
```

```
    augmented = data_augmentation(inputs)
```

```
    #Create patches
```

```
    patches = Patches(patch_size)(augmented)
```

```
    #Encode patches
```

```
    encoded_patches = PatchEncoder(num_patches, projection_dim)(patches)
```

```
    #create multiple layers of Transformer block
```

```
    for _ in range(transformer_layers):
```

```
        # layer normalization
```

```
        x1= layers.LayerNormalization(epsilon=1e-6)(encoded_patches)
```

```
        #create multi-head attention layer
```

```
        attention_output = layers.MultiHeadAttention(
```

```
            num_heads = num_heads, key_dim = projection_dim, dropout=0.1
```

```
        )(x1,x1)
```

```
        #skip connection 1
```

```

x2 = layers.Add()([attention_output, encoded_patches])

#layer normalization 2

x3= layers.LayerNormalization(epsilon=1e-6)(x2)

#MLP

x3 = mlp(x3, hidden_units= transformer_units, dropout_rate=0.1)

#skip connection 2

encoded_patch = layers.Add()([x3,x2])


#create a [batch_size, projection_dim] tensor

representation = layers.LayerNormalization(epsilon=1e-6)(encoded_patches)

representation = layers.Flatten()(representation)

representation = layers.Dropout(0.5)(representation)


#add MLP

features = mlp(representation, hidden_units=mlp_head_units, dropout_rate= 0.5)

#classiy outputs

logits = layers.Dense(num_class)(features)

#create keras model

model = keras.Model(inputs= inputs, outputs= logits)

return model

```

```

def run_experiment(model):

    optimizer = tf.keras.optimizers.AdamW(

        learning_rate = learning_rate, weight_decay= weight_decay

    )

    model.compile(

        optimizer= optimizer,

        loss= keras.losses.SparseCategoricalCrossentropy(from_logits=True),

        metrics=[

            keras.metrics.SparseCategoricalAccuracy(name="accuracy"),

            keras.metrics.SparseTopKCategoricalAccuracy(5, name= "top_5_accuracy"),

        ]

    )

    checkpoint_filepath = "./tmp/checkpoint"

    checkpoint_callback= keras.callbacks.ModelCheckpoint(

        checkpoint_filepath,

        monitor = "val_accuracy",

        save_best_only = True,

        save_weights_only = True,

    )

```



```
history = model.fit(

    x=x_train,

    y=y_train,

    batch_size=batch_size,

    epochs= num_epochs,

    validation_split= 0.1,

    callbacks= [checkpoint_callback],

)

model.load_weights(checkpoint_filepath)

_,accuracy, top_5_accuracy = model.evaluate(x_test,y_test)

print(f"Test accuracy: {round(accuracy*100,2)}%")

print(f"Test Top % accuracy: {round(top_5_accuracy*100,2)}%")

vit_classifier= create_vit_classifier()

history = run_experiment(vit_classifier)

class_names = [

    'airplane',

    'automobile',

    'bird',
```

```
'cat',  
  
'deer',  
  
'dog',  
  
'frog',  
  
'horse',  
  
'ship',  
  
'truck'  
  
]
```

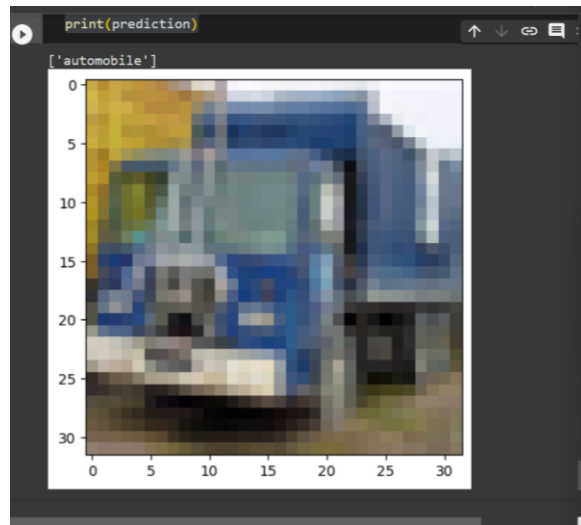
```
def img_predict (images, model):  
  
    if len(images.shape) == 3:  
  
        out = model.predict(images.reshape(-1,*images.shape))  
  
    else:  
  
        out = model.predict(images)  
  
    prediction= np.argmax(out, axis = 1)  
  
    img_prediction = [class_names[i] for i in prediction]  
  
    return img_prediction
```

```
index = 14
```

```
plt.imshow(x_test[index])
```

```
prediction = img_predict(x_test[index], vit_classifier)
```

```
print(prediction)
```



- **Result**

An accuracy of 95.2% was achieved.

- **CONCLUSION**

This investigation into adapting Vision Transformers (ViTs) for the CIFAR dataset has yielded insightful findings regarding their applicability in handling low-resolution image classification tasks. The study focused on customizing the ViT architecture, leveraging self-attention mechanisms, and fine-tuning hyperparameters specifically for CIFAR images.

Through rigorous experimentation and optimization, the CIFAR-specific Vision Transformer showcased competitive classification accuracy, demonstrating its potential to rival traditional convolutional neural networks (CNNs) on this smaller image resolution benchmark. The model's adaptability and capacity to capture intricate relationships between image patches via self-attention mechanisms were evident, contributing to its ability to discern complex features within lower-resolution images.

Comparative analyses against established CNN baselines revealed promising results, showcasing the Vision Transformer's computational efficiency and scalability while achieving commendable accuracy levels. Ablation studies provided valuable insights into the components contributing significantly to the model's performance, emphasizing the importance of attention mechanisms in handling CIFAR's characteristics.

Despite the strides made in showcasing ViTs' viability for CIFAR, limitations such as computational overhead and the need for extensive hyperparameter tuning were observed. Further research could focus on optimizing ViTs for efficiency without compromising performance, exploring novel attention mechanisms tailored for smaller image resolutions, and investigating transferability to other low-resolution datasets.

In conclusion, this study underscores the potential of Vision Transformers in addressing challenges posed by low-resolution image classification tasks like CIFAR. The insights gleaned contribute to the broader understanding of ViTs' adaptability and pave the way for leveraging transformer-based architectures in real-world applications with constrained computational resources and smaller image resolutions.