

# CMPE 597 Sp. Tp. Deep Learning - Term Project

## *DINO Architecture - CIFAR100 Dataset*

*Anil Turgut - 2022702072*

*Selahattin Seha Cirit - 2023705006*

In this Jupyter Notebook, a study was carried out by finetuning the classification model with the image embeddings obtained using the pre-trained self-supervised learning model that we proposed in the project. Image embeddings in this notebook were created using the **DINOv2vits14** model developed by MetaResearch. As the output of this model, there is an embedding list output with 384 dimensions for each image. The results of these embeddings resulting from the pretext task were analyzed using *Single Layer MLP* and *SupportVectorClassifier (SVC)* models.

Moreover, this notebook includes the analysis of *self-supervised learning* (specifically DINO) with **CIFAR10** dataset. In the following tasks, we are introduce our work in detail. Let's move on.

### Task 1: Importing Libraries

```
In [22]: # Importing necessary libraries
import copy
import cv2
import glob
import json
import keras
import matplotlib.pyplot as plt
import numpy as np
import os
import torch
import torch.nn as nn
import torch.optim as optim
import torch.utils.data as data
import torchvision.transforms as T
import torchvision.datasets as datasets
import tqdm
from tqdm.notebook import tqdm as tqdm_note
import zipfile
from copy import deepcopy
from PIL import Image
import pandas as pd
from sklearn.metrics import accuracy_score, classification_report, f1_score, precision_s
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import OneHotEncoder
import socket
from tensorflow.keras import datasets as tfdatasets, layers, models
```

The libraries to be used have been imported as in the cell above. *TensorFlow Keras, Torch* and *ScikitLearn* libraries were used when establishing classification models. Other libraries are also used for different purposes.

## Task 2.1: Extracting CIFAR100 Dataset

```
In [25]: class_names = pd.read_csv("cifar100_classnames.csv")

if 'Unnamed: 0' in class_names.columns:
    class_names.drop(columns=['Unnamed: 0'], inplace=True)

(train_images, train_labels), (test_images, test_labels) = tfdatasets.cifar100.load_data
class_names.head()
```

```
Out[25]:
```

	ClassName
0	apple
1	aquarium_fish
2	baby
3	bear
4	beaver

The *CIFAR-100* dataset consists of 60000 32x32 colour images in 100 subclasses, with 600 images per class. There are 50000 training images and 10000 test images.

This dataset is just like the CIFAR-10, except it has 100 classes containing 600 images each. There are 500 training images and 100 testing images per class. The 100 classes in the CIFAR-100 are grouped into 20 superclasses. Each image comes with a "fine" label (the class to which it belongs) and a "coarse" label (the superclass to which it belongs).

Reference: <https://www.cs.toronto.edu/~kriz/cifar.html>

*class\_names* list is read from csv file. Each class has a label which is its index.

## Task 2.2: Pretext Task - Preparing DINO model

```
In [167... torch.manual_seed(42)
dinov2_vits14 = torch.hub.load("facebookresearch/dinov2", "dinov2_vits14")
```

Using cache found in C:\Users\anil.turgut/.cache/torch/hub/facebookresearch\_dinov2\_main

DINO, a new self supervised system by Facebook AI, is able to learn incredible representations from unlabeled data. It was introduced in their paper "Emerging Properties in Self-Supervised Vision Transformers". This architecture mainly uses the vision transformers (ViT) to extract information/representation from unlabeled dataset. In other words, A Student ViT learns to predict global features in an image from local patches supervised by the cross entropy loss from a momentum Teacher ViT's embeddings while doing centering and sharpening to prevent mode collapse.

DINO is currently the state of art architecture in Self-Supervised Image Classification - ImageNet1K. That's why we have used this architecture as the backbone of our project. There are 4 kinds of pretrained backbone models of this architecture provided by *facebookresearch* and we have used the **ViT-S/14 distilled** model with almost **21M** parameters.

Using **Torch.hub**, we have loaded the model to our working environment to compute image embeddings as pretext task.

## References:

- <https://github.com/facebookresearch/dinov2>
- <https://arxiv.org/abs/2104.14294>
- <https://towardsdatascience.com/dino-emerging-properties-in-self-supervised-vision-transformers-summary-ab91df82cc3c>

```
In [119]: device = torch.device('cuda' if torch.cuda.is_available() else "cpu") # sometimes not en
#device = "cpu"
device
```

```
Out[119]: device(type='cuda')
```

Device is mainly selected as Cuda due to its performance. However, our local machines have not powerful gpu (*NVIDIA GeForce MX330 2GB*), sometimes *CPU* is selected intentionally.

```
In [29]: dinov2_vits14.to(device)
```

```
Out[29]: DinoVisionTransformer(
  (patch_embed): PatchEmbed(
    (proj): Conv2d(3, 384, kernel_size=(14, 14), stride=(14, 14))
    (norm): Identity()
  )
  (blocks): ModuleList(
    (0-11): 12 x NestedTensorBlock(
      (norm1): LayerNorm((384,), eps=1e-06, elementwise_affine=True)
      (attn): MemEffAttention(
        (qkv): Linear(in_features=384, out_features=1152, bias=True)
        (attn_drop): Dropout(p=0.0, inplace=False)
        (proj): Linear(in_features=384, out_features=384, bias=True)
        (proj_drop): Dropout(p=0.0, inplace=False)
      )
      (ls1): LayerScale()
      (drop_path1): Identity()
      (norm2): LayerNorm((384,), eps=1e-06, elementwise_affine=True)
      (mlp): Mlp(
        (fc1): Linear(in_features=384, out_features=1536, bias=True)
        (act): GELU(approximate='none')
        (fc2): Linear(in_features=1536, out_features=384, bias=True)
        (drop): Dropout(p=0.0, inplace=False)
      )
      (ls2): LayerScale()
      (drop_path2): Identity()
    )
  )
  (norm): LayerNorm((384,), eps=1e-06, elementwise_affine=True)
  (head): Identity()
)
```

We moved our DINO model to the device to ensure that all computations involving the model parameters and inputs will be performed on the specified device. Also, all the computed images' embeddings will have **384** dimensions with this DINO model.

```
In [120]: transform_image = T.Compose([T.ToTensor(), T.Resize(244), T.CenterCrop(224), T.Normalize

def load_image(img: str) -> torch.Tensor:
    """
    Load an image and return a tensor that can be used as an input to DINOv2.
    """
    img = Image.open(img)
```

```

transformed_img = transform_image(img)[:3].unsqueeze(0)

return transformed_img

def plot_image(tensor_image):
    # Convert tensor to NumPy array and transpose dimensions
    numpy_img = tensor_image.squeeze().permute(1, 2, 0).cpu().numpy()

    # Plot the image
    plt.imshow(numpy_img)
    plt.axis('off')
    plt.show()

def compute_embeddings(images: list) -> list:
    """
    Create an index that contains all of the images in the specified list of files.
    """
    all_embeddings = []

    with torch.no_grad():

        for image in tqdm_note(images):

            image = transform_image(image)[:3].unsqueeze(0)
            embeddings = dinov2_vits14(image.to(device))
            all_embeddings.append(np.array(embeddings[0].cpu().numpy()).reshape(1, -1).tolist)

    return all_embeddings

```

Cell above have 3 functions to help while transforming image to the shape that dino can understand. *load\_image* and *plot\_image* functions basically loads the *.jpg* or *.png* format images, transforms to the tensor and plot the loaded image respectively.

*compute\_embeddings* function is defined to compute image embeddings from the given image list using DINO model. In our project, CIFAR training and test image datasets will be executed by this model and output embeddings will be an input for our downstream task -*Classification*-.

## 2.3 Computing/Loading Embeddings

Using *compute\_embeddings* function above, we will compute the each image embeddings in train/test dataset. Then, we are going to store this embeddings as JSON file not to recalculate again and again.

```

In [ ]: embeddings = compute_embeddings(train_images)

with open("CIFAR100Embeddings/_dinocifar100_all_embeddings.json", "w") as f:
    f.write(json.dumps(embeddings))

test_embeddings = compute_embeddings(test_images)

with open("CIFAR100Embeddings/_dinocifar100_all_embeddings_test.json", "w") as f:
    f.write(json.dumps(test_embeddings))

```

```

In [31]: with open('CIFAR100Embeddings/_dinocifar100_all_embeddings.json') as f:
    embeddings = json.load(f)

```

```
with open('CIFAR100Embeddings/_dinocifar100_all_embeddings_test.json') as f:
    test_embeddings = json.load(f)
```

**embeddings** are the computed embeddings for the CIFAR10 training images (50000 records) and **test\_embeddings** are the computed embeddings for the CIFAR10 training images (10000 records). We do not need to compute the embeddings for label (y values) since the pretext task is unlabeled.

```
In [32]: copied_training_embeddings = embeddings.copy()
copied_training_embeddings = np.array(copied_training_embeddings).reshape(-1, 384)
copied_test_embeddings = test_embeddings.copy()
copied_test_embeddings = np.array(copied_test_embeddings).reshape(-1, 384)
```

Copied embeddings will be used in the SVC model in below.

## 2.4 Preparing Dataset for Training

```
In [33]: X_train = embeddings
X_train = np.array(X_train).reshape(-1, 384)
y_train = train_labels
ohe = OneHotEncoder(handle_unknown='ignore', sparse_output=False).fit(y_train)
y_train = ohe.transform(y_train)
X_test = test_embeddings
X_test = np.array(X_test).reshape(-1, 384)
y_test = test_labels
ohe = OneHotEncoder(handle_unknown='ignore', sparse_output=False).fit(y_test)
y_test = ohe.transform(y_test)
```

```
In [34]: # convert pandas DataFrame (X) and numpy array (y) into PyTorch tensors
X_train = torch.tensor(X_train, dtype=torch.float32)
y_train = torch.tensor(y_train, dtype=torch.float32)
X_test = torch.tensor(X_test, dtype=torch.float32)
y_test = torch.tensor(y_test, dtype=torch.float32)
```

```
In [35]: print("Shape of X_train:", X_train.shape)
print("Shape of y_train:", y_train.shape)
print("Shape of X_test:", X_test.shape)
print("Shape of y_test:", y_test.shape)
```

```
Shape of X_train: torch.Size([50000, 384])
Shape of y_train: torch.Size([50000, 100])
Shape of X_test: torch.Size([10000, 384])
Shape of y_test: torch.Size([10000, 100])
```

In preparing dataset section, we reshaped the computed embeddings of images as (-1,384). Also labels are redefined as one-hot-encoded list. Thus, each label record consists of 100 dimensions and including one 1 rest is 0.

All the dataset items are converted to *tensor* to be used in perceptron model.

## 3.1 Downstream Task - Single Layer Classification Perceptron Model

```
In [36]: class DinoVisionTransformerClassifier(nn.Module):
    def __init__(self):
        super().__init__()
        self.hidden = nn.Linear(384, 256)
        self.act = nn.ReLU()
```

```

        self.output = nn.Linear(256, 100)

    def forward(self, x):
        x = self.act(self.hidden(x))
        x = self.output(x)
        return x

# loss metric and optimizer
model = DinoVisionTransformerClassifier()
loss_fn = nn.CrossEntropyLoss()
optimizer = optim.Adam(model.parameters(), lr=0.001)

# prepare model and training parameters
n_epochs = 25
batch_size = 64
batches_per_epoch = len(X_train) // batch_size

```

Classifier is designed to classify image labels using embeddings using Single Hidden Layer perceptron model. 384 image embeddings for each image will be corresponded by the hidden layer and 256 hidden units in that layer. Output will be 100 class values through FC layer.

Loss is selected as CrossEntropyLoss since the problem is multi-class classification and optimizer is selected as *Adam* optimizer to reduce the effects of hyperparameters and including momentum and SGD mechanism.

After analyzing the model training process, epoch numbers and batch sizes selected accordingly.

```

In [37]: best_acc = - np.inf    # init to negative infinity
         best_weights = None
         train_loss_hist = []
         train_acc_hist = []
         test_loss_hist = []
         test_acc_hist = []

```

```

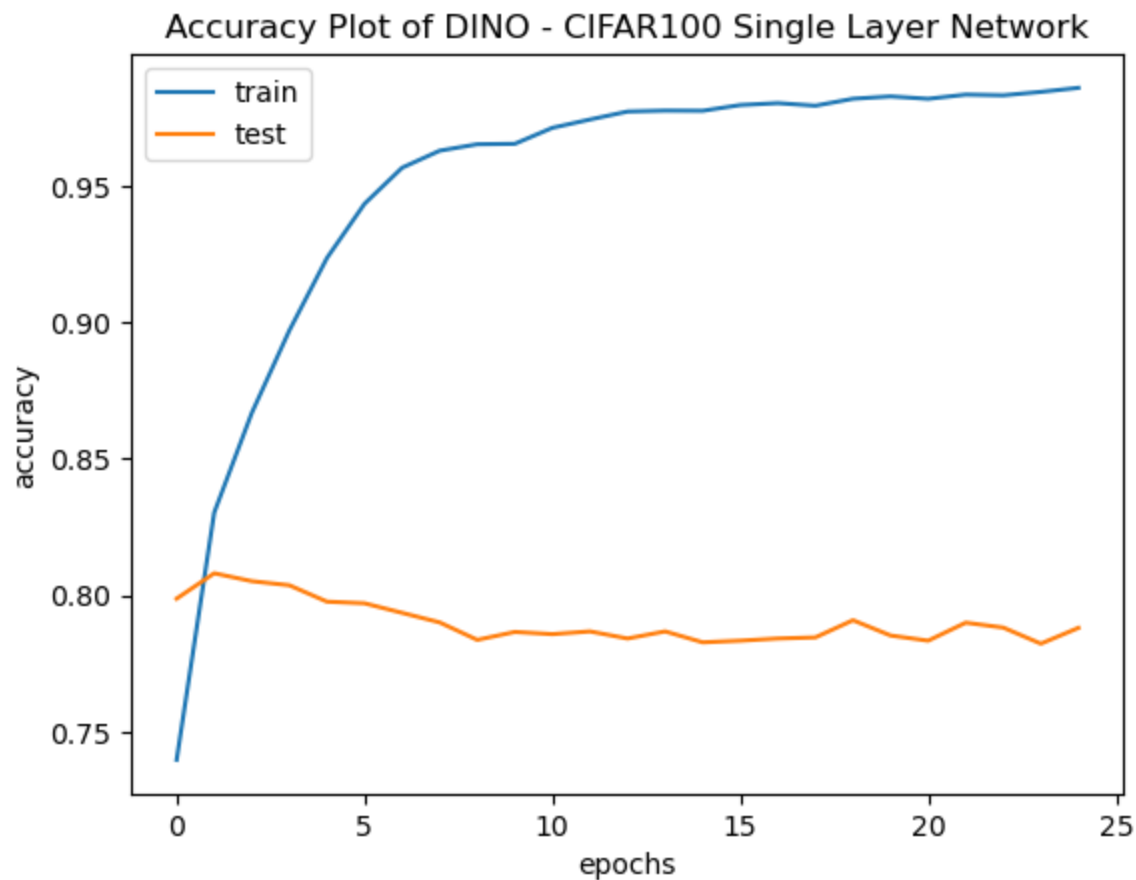
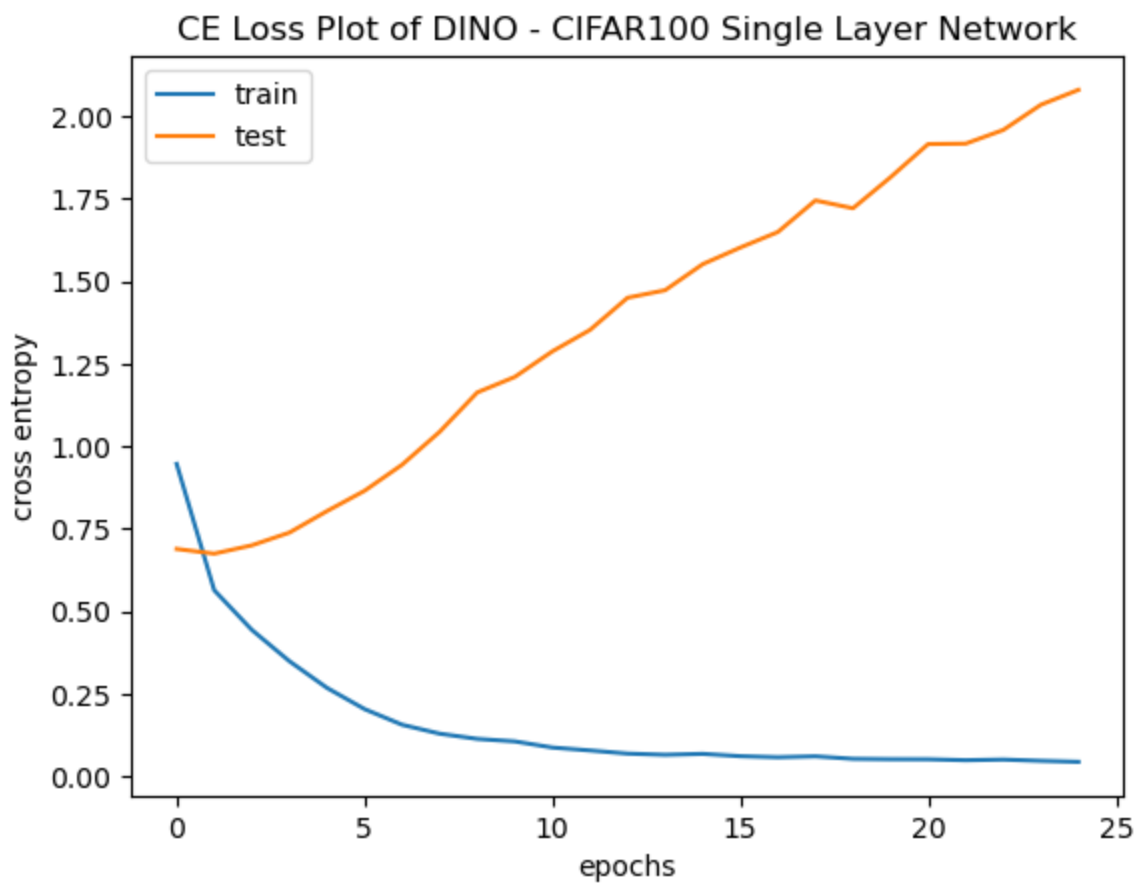
In [38]: # training loop
         for epoch in range(n_epochs):
             epoch_loss = []
             epoch_acc = []
             # set model in training mode and run through each batch
             model.train()
             with tqdm.trange(batches_per_epoch, unit="batch", mininterval=0) as bar:
                 bar.set_description(f"Epoch {epoch}")
                 for i in bar:
                     # take a batch
                     start = i * batch_size
                     X_batch = X_train[start:start+batch_size]
                     y_batch = y_train[start:start+batch_size]
                     # forward pass
                     y_pred = model(X_batch)
                     loss = loss_fn(y_pred, y_batch)
                     # backward pass
                     optimizer.zero_grad()
                     loss.backward()
                     # update weights
                     optimizer.step()
                     # compute and store metrics
                     acc = (torch.argmax(y_pred, 1) == torch.argmax(y_batch, 1)).float().mean()
                     epoch_loss.append(float(loss))
                     epoch_acc.append(float(acc))
                     bar.set_postfix(
                         loss=float(loss),
                         acc=float(acc)

```









Accuracy in training is continuously improving as expected, whereas test accuracy kind of oscillates but improved. Double depth phenomenon can be seen in second plot. But CE Loss graph shows us that improving training performance might result with decrease in test performance as well.

## 3.2 Downstream Task - TensorFlow Trainer Classification

In this part, we have used the *TensorFlow Trainer* module to train classification network. Architecture is similar to the previous model.

```
In [41]: from tensorflow.keras import models, layers

tf_model = models.Sequential()

tf_model.add(layers.Dense(256, activation='relu', input_shape=(384,)))

tf_model.add(layers.Dense(100, activation='softmax')) # Assuming 10 classes, and using

tf_model.summary()
```

Model: "sequential\_1"

Layer (type)	Output Shape	Param #
dense_2 (Dense)	(None, 256)	98,560
dense_3 (Dense)	(None, 100)	25,700

Total params: 124,260 (485.39 KB)

Trainable params: 124,260 (485.39 KB)

Non-trainable params: 0 (0.00 B)

```
In [42]: tf_model.compile(optimizer='adam',
                        loss=keras.losses.SparseCategoricalCrossentropy(from_logits=True),
                        metrics=['accuracy'])
```

Optimizer and loss function is selected as *adam* and *SparseCategoricalCrossEntropy* to see whether there is any improvement in changing loss function.

```
In [43]: history = tf_model.fit(copied_training_embeddings, train_labels, epochs=25,
                               validation_data = (copied_test_embeddings, test_labels))
```

Epoch 1/25

```
C:\Users\anil.turgut\AppData\Local\anaconda3\Lib\site-packages\keras\src\backend\tensorflow\nn.py:599: UserWarning: "`sparse_categorical_crossentropy` received `from_logits=True`, but the `output` argument was produced by a Softmax activation and thus does not represent logits. Was this intended?
```

```
    output, from_logits = _get_logits(
1563/1563 ————— 3s 1ms/step - accuracy: 0.6180 - loss: 1.58
43 - val_accuracy: 0.7810 - val_loss: 0.7621
```

Epoch 2/25

```
1563/1563 ————— 2s 1ms/step - accuracy: 0.8309 - loss: 0.55
83 - val_accuracy: 0.7854 - val_loss: 0.7604
```

Epoch 3/25

```
1563/1563 ————— 2s 1ms/step - accuracy: 0.8716 - loss: 0.41
36 - val_accuracy: 0.7888 - val_loss: 0.7964
```

Epoch 4/25

```
1563/1563 ————— 2s 1ms/step - accuracy: 0.8996 - loss: 0.31
21 - val_accuracy: 0.7902 - val_loss: 0.8450
```

Epoch 5/25

```
1563/1563 ————— 2s 1ms/step - accuracy: 0.9228 - loss: 0.24
40 - val_accuracy: 0.7861 - val_loss: 0.9435
```

Epoch 6/25

```
1563/1563 ————— 2s 1ms/step - accuracy: 0.9390 - loss: 0.18
57 - val_accuracy: 0.7877 - val_loss: 1.0087
```

Epoch 7/25

```
1563/1563 ————— 2s 1ms/step - accuracy: 0.9505 - loss: 0.15
```

```

27 - val_accuracy: 0.7840 - val_loss: 1.1090
Epoch 8/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9577 - loss: 0.12
80 - val_accuracy: 0.7896 - val_loss: 1.1405
Epoch 9/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9612 - loss: 0.11
51 - val_accuracy: 0.7834 - val_loss: 1.2787
Epoch 10/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9663 - loss: 0.09
77 - val_accuracy: 0.7817 - val_loss: 1.4193
Epoch 11/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9712 - loss: 0.08
59 - val_accuracy: 0.7824 - val_loss: 1.4587
Epoch 12/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9709 - loss: 0.08
52 - val_accuracy: 0.7848 - val_loss: 1.5297
Epoch 13/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9769 - loss: 0.06
92 - val_accuracy: 0.7780 - val_loss: 1.6645
Epoch 14/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9742 - loss: 0.07
61 - val_accuracy: 0.7747 - val_loss: 1.7504
Epoch 15/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9779 - loss: 0.06
47 - val_accuracy: 0.7828 - val_loss: 1.7701
Epoch 16/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9795 - loss: 0.05
92 - val_accuracy: 0.7791 - val_loss: 1.9190
Epoch 17/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9783 - loss: 0.06
56 - val_accuracy: 0.7792 - val_loss: 1.9297
Epoch 18/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9807 - loss: 0.05
95 - val_accuracy: 0.7828 - val_loss: 2.0126
Epoch 19/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9801 - loss: 0.06
14 - val_accuracy: 0.7834 - val_loss: 2.0507
Epoch 20/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9813 - loss: 0.05
42 - val_accuracy: 0.7826 - val_loss: 2.1352
Epoch 21/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9826 - loss: 0.05
40 - val_accuracy: 0.7839 - val_loss: 2.2220
Epoch 22/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9823 - loss: 0.05
41 - val_accuracy: 0.7825 - val_loss: 2.2998
Epoch 23/25
1563/1563 ————— 2s 2ms/step - accuracy: 0.9834 - loss: 0.05
03 - val_accuracy: 0.7808 - val_loss: 2.3085
Epoch 24/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9849 - loss: 0.04
71 - val_accuracy: 0.7851 - val_loss: 2.3866
Epoch 25/25
1563/1563 ————— 2s 1ms/step - accuracy: 0.9849 - loss: 0.05
01 - val_accuracy: 0.7829 - val_loss: 2.4561

```

```
In [44]: y_train.shape
```

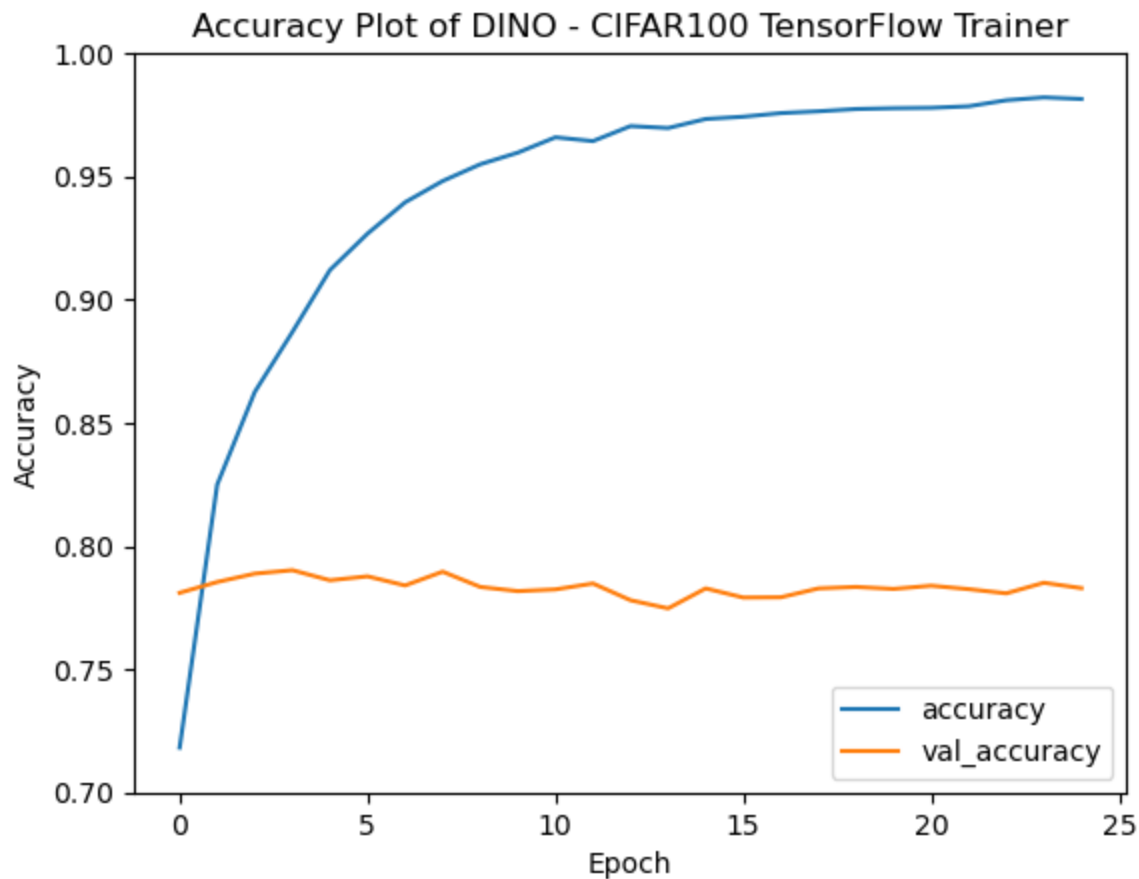
```
Out[44]: torch.Size([50000, 100])
```

```
In [58]: plt.plot(history.history['accuracy'], label='accuracy')
plt.plot(history.history['val_accuracy'], label = 'val_accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.title("Accuracy Plot of DINO - CIFAR100 TensorFlow Trainer")
```

```
plt.ylim([0.7, 1])
plt.legend(loc='lower right')
```

Out[58]:

<matplotlib.legend.Legend at 0x19c0bddae50>



Accuracy plot shows that the results of *TF Trainer* and *Single Layer Perceptron* models are similar and satisfied. We can cover almost 80% of the validation image dataset.

Let's finally analyze the *Support Vector Classifier (SVC)* model to analyze the results.

## 3.3 Downstream Task - Classification Using Support Vector Classifier (SVC)

*SVC with gamma parameter* is created as the downstream classification model

In [48]:

```
from sklearn import svm

clf = svm.SVC(gamma='scale')

print(len(embeddings))

clf.fit(np.array(embeddings).reshape(-1, 384), train_labels)
```

50000

```
C:\Users\anil.turgut\AppData\Local\anaconda3\Lib\site-packages\sklearn\utils\validation.py:1184: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().
  y = column_or_1d(y, warn=True)
```

Out[48]:

▼ SVC  
SVC()

Training is completed between training image embeddings and training labels (without one hot encoding)

```
In [49]: # train metrics
y_pred_train_list = []

for embed in tqdm_note(embeddings):

    y_pred = clf.predict(np.array(embed).reshape(1, -1))
    y_pred_train_list.append(y_pred[0])

0%|          | 0/50000 [00:00<?, ?it/s]
```

```
In [50]: y_pred_test_list = []

for embed in tqdm_note(test_embeddings):

    y_pred = clf.predict(np.array(embed).reshape(1, -1))
    y_pred_test_list.append(y_pred[0])

0%|          | 0/10000 [00:00<?, ?it/s]
```

```
In [54]: with open("CIFAR100Predictions/dinocifar100trainingpreds.json", "w") as f:
        f.write(json.dumps(np.array(y_pred_train_list).tolist()))

with open("CIFAR100Predictions/dinocifar100testpreds.json", "w") as f:
    f.write(json.dumps(np.array(y_pred_test_list).tolist()))
```

```
In [144... with open('CIFAR100Predictions/dinocifar100trainingpreds.json') as f:
    y_pred_train_list = json.load(f)

with open('CIFAR100Predictions/dinocifar100testpreds.json') as f:
    y_pred_test_list = json.load(f)
```

```
In [55]: # Calculate evaluation metrics on the training data
accuracy_train = accuracy_score(train_labels, y_pred_train_list)
precision_train = precision_score(train_labels, y_pred_train_list, average='weighted')
recall_train = recall_score(train_labels, y_pred_train_list, average='weighted')
f1_train = f1_score(train_labels, y_pred_train_list, average='weighted')

print("Training Accuracy:", round(accuracy_train,6))
print("Training Precision:", round(precision_train,6))
print("Training Recall:", round(recall_train,6))
print("Training F1-score:", round(f1_train,6))

Training Accuracy: 0.942
Training Precision: 0.943095
Training Recall: 0.942
Training F1-score: 0.942222
```

```
In [56]: # Calculate evaluation metrics on test data
accuracy = accuracy_score(test_labels, y_pred_test_list)
precision = precision_score(test_labels, y_pred_test_list, average='weighted')
recall = recall_score(test_labels, y_pred_test_list, average='weighted')
f1 = f1_score(test_labels, y_pred_test_list, average='weighted')

print("Test Accuracy:", round(accuracy,4))
print("Test Precision:", round(precision,4))
print("Test Recall:", round(recall,4))
print("Test F1-score:", round(f1,4))

Test Accuracy: 0.8316
Test Precision: 0.836
Test Recall: 0.8316
Test F1-score: 0.8325
```

```
In [59]: train_metrics = [accuracy_train, precision_train, recall_train, f1_train]
test_metrics = [accuracy, precision, recall, f1]
metrics_labels = ['Accuracy', 'Precision', 'Recall', 'F1 Score']
num_metrics = len(metrics_labels)

index = np.arange(num_metrics)

bar_width = 0.35

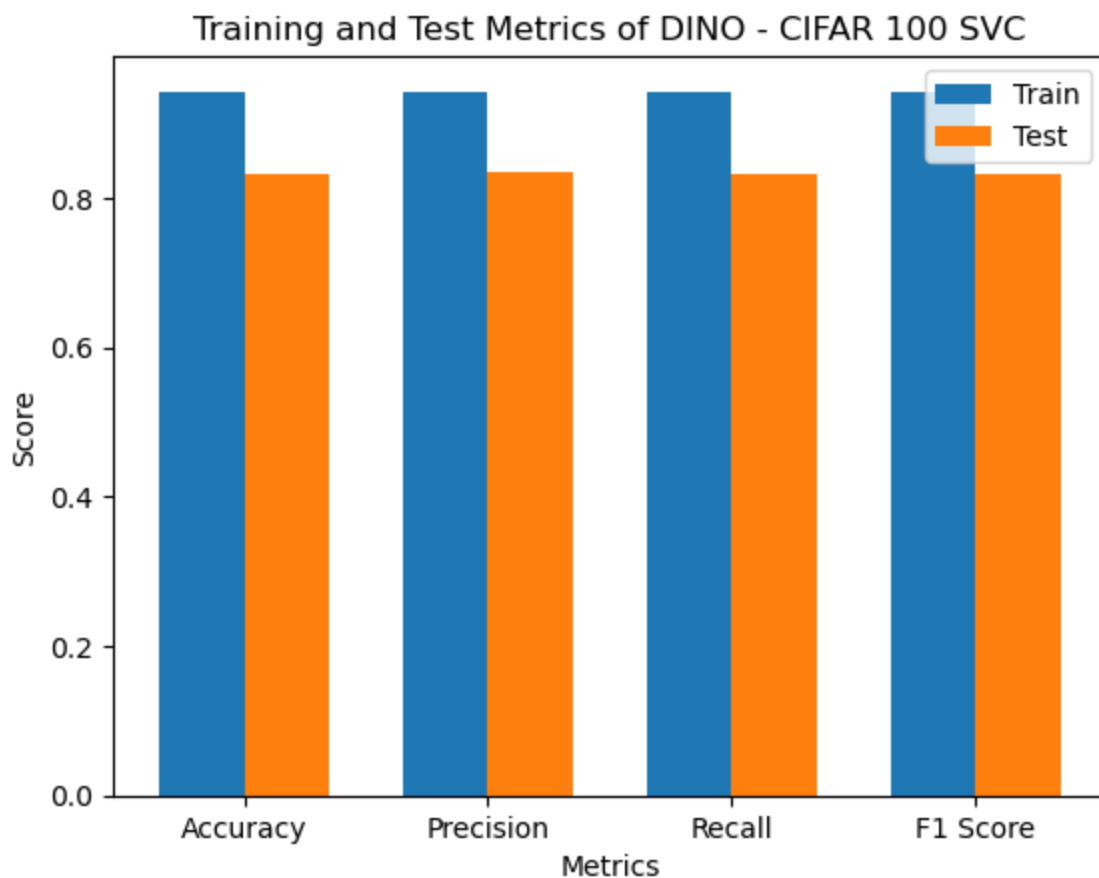
fig, ax = plt.subplots()

trainBars = ax.bar(index, train_metrics, bar_width, label='Train')

testBars = ax.bar(index + bar_width, test_metrics, bar_width, label='Test')

ax.set_xlabel('Metrics')
ax.set_ylabel('Score')
ax.set_title('Training and Test Metrics of DINO - CIFAR 100 SVC')
ax.set_xticks(index + bar_width / 2)
ax.set_xticklabels(metrics_labels)
ax.legend()

plt.show()
```



SVC has performed very well with *CIFAR100* embeddings, it has better performance than the neural network classifiers. Test accuracy and other metrics are around 84%.

## 4.1 Testing Classifier Model Performance on Image from different Dataset

In this part, we will be testing our three downstream models with different dataset(different than CIFAR)

```
In [121... input_file = "C:\\Users\\anil.turgut\\Desktop\\CMPE597\\Project\\Code\\ExampleTestImages\\
new_image = load_image(input_file)

plot_image(new_image)
```

Clipping input data to the valid range for imshow with RGB data ([0..1] for floats or [0..255] for integers).



Image is uploaded using *load\_image* function

```
In [124... dinov2_vits14.to(device)
```

```
Out[124]: DinoVisionTransformer(
  (patch_embed): PatchEmbed(
    (proj): Conv2d(3, 384, kernel_size=(14, 14), stride=(14, 14))
    (norm): Identity()
  )
  (blocks): ModuleList(
    (0-11): 12 x NestedTensorBlock(
      (norm1): LayerNorm((384,), eps=1e-06, elementwise_affine=True)
      (attn): MemEffAttention(
        (qkv): Linear(in_features=384, out_features=1152, bias=True)
        (attn_drop): Dropout(p=0.0, inplace=False)
        (proj): Linear(in_features=384, out_features=384, bias=True)
        (proj_drop): Dropout(p=0.0, inplace=False)
      )
      (ls1): LayerScale()
      (drop_path1): Identity()
      (norm2): LayerNorm((384,), eps=1e-06, elementwise_affine=True)
      (mlp): Mlp(
        (fc1): Linear(in_features=384, out_features=1536, bias=True)
        (act): GELU(approximate='none')
        (fc2): Linear(in_features=1536, out_features=384, bias=True)
        (drop): Dropout(p=0.0, inplace=False)
      )
      (ls2): LayerScale()
      (drop_path2): Identity()
    )
  )
  (norm): LayerNorm((384,), eps=1e-06, elementwise_affine=True)
```

```
(head): Identity()
```

## 4.1.1 - SVC

```
In [68]: with torch.no_grad():
         embedding = dinov2_vits14(new_image.to(device))

         prediction = clf.predict(np.array(embedding[0].cpu()).reshape(1, -1))

         print("Predicted class: " + class_names['ClassName'].iloc[prediction[0]])
```

Predicted class: bridge

## 4.1.2 - Single Layer Torch Network

```
In [164]: with torch.no_grad():

         model.eval()
         prediction = model(embedding)
         prediction = prediction.cpu().detach().numpy()
         #max_value, index = torch.max(prediction, dim = 1)
         max_indices = np.argsort(prediction, axis=1)[:,-2:]
         index = np.argmax(prediction, axis = 1)

         print("Predicted class: " + class_names['ClassName'].iloc[index[0]] +
               " with value " + str(prediction[0][max_indices[0][1]]))
         print("Second Predicted class: " + class_names['ClassName'].iloc[max_indices[0][0]]
               + " with value " + str(prediction[0][max_indices[0][0]]))
```

Predicted class: willow\_tree with value 1.3158058  
Second Predicted class: bridge with value 1.284671

## 4.1.3 - TF Trainer

```
In [103]: prediction = tf_model.predict(embedding)

         print("Predicted class: " + class_names['ClassName'].iloc[np.argmax(prediction,axis = 1)]
```

1/1 ————— 0s 21ms/step

Predicted class: bridge

All of 3 classifiers can can correctly classify images which are not from neither test nor training dataset

# 5. Results/Analysis

*DINO* is a self-supervised learning (SSL) model architecture which is developed by *FacebookAI* is the state of art architecture in SSL-Image Classification context. In this notebook, we have analyzed ImageNet1k fed pretrained *dino* model with the **CIFAR100** dataset. Image embeddings obtained from *DINO* model are used as an input to train 3 different classification model/networks as following: *Single Layer Network*, *TensorFlow Trainer* and *Support Vector Classifier*.

All of the 3 models performed significantly well in both *CIFAR* training and test dataset almost over **84%** accuracy. In other words, without memorizing, output embeddings are fit to these models to classify unknown images' labels.

However, SVC model is the best performed model in terms of test dataset performance metrics.



In terms of performance/memory of the algorithm analysis, *neural network classifiers* are more efficient than the SVC in terms of both speed and the memory usage.

## References

- <https://www.cs.toronto.edu/~kriz/cifar.html>
- <https://github.com/facebookresearch/dinov2>
- <https://arxiv.org/abs/2104.14294>
- <https://towardsdatascience.com/dino-emerging-properties-in-self-supervised-vision-transformers-summary-ab91df82cc3c>
- <https://arxiv.org/abs/2010.11929v2>
- <https://arxiv.org/abs/1706.03762>
- <https://paperswithcode.com/sota/self-supervised-image-classification-on>