# Computer Vision HW2 Report
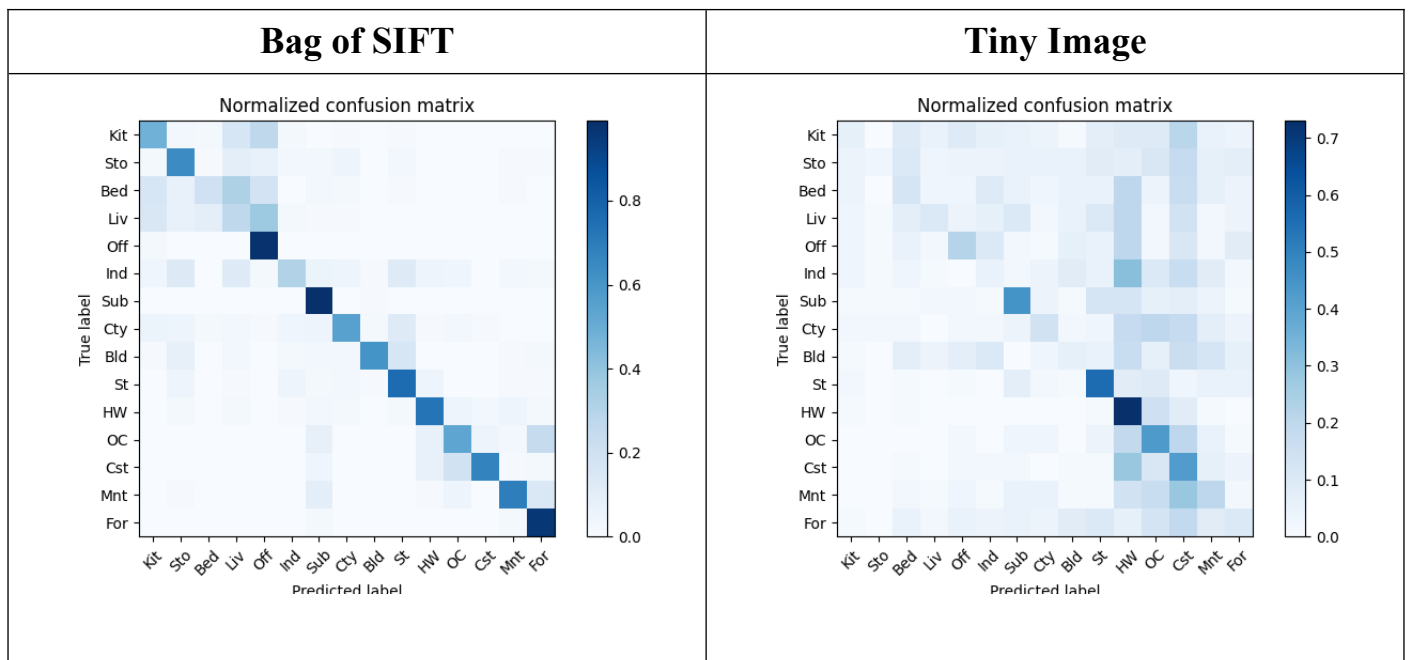
Student ID: B10901151
Name: 林祐群

## Part 1. (10%)

- **Plot confusion matrix of two settings. (i.e. Bag of sift and tiny image) (5%)**

**Ans:**

| Bag of SIFT | Tiny Image |
|:---:|:---:|
| Normalized confusion matrix | Normalized confusion matrix |

- **Compare the results/accuracy of both settings and explain the result. (5%)**

**Ans:**

| Categories | Results/Accuracy |
|:---:|:---:|
| Bag of SIFT | 0.6246666666666667 |
| Tiny Image | 0.24466666666666667 |

It's clear that the accuracy of the setting of "tiny image" is a lot lower than that of "bag of SIFT". The result is obvious because tiny image can only capture basic visual features such as color distribution and brightness. On the other hands, bag of SIFT can acquire finer visual concepts such as texture and local shape.

Follow the instructions in the provided python files, I implement the two setting with proper manner such that the accuracy of both setting are great. I've also tried different metric for the cdist, and found the setting perform suitably for the task. Further study can be done to determine whether L2 normalization, other k-means' number, or vocab_size will lead to even better results.

# Part 2. (25%)

- **Report accuracy of both models on the validation set. (2%)**

**Ans:**

|          | A (Mynet) | B (Resnet) |
|----------|-----------|------------|
| accuracy | 0.8642    | 0.9034     |

- **Print the network architecture & number of parameters of both models. What is the main difference between ResNet and other CNN architectures? (5%)**

**Ans:**

## Mynet (Number of parameters: **3841610**)

## Architecture:

MyNet(

 (features): Sequential(

  (0): Conv2d(3, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))

  (1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

  (2): SiLU(inplace=True)

  (3): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))

  (4): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

  (5): SiLU(inplace=True)

  (6): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)

  (7): Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))

  (8): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

  (9): SiLU(inplace=True)

  (10): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))

  (11): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

  (12): SiLU(inplace=True)

  (13): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)

  (14): Conv2d(128, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))

  (15): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

  (16): SiLU(inplace=True)

  (17): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))

  (18): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

  (19): SiLU(inplace=True)

  (20): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))

  (21): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

  (22): SiLU(inplace=True)

  (23): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)

 )

```
(classifier): Sequential(

  (0): Linear(in_features=4096, out_features=512, bias=True)

  (1): BatchNorm1d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

  (2): SiLU(inplace=True)

  (3): Dropout(p=0.5, inplace=False)

  (4): Linear(in_features=512, out_features=10, bias=True)

 )

)
```

# Resnet (Number of parameters: 11173962)

# Architecture:

```
ResNet18(

 (resnet): ResNet(

  (conv1): Conv2d(3, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)

  (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

  (relu): ReLU(inplace=True)

  (maxpool): Identity()

  (layer1): Sequential(

   (0): BasicBlock(

    (conv1): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)

    (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

    (relu): ReLU(inplace=True)

    (conv2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)

    (bn2): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

   )

   (1): BasicBlock(

    (conv1): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)

    (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

    (relu): ReLU(inplace=True)

    (conv2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)

    (bn2): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

   )

  )

  (layer2): Sequential(

   (0): BasicBlock(

    (conv1): Conv2d(64, 128, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), bias=False)

    (bn1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

    (relu): ReLU(inplace=True)

    (conv2): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)

    (bn2): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

    (downsample): Sequential(

     (0): Conv2d(64, 128, kernel_size=(1, 1), stride=(2, 2), bias=False)

     (1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
```

```
    )
  )
  (1): BasicBlock(
    (conv1): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
    (conv2): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  )
)
(layer3): Sequential(
  (0): BasicBlock(
    (conv1): Conv2d(128, 256, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), bias=False)
    (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
    (conv2): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (downsample): Sequential(
      (0): Conv2d(128, 256, kernel_size=(1, 1), stride=(2, 2), bias=False)
      (1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    )
  )
  (1): BasicBlock(
    (conv1): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
    (conv2): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  )
)
(layer4): Sequential(
  (0): BasicBlock(
    (conv1): Conv2d(256, 512, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), bias=False)
    (bn1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
    (conv2): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (downsample): Sequential(
      (0): Conv2d(256, 512, kernel_size=(1, 1), stride=(2, 2), bias=False)
      (1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    )
  )
```

```
    (1): BasicBlock(
      (conv1): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      (bn1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (relu): ReLU(inplace=True)
      (conv2): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      (bn2): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    )
  )
  (avgpool): AdaptiveAvgPool2d(output_size=(1, 1))
  (fc): Linear(in_features=512, out_features=10, bias=True)
 )
)
```

The main difference between ResNet and other CNN architectures is the introduction of residual connections (or skip connections), which fundamentally addresses the vanishing gradient problem in deep neural networks.
In traditional deep neural networks, as the network becomes deeper: training becomes increasingly difficult, gradient signals struggle to propagate through many layers, and performance degrades with increased depth, contrary to expectations that more layers should improve performance.
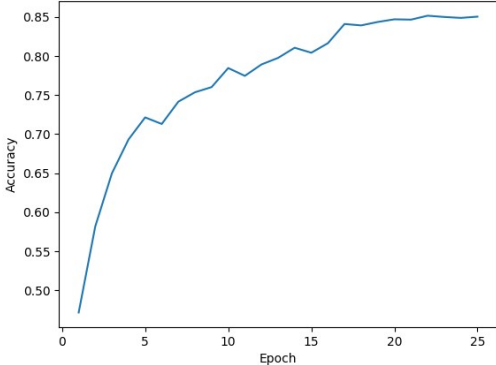
ResNet introduces residual blocks that allow direct shortcut connections that skip one or more layers. Also, the identity mapping where the input can be directly added to the output of the convolutional layers

By introducing the aforementioned method, ResNet successfully solves the critical limitation in deep neural networks, which allow for much deeper and more effective network architectures.
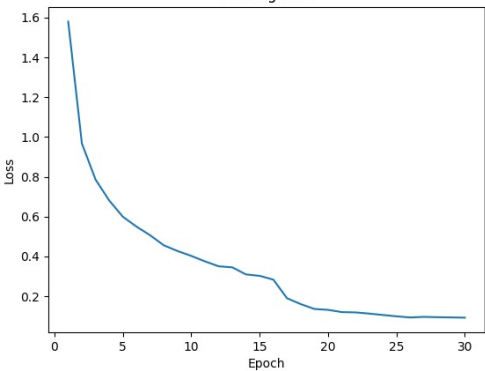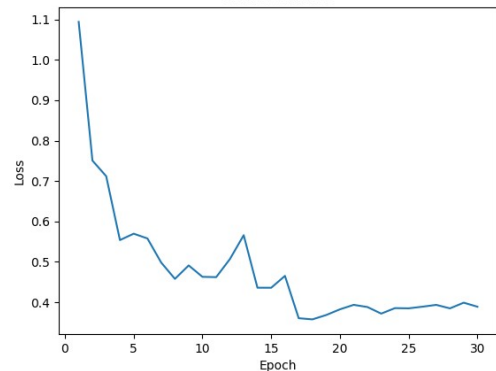
• **Plot four learning curves (loss & accuracy) of the training process (train/validation)**
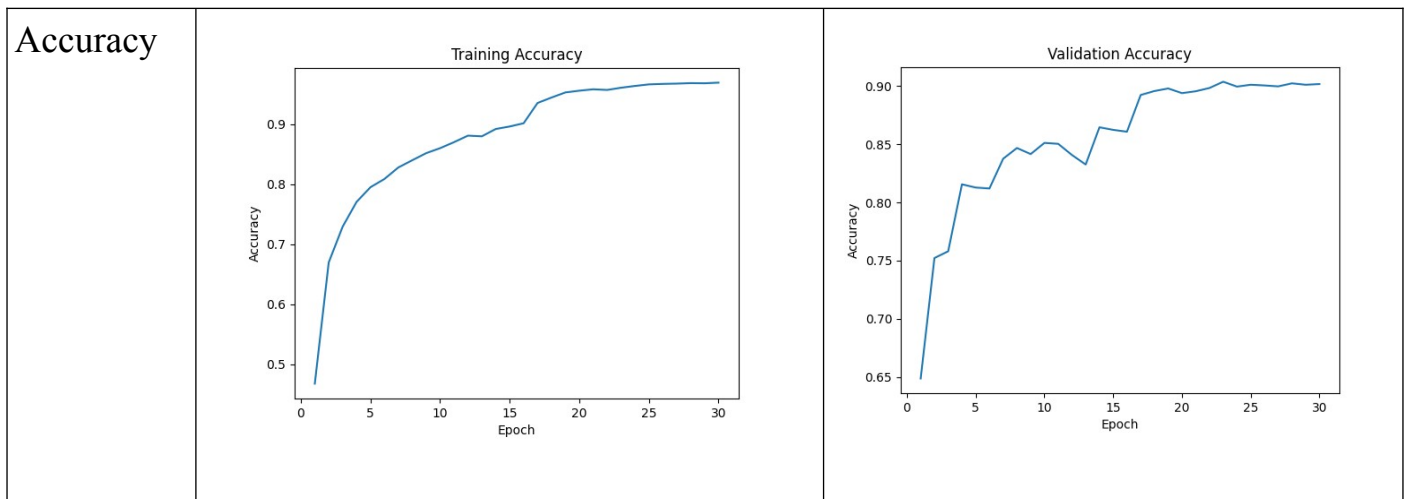**for both models. Total 8 plots. (8%)**
**Ans:**

## Mynet:

|  | Train | Validation |
|---|---|---|
| Loss |  Training Loss |  Validation Accuracy |
| Accuracy |  Training Accuracy |  Validation Accuracy |

## Resnet:

|  | Train | Validation |
|---|---|---|
| Loss |  Training Loss |  Validation Loss |

| Accuracy |  |
| --- | --- |

- **Briefly describe what method do you apply on your best model? (e.g. data augmentation, model architecture, loss function, etc) (10%)**

**Ans:**

In my best approach of ResNet, I utilize a 3×3 kernel (stride=1) and remove the initial max pooling layer so that more spatial information is preserved for the small 32×32 images. I then use aggressive data augmentation—including random cropping with padding, random horizontal flipping, and random rotation—to improve generalization. (I've also tested the result without using random cropping, but find the method a feasible way to improve accuracy, prevent the model from overfitting, and possibly enhance the robustness of the model.)

The network is trained with cross-entropy loss and Soft Gradient Decent(I've as well tested with weight decay using an Adam optimizer, but find it hard to outperform the baseline. Maybe using AdamW will be a possible way).

I compared this with my custom VGG-style MyNet (using SiLU activations and batch normalization), but the ResNet18 still yielded superior performance even though I've tested using different classifier, different activation function, and semi-supervise learning on both MyNet and ResNet. However, modified ResNet + 3 augmentation techniques (+ resize and normalize) + SGD + Cross-entropy loss + (Semi-supervise Learning) is still my best setting where the part inside the brackets are optional which gives similar results.