

# ECE 60146 - Homework 4

Ali Almualllem, aalmuall@purdue.edu

February 2025

## 1 DL Studio

### 1.1 Install DLStudio

I installed the DLStudio [2] by following the instructions on the website by running the following command

```
1 python setup.py install
```

but this has resulted in issues in my machine so I used the build commands and then imported the library like below. The system path should include the directory where DLStudio is located.

```
1 import sys
2 sys.path.append( "../" )
3 from DLStudio import *
```

### 1.2 Explore Example Networks

I ran both networks with the default parameters as noted below. Notice how the second network is much more complex which we anticipate will perform better in image classification tasks. **Net()** has only 6 layers of filters, while **Net2()** has 128 of those which would translate to more capacity to learn complex features.

#### 1.2.1 Net() default model architecture

```
1 Net(
2   (conv1): Conv2d(3, 6, kernel_size=(5, 5), stride=(1, 1))
3   (conv2): Conv2d(6, 16, kernel_size=(5, 5), stride=(1, 1))
4   (fc1): Linear(in_features=400, out_features=120, bias=True)
5   (fc2): Linear(in_features=120, out_features=84, bias=True)
6   (fc3): Linear(in_features=84, out_features=10, bias=True)
7 )
```

#### 1.2.2 Net2() default model architecture

```
1 Net2(
2   (relu): ReLU()
3   (conv1): Conv2d(3, 128, kernel_size=(5, 5), stride=(1, 1), padding=(2, 2))
4   (pool1): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
5   (conv2): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
6   (pool2): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
7   (conv3): Conv2d(128, 128, kernel_size=(2, 2), stride=(1, 1), padding=(1, 1))
8   (pool3): MaxPool2d(kernel_size=2, stride=1, padding=0, dilation=1, ceil_mode=False)
9   (fc1): Linear(in_features=8192, out_features=150, bias=True)
10  (fc2): Linear(in_features=150, out_features=100, bias=True)
11  (fc3): Linear(in_features=100, out_features=10, bias=True)
12 )
```

### 1.2.3 Results for different parameters

We run the `Net()` with different parameters and we report the results below.

#### Original Architecture

```
1 Net(  
2   (conv1): Conv2d(3, 6, kernel_size=(5, 5), stride=(1, 1))  
3   (conv2): Conv2d(6, 16, kernel_size=(5, 5), stride=(1, 1))  
4   (fc1): Linear(in_features=400, out_features=120, bias=True)  
5   (fc2): Linear(in_features=120, out_features=84, bias=True)  
6   (fc3): Linear(in_features=84, out_features=10, bias=True)  
7 )
```

**Modified Architecture** We changed the filters to 3 in the first convolution layer instead of 6, and 8 in the second convolutional layer instead of 16.

```
1 Net(  
2   (conv1): Conv2d(3, 3, kernel_size=(5, 5), stride=(1, 1))  
3   (conv2): Conv2d(3, 8, kernel_size=(5, 5), stride=(1, 1))  
4   (fc1): Linear(in_features=200, out_features=120, bias=True)  
5   (fc2): Linear(in_features=120, out_features=84, bias=True)  
6   (fc3): Linear(in_features=84, out_features=10, bias=True)  
7 )
```

Class	Unmodified Net()	Modified Net() with less filters
plane	<b>53%</b>	51%
car	<b>72%</b>	60%
bird	<b>50%</b>	38%
cat	33%	<b>47%</b>
deer	32%	<b>39%</b>
dog	<b>55%</b>	20%
frog	<b>67%</b>	49%
horse	46%	<b>48%</b>
ship	<b>77%</b>	65%
truck	33%	<b>49%</b>
overall	<b>52%</b>	47%

Table 1: The accuracies for the original `Net()` and the modified version with less number of filters. As anticipated, the original one with more filters achieved higher overall accuracy as those filters enabled more features discovery. The **boldface** values are higher for each category.

## 1.3 Run CIFAR-10 Example

We run the following to run the CIFAR-10 [3] example.

**Note:** I had to modify the provided code because it caused issues with the multi-process on my MAC computer with Apple Silicon.

The edit was basically wrapping the functions inside the `playing_with_cifar10.py` into a main function

```
1 if __name__ == '__main__':  
2     dls = DLStudio(  
3         #  
4             dataroot = "/home/kak/ImageDatasets/CIFAR-10/",  
5             dataroot = "./data/CIFAR-10/",  
6             image_size = [32,32],  
7             path_saved_model = "./saved_model",  
8             momentum = 0.9,  
9             learning_rate = 1e-3,  
10            epochs = 2,  
11            batch_size = 4,  
12            classes = ('plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse',  
13                       'ship', 'truck'),  
14            use_gpu = True,  
15            )  
16     ...  
17     ...
```

## Input

```
1 python playing_with_cifar10.py
```

## Output

```
1 ...
2 ...
3 [epoch:2/2 iter=11000 elapsed_time= 164 secs] Predicted Labels:      cat
   ship      dog      dog
4 [epoch:2/2 iter=11000 elapsed_time= 164 secs] Loss: 1.286
5
6
7 [epoch:2/2 iter=12000 elapsed_time= 169 secs] Ground Truth:      car
   car      bird      cat
8 [epoch:2/2 iter=12000 elapsed_time= 169 secs] Predicted Labels:  car
   frog      deer      cat
9 [epoch:2/2 iter=12000 elapsed_time= 169 secs] Loss: 1.276
10
11 Finished Training
```

## CIFAR10 Example - Loss with Net()

The loss as a function of iterations for the example code for training CIFAR10. The loss has monotonically decreased which suggests the network is learning successfully.

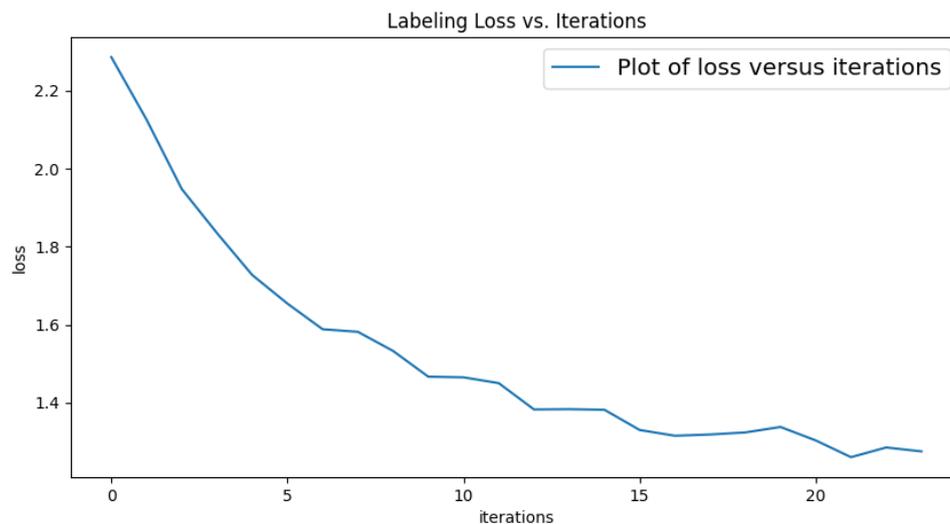


Figure 1: CIFAR10 loss graph at the end of the training.

## CIFAR10 Example - Accuracy with Net()

```
1 Prediction accuracy for plane : 53 %
2 Prediction accuracy for car : 72 %
3 Prediction accuracy for bird : 50 %
4 Prediction accuracy for cat : 33 %
5 Prediction accuracy for deer : 32 %
6 Prediction accuracy for dog : 55 %
7 Prediction accuracy for frog : 67 %
8 Prediction accuracy for horse : 46 %
9 Prediction accuracy for ship : 77 %
10 Prediction accuracy for truck : 33 %
11
12 Overall accuracy of the network on the 10000 test images: 52 %
```

## CIFAR10 Example - Confusion Matrix with Net()

Notice how the diagonal entries are higher which suggests that the network predicts the correct label with higher accuracy.

```
1
2 Displaying the confusion matrix:
3
4      plane      car      bird      cat      deer      dog      frog      horse      ship      truck
5
6 plane:  53.90   3.30   7.70   2.00   1.80   0.70   2.80   0.30  27.10   0.40
7   car:   3.00  72.90   1.10   1.30   0.50   0.40   1.40   0.00  14.70   4.70
8   bird:  6.80   1.40  50.70   7.90   7.80   9.20   8.80   1.00   5.60   0.80
9   cat:   2.70   2.50  12.30  33.40   4.70  25.90   8.90   2.00   5.60   2.00
10  deer:   3.90   1.80  24.90   7.00  32.70   7.20  11.80   6.60   3.70   0.40
11  dog:    1.60   1.50  12.00  14.40   3.80  55.60   4.50   3.50   2.90   0.20
12  frog:   0.90   1.60   8.80   7.20   3.60   6.70  67.20   0.30   2.90   0.80
13  horse:  3.40   0.80  11.90   8.80   8.00  13.90   2.00  46.10   2.80   2.30
14  ship:  11.40   3.20   2.40   1.70   1.20   1.20   0.50   0.20  77.90   0.30
15  truck:   5.00  27.70   2.80   4.20   0.40   1.10   2.30   0.70  22.60  33.20
```

## CIFAR10 Example - Loss with Net2()

The loss as a function of iterations for the example code for training CIFAR10 and with the network **Net2()**. The loss appears to be lower than the first trial above in Fig. 1 which used a simpler network.

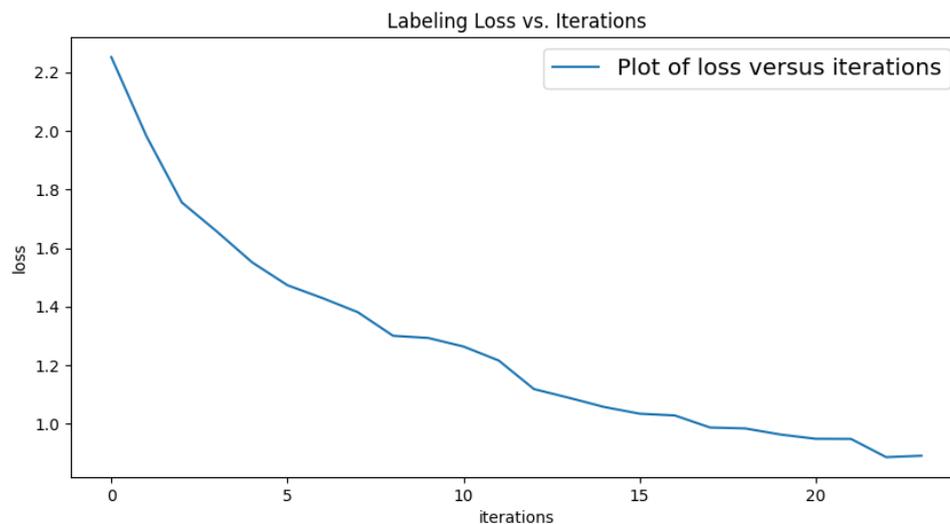


Figure 2: CIFAR10 loss graph at the end of the training with the more complex network Net2()

## CIFAR10 Example - Accuracy with Net2()

Notice how the overall accuracy is **69%** while the previous network only achieved **52%**. A substantial improvement.

```
1 Prediction accuracy for plane : 69 %
2 Prediction accuracy for car : 81 %
3 Prediction accuracy for bird : 57 %
4 Prediction accuracy for cat : 54 %
5 Prediction accuracy for deer : 61 %
6 Prediction accuracy for dog : 33 %
7 Prediction accuracy for frog : 89 %
8 Prediction accuracy for horse : 68 %
9 Prediction accuracy for ship : 75 %
10 Prediction accuracy for truck : 80 %
11
12
13
```

```
14 Overall accuracy of the network on the 10000 test images: 67 %
```

## CIFAR10 Example - Confusion Matrix with Net2()

Notice how the diagonal entries are higher which suggests that the network predicts the correct label with higher accuracy.

But also compared with the previous experiment, the diagonal entries which correspond to accuracies are much higher.

```
1 Displaying the confusion matrix:
2
3
4
5
6
7
8
9
10
11
12
13
14
```

	plane	car	bird	cat	deer	dog	frog	horse	ship	truck
plane:	69.40	2.20	12.00	2.50	2.30	0.10	2.10	1.30	4.70	3.40
car:	1.50	81.10	1.20	0.10	0.40	0.10	2.20	0.10	1.00	12.30
bird:	4.10	0.50	57.10	7.10	10.00	1.80	14.80	2.40	0.60	1.60
cat:	0.90	0.90	7.30	54.10	5.90	4.10	20.60	2.90	1.10	2.20
deer:	1.30	0.20	5.70	7.20	61.80	0.70	13.80	8.00	0.90	0.40
dog:	0.90	0.70	7.40	33.20	6.10	33.20	13.00	3.30	1.10	1.10
frog:	0.30	0.40	3.90	2.60	2.70	0.00	89.20	0.30	0.20	0.40
horse:	0.70	0.60	4.40	7.30	10.90	1.60	3.00	68.80	0.40	2.30
ship:	7.60	5.10	3.10	1.30	0.90	0.10	1.30	1.00	75.90	3.70
truck:	3.50	8.30	0.60	1.20	0.70	0.00	1.90	1.10	1.80	80.90

## 1.4 Understanding Network Architecture

As stated in the instructions, we should expect the network architecture to have a final layer with a dimension of 10 to match the 10 classes we are predicting. Upon inspecting the `Net()` and `Net2()` architecture, the final layer indeed outputs 10 as seen below.

`Net()` final linear layer outputs 10 items.

```
1 ...
2 ...
3 self.fc3 = nn.Linear(84, 10)
```

`Net2()` final linear layer outputs 10 items.

```
1 ...
2 ...
3 self.fc3 = nn.Linear(100, 10)
```

## 1.5 Analyze Resolution Changes

The input images are of size  $32 \times 32$ .

`Net()`: Quickly reduces resolution due to lack of padding, which limits its ability to retain spatial detail.

`Net2()`: Uses padding and pooling effectively to control resolution changes, gradually reducing spatial dimensions while preserving spatial features for higher accuracy on complex tasks.

## 1.6 Kernel Size and Padding

- **Larger Kernels ( $5 \times 5$ ):**

- **Net() architecture**

- \* Both convolutional layers use  $5 \times 5$  kernels, capturing larger spatial patterns but requiring more parameters.
- \* Rapid spatial resolution reduction can hinder the preservation of fine-grained spatial information.

- **Net2() architecture**

- \* Conv1 uses a  $5 \times 5$  kernel with padding to preserve spatial dimensions.
- \* Subsequent layers use smaller kernels ( $3 \times 3$ ,  $2 \times 2$ ) to refine features and reduce computational complexity.

- **Smaller Kernels ( $3 \times 3, 2 \times 2$ ):**
  - They focus on smaller, local patterns, making them computationally efficient.
  - **Net2() architecture** Uses these kernels in deeper layers to progressively refine features that have been extracted in earlier layers, which improves performance.

## Padding

- **Without Padding:**
  - **Net() architecture**
    - \* The lack of padding reduces spatial resolution after each convolution.
    - \* Example: Input ( $32 \times 32$ )  $\rightarrow$  After Conv1 ( $28 \times 28$ )  $\rightarrow$  After Conv2 ( $24 \times 24$ ).
    - \* **Consequences?**
      - Loss of valuable edge or boundary information as they are lost when convolving.
- **With Padding:**
  - **Net2() architecture**
    - \* Padding make sures the spatial dimensions are preserved after convolutions.
    - \* Example: Input ( $32 \times 32$ )  $\rightarrow$  After Conv1 ( $32 \times 32$ ).
    - \* **Consequences?**
      - Preserve spatial information for deeper layers.

## Performance Implications

- **Feature Extraction:**
  - **Net() architecture** Limited filters and no padding lead to poor feature extraction.
  - **Net2() architecture** Using padding and progressive refinement improve feature extraction.

## 2 Programming Tasks

### 2.1 Dataset Creation

**NOTE:** Some parts of the code below may have been adopted from the code snippet provided in the homework. [1]

The code below specify the paths for the dataset and annotations, specify the augmentation transformation and define a function to create the 3 datasets:

- Single instance
- Multiple instances, same
- Multiple instances, different

If the data were not sufficient to extract the required dataset (500 images in total per category for each dataset, 400 of which are for training, 100 for validation), then the code will use augmentation to transform some images to meet the required threshold.

The categories are as follows:

- airplane
- train
- bird
- giraffe
- clock

```
1 import os
2 import shutil
3 from pycocotools.coco import COCO
4 from PIL import Image
5 import numpy as np
6 import torch
7 import torch.nn as nn
8 import torch.optim as optim
9 from torch.utils.data import Dataset, DataLoader
10 import torchvision.transforms as transforms
11 from sklearn.metrics import confusion_matrix
12 import matplotlib.pyplot as plt
13 import seaborn as sns
14 import random
15 from tabulate import tabulate
16
17 # Set paths
18 data_dir = "data"
19 ann_file = os.path.join(data_dir, "annotations", "instances_train2014.json")
20 image_dir = os.path.join(data_dir, "train2014")
21 output_dir = "output_datasets"
22 os.makedirs(output_dir, exist_ok=True)
23
24 # Selected classes (5 from the single-instance categories)
25 classes = ['airplane', 'train', 'bird', 'giraffe', 'clock']
26 datasets = ['single_instance', 'multi_instance_same', 'multi_instance_diff']
27
28 # Initialize COCO API
29 coco = COCO(ann_file)
30
31 # Augmentation transformation
32 augmentation_transform = transforms.Compose([
33     transforms.RandomHorizontalFlip(),
34     transforms.RandomRotation(degrees=30),
35     transforms.ColorJitter(brightness=0.2, contrast=0.2, saturation=0.2, hue=0.1),
36     transforms.RandomResizedCrop((64, 64), scale=(0.8, 1.0)),
37     transforms.ToTensor(),
38     transforms.ToPILImage(),
39 ])
40
```

```

41
42
43 #When the images are deficient (did not meet the training and validation threshold of
    500 images, augment some images)
44 def augment_images(input_dir, output_dir, required_count):
45     """
46     Augments images in the input directory until the required count is met.
47     """
48     existing_images = os.listdir(input_dir)
49     augmented_count = len(existing_images)
50
51     while augmented_count < required_count:
52         # randomly select an image for augmentation
53         image_name = random.choice(existing_images)
54         img_path = os.path.join(input_dir, image_name)
55         img = Image.open(img_path)
56
57         #Apply the augmentation transform
58         augmented_img = augmentation_transform(img)
59
60         #save the augmented image
61         augmented_img.save(os.path.join(output_dir, f"aug_{augmented_count}.jpg"))
62         augmented_count += 1
63
64 #A function to extract the images
65 def extract_images_for_class(class_name, dataset_type, min_instances, max_instances,
    multiple_categories=False):
66     cat_ids = coco.getCatIds(catNms=[class_name])
67     img_ids = coco.getImgIds(catIds=cat_ids)
68
69     #To remove the duplicates
70     img_ids = list(set(img_ids))
71
72     #Create the directories for training and validation
73     train_dir = os.path.join(output_dir, dataset_type, 'train', class_name)
74     val_dir = os.path.join(output_dir, dataset_type, 'val', class_name)
75     os.makedirs(train_dir, exist_ok=True)
76     os.makedirs(val_dir, exist_ok=True)
77
78     extracted = 0
79     for img_id in img_ids:
80         if extracted >= 500: # because the total images are 500. 400 for training and
            100 for validation
81             break
82         img_info = coco.loadImgs(img_id)[0]
83         ann_ids = coco.getAnnIds(imgIds=img_id, iscrowd=False)
84         anns = coco.loadAnns(ann_ids)
85
86
87         obj_counts = {}
88         for ann in anns:
89             cat_name = coco.loadCats(ann['category_id'])[0]['name']
90             obj_counts[cat_name] = obj_counts.get(cat_name, 0) + 1
91
92         if multiple_categories:
93             valid = (obj_counts.get(class_name, 0) >= min_instances) and (len(obj_counts)
                >= 2)
94         else:
95             valid = (obj_counts.get(class_name, 0) >= min_instances)
96             if max_instances is not None:
97                 valid = valid and (obj_counts.get(class_name, 0) <= max_instances)
98
99         if valid:
100             # Saving
101             save_dir = train_dir if extracted < 400 else val_dir
102             img_path = os.path.join(image_dir, img_info['file_name'])
103             img = Image.open(img_path).resize((64, 64))
104             img.save(os.path.join(save_dir, img_info['file_name']))
105             extracted += 1
106
107     # If the threshold is not met, then augment images.
108     if extracted < 400:
109         augment_images(train_dir, train_dir, 400)

```

```

110     if extracted < 500:
111         augment_images(val_dir, val_dir, 100)

```

## Extracting the images

```

1 # Create datasets
2 for object_class in classes:
3     # Dataset 1: Single instance
4     extract_images_for_class(object_class, 'single_instance', 1, 1)
5     # Dataset 2: Multi instance of same object
6     extract_images_for_class(object_class, 'multi_instance_same', 2, None)
7     # Dataset 3: Multi instance of different objects
8     extract_images_for_class(object_class, 'multi_instance_diff', 1, None,
9                               multiple_categories=True)
10 print ("Successfully extracted the images")

```

## Creting the Dataset and Dataloader

The following creates a dataset to be used by the dataloader.

```

1 # Dataset creation
2 class COCODataset(Dataset):
3     def __init__(self, parent_dir, split, transform=None):
4         self.transform = transform
5         self.image_paths = []
6         self.labels = []
7         class_to_idx = {object_class: i for i, object_class in enumerate(classes)}
8
9         for object_class in classes:
10            cls_dir = os.path.join(parent_dir, split, object_class)
11            for img_file in os.listdir(cls_dir):
12                self.image_paths.append(os.path.join(cls_dir, img_file))
13                self.labels.append(class_to_idx[object_class])
14
15            def __len__(self):
16                return len(self.image_paths)
17
18            def __getitem__(self, idx):
19                img = Image.open(self.image_paths[idx]).convert('RGB')
20                if self.transform:
21                    img = self.transform(img)
22                return img, self.labels[idx]

```

## 2.2 Verifying that each dataset has 400 training and 100 testing images

### Training and validation count

To ensure that the training and validation numbers are correct (400 images per class for training, 100 images per class for validation), we add a statement that perform augmentation whenever there is deficiency in those thresholds. The augmentation basically make new versions of the same images but with some random rotation, etc. The augmentation code has been included in the **Dataset Creation** subsection previously.

Table. 2 and Table. 3 show that we have successfully extracted the correct count for each class.

We use the following ode to calculate the dataset:

```

1 #Print the classes counts
2 def create_count_tables(datasets):
3     train_table = []
4     val_table = []
5
6     for dataset in datasets:
7         train_row = [dataset]
8         val_row = [dataset]
9
10        for cls in classes:

```

```

11     # Count training images
12     train_path = os.path.join(output_dir, dataset, 'train', cls)
13     train_count = len(os.listdir(train_path))
14     train_row.append(train_count)
15
16     # Count validation images
17     val_path = os.path.join(output_dir, dataset, 'val', cls)
18     val_count = len(os.listdir(val_path))
19     val_row.append(val_count)
20
21     train_table.append(train_row)
22     val_table.append(val_row)
23
24     # Create headers
25     headers = ["Dataset Type"] + classes
26
27     print("\nTraining Data Counts:")
28     print(tabulate(train_table, headers=headers, tablefmt="grid"))
29
30     print("\nValidation Data Counts:")
31     print(tabulate(val_table, headers=headers, tablefmt="grid"))
32
33 # Generate tables
34 create_count_tables(datasets)

```

Dataset Type	airplane	train	bird	giraffe	clock
single_instance	400	400	400	400	400
multi_instance_same	400	400	400	400	400
multi_instance_diff	400	400	400	400	400

Table 2: Training Data Counts

Dataset Type	airplane	train	bird	giraffe	clock
single_instance	100	100	100	100	100
multi_instance_same	100	100	100	100	100
multi_instance_diff	100	100	100	100	100

Table 3: Validation Data Counts

## 2.3 Data visualization

to visualize the datasets, we plot 3 images from all the five categories for all the 3 datasets.

```
1 def plot_dataset_samples(datasets):
2     # Create one plot per dataset type
3     for dataset_idx, dataset in enumerate(datasets):
4         fig, axes = plt.subplots(len(classes), 3, figsize=(15, 20))
5
6         # Main title for each plot
7         dataset_names = {
8             'single_instance': 'Single Instance Dataset',
9             'multi_instance_same': 'Multi-Instance (Same) Dataset',
10            'multi_instance_diff': 'Multi-Instance (Different) Dataset'
11        }
12        fig.suptitle(f"{dataset_names[dataset]} Samples\n"
13                    "(Rows: Object Categories, Columns: Random Samples)",
14                    fontsize=14,
15                    fontweight='bold',
16                    y=1.02)
17
18        # Configure grid layout
19        plt.subplots_adjust(hspace=0.3, wspace=0.1)
20
21        # Remove empty axes
22        for ax in axes.flat:
23            ax.axis('off')
24            ax.set_xticks([])
25            ax.set_yticks([])
26
27        # Plot images for each category
28        for row_idx, class_name in enumerate(classes):
29            img_dir = os.path.join(output_dir, dataset, 'train', class_name)
30            all_images = [f for f in os.listdir(img_dir) if f.endswith('.jpg')]
31            random.shuffle(all_images)
32
33            # Plot 3 samples across columns
34            for col_idx in range(3):
35                ax = axes[row_idx, col_idx]
36                if col_idx < len(all_images):
37                    img = Image.open(os.path.join(img_dir, all_images[col_idx]))
38                    ax.imshow(img)
39                    ax.set_title(f"Sample {col_idx+1}", fontsize=9, pad=4)
40
41            # Add row labels only to first column
42            if col_idx == 0:
43                ax.text(-0.3, 0.5, class_name,
44                       rotation=0,
45                       fontsize=12,
46                       va='center',
47                       ha='right',
48                       transform=ax.transAxes)
49
50        plt.tight_layout()
51        plt.show()
52
53 datasets = ['single_instance', 'multi_instance_same', 'multi_instance_diff']
54 plot_dataset_samples(datasets)
```

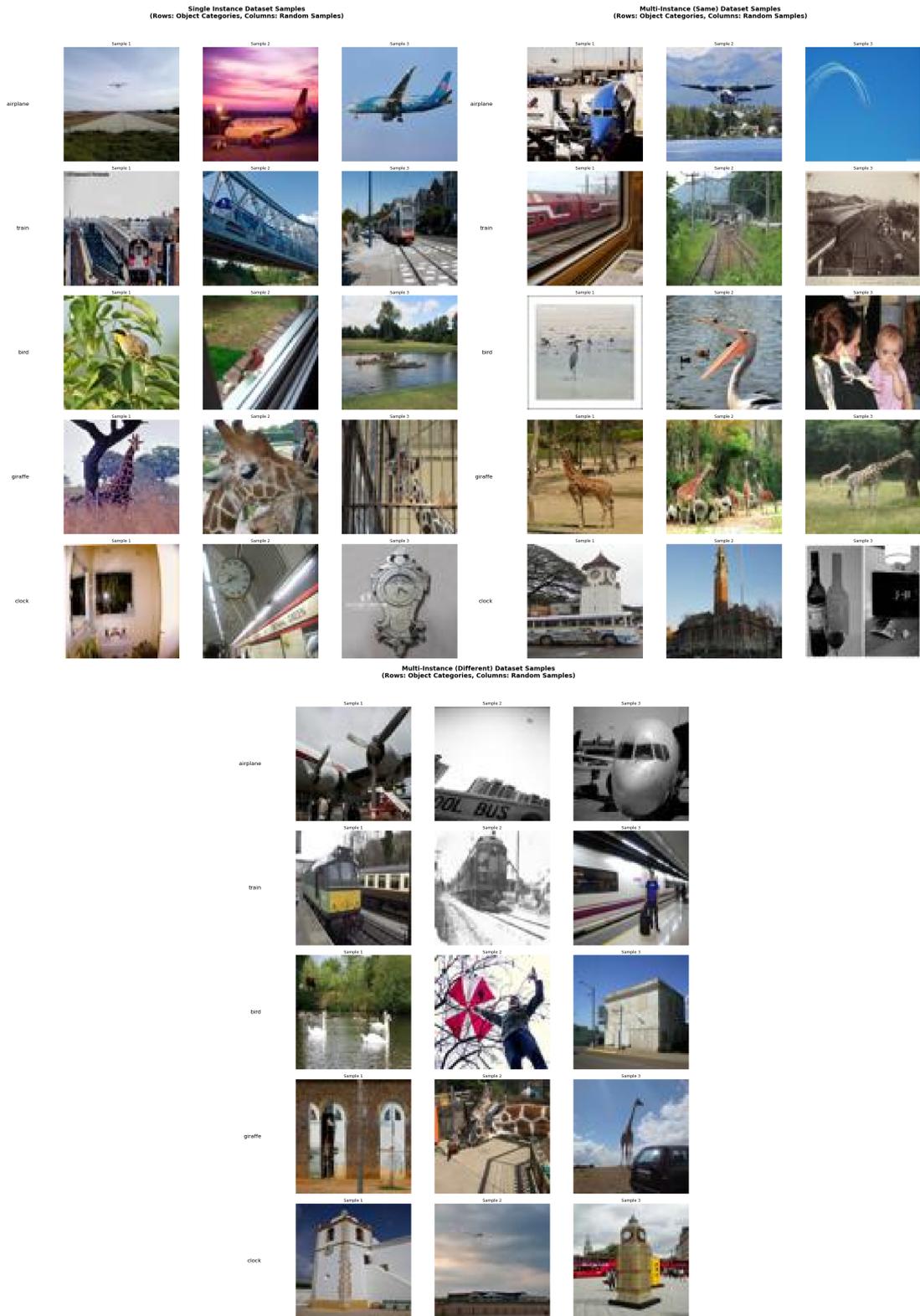


Figure 3: From left to right, examples of the **Single instance**, followed by **Multiple instances, same** and **Multiple instances, different** images.

## 2.4 Tables for Train and Validation Count

The count tables for training and validation have been included in Table. 2 and Table. 3 respectively.

## 2.5 Network implementation

The following is the network implementation as instructed.

```
1 # Model Definition
2 class HW4Net(nn.Module):
3     def __init__(self):
4         super().__init__()
5         self.conv1 = nn.Conv2d(3, 16, 3)
6         self.pool = nn.MaxPool2d(2, 2)
7         self.conv2 = nn.Conv2d(16, 32, 3)
8         self.fc1 = nn.Linear(32 * 14 * 14, 64)
9         self.fc2 = nn.Linear(64, 5)
10
11     def forward(self, x):
12         x = self.pool(nn.functional.relu(self.conv1(x)))
13         x = self.pool(nn.functional.relu(self.conv2(x)))
14         x = x.view(x.size(0), -1)
15         x = nn.functional.relu(self.fc1(x))
16         x = self.fc2(x)
17         return x
```

### Model Parameters

To calculate the model parameters, we use the following code:

```
1 def get_num_parameters(model):
2     total_params = 0
3     trainable_params = 0
4     for param in model.parameters():
5         num_params = param.numel()
6         total_params += num_params
7         #Check if the parameters is trainable.
8         #Our model do not have non-trainable paremeters, so we expect it to be the same
9         as total
10        if param.requires_grad:
11            trainable_params += num_params
12    return total_params, trainable_params
```

	Number of Conv. Layers	Number of parameters
HW4Net (baseline)	2	406885
HWNNet1 (bonus)	5	1119877
HWNNet2 (bonus)	10	1706133

Table 4: The number of parameters (trainable) for each of the networks we will be using. The baseline network is the first network **HW4Net**. The other two are the network used for the bonus part.

### Dataset and Dataloader

Using the Dataset and Dataloader classes from *torch.utils.data.Dataset* and *torch.utils.data.Dataloader*, I wrote the following classes to create the dataset and load it.

```
1 # Dataset creation
2 class COCODataset(Dataset):
3     def __init__(self, parent_dir, split, transform=None):
4         self.transform = transform
5         self.image_paths = []
6         self.labels = []
7         class_to_idx = {object_class: i for i, object_class in enumerate(classes)}
8
9         for object_class in classes:
10            cls_dir = os.path.join(parent_dir, split, object_class)
11            for img_file in os.listdir(cls_dir):
12                self.image_paths.append(os.path.join(cls_dir, img_file))
13                self.labels.append(class_to_idx[object_class])
14
15    def __len__(self):
```

```

16         return len(self.image_paths)
17
18     def __getitem__(self, idx):
19         img = Image.open(self.image_paths[idx]).convert('RGB')
20         if self.transform:
21             img = self.transform(img)
22         return img, self.labels[idx]
23
24     # Transforms
25     transform = transforms.Compose([
26         transforms.ToTensor(),
27         transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5))
28     ])
29
30     model_names = ["HW1", "HW2", "HW4"]
31     results = {} #To hold the accuracies and taining loss
32     epochs = 100
33
34     for ds in datasets:
35         train_dataset = COCODataset(os.path.join(output_dir, ds), 'train', transform)
36         val_dataset = COCODataset(os.path.join(output_dir, ds), 'val', transform)
37         train_loader = DataLoader(train_dataset, batch_size=32, shuffle=True)
38         val_loader = DataLoader(val_dataset, batch_size=32, shuffle=False)
39
40         results[ds]={}
41         for model_name in model_names:
42
43             ....
44             ....
45             ....

```

## Training routine

```

1
2
3     # Training Function
4     def train_model(model, train_loader, val_loader, dataset_name, epochs, model_name):
5         device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
6         model.to(device)
7         criterion = nn.CrossEntropyLoss()
8         optimizer = optim.Adam(model.parameters(), lr=1e-3, betas=(0.9, 0.99))
9
10        train_losses = []
11        val_accuracies = []
12
13        for epoch in range(epochs):
14            model.train()
15            running_loss = 0.0
16            num_batches = len(train_loader)
17            for i, (inputs, labels) in enumerate(train_loader):
18                inputs, labels = inputs.to(device), labels.to(device)
19
20                #Clearing the gradients
21                optimizer.zero_grad()
22
23                #Prediction
24                outputs = model(inputs)
25                loss = criterion(outputs, labels)
26
27                #Propagating the loss
28                loss.backward()
29                optimizer.step()
30                running_loss += loss.item()
31
32                #print("i = ", i)
33                #train_losses.append(loss)
34
35                #if (i+1) % 100 == 0:
36                #    avg_loss = running_loss / 100
37                #    train_losses.append(avg_loss)
38                #    running_loss = 0.0

```

```

39     avg_loss = running_loss / num_batches
40     train_losses.append(avg_loss)
41
42
43
44     # Validation
45     correct = 0
46     total = 0
47     model.eval()
48     with torch.no_grad():
49         for inputs, labels in val_loader:
50             inputs, labels = inputs.to(device), labels.to(device)
51             outputs = model(inputs)
52             _, predicted = torch.max(outputs.data, 1)
53             total += labels.size(0)
54             correct += (predicted == labels).sum().item()
55     accuracy = 100 * correct / total
56     val_accuracies.append(accuracy)
57
58     if (epoch%10==0):
59
60         print(f"Model: {model_name}. Dataset: {dataset_name} - Epoch {epoch+1},
61               Training loss: {avg_loss:.4f} | Val Acc: {accuracy:.2f}%")
62
63     return train_losses, val_accuracies
64
65 #Plot Confusion Matrix
66 def evaluate(model, val_loader, dataset_name, model_name):
67     device = next(model.parameters()).device
68     all_labels = []
69     all_preds = []
70     model.eval()
71     with torch.no_grad():
72         for inputs, labels in val_loader:
73             inputs = inputs.to(device)
74             outputs = model(inputs)
75             _, preds = torch.max(outputs, 1)
76             all_labels.extend(labels.cpu().numpy())
77             all_preds.extend(preds.cpu().numpy())
78
79     cm = confusion_matrix(all_labels, all_preds)
80     plt.figure(figsize=(10,7))
81     sns.heatmap(cm, annot=True, fmt='d', xticklabels=classes, yticklabels=classes)
82     plt.title('Model: ' + model_name + ' - Dataset: ' + dataset_name + ' - Confusion
83               Matrix')
84     plt.show()
85
86 model_names = ["HW1", "HW2", "HW4"]
87 results = {} #To hold the accuracies and taining loss
88 epochs = 100
89
90 for ds in datasets:
91     train_dataset = COCODataset(os.path.join(output_dir, ds), 'train', transform)
92     val_dataset = COCODataset(os.path.join(output_dir, ds), 'val', transform)
93     train_loader = DataLoader(train_dataset, batch_size=32, shuffle=True)
94     val_loader = DataLoader(val_dataset, batch_size=32, shuffle=False)
95
96     results[ds]={}
97     for model_name in model_names:
98
99         if model_name == "HW1":
100             model = HWNet1()
101         elif model_name == "HW2":
102             model = HWNet2()
103         elif model_name == "HW4":
104             model = HW4Net()
105
106     total_params, trainable_params = get_num_parameters(model)
107     print ("Current model: ", model_name, ". Total params: ", str(total_params))
108     print ("Current model: ", model_name, ". Trainable params: ", str(

```

```

109     print("Model: ", model_name, ". Training on ", ds, " dataset.")
110
111     #Training
112     train_loss, val_acc = train_model(model, train_loader, val_loader, ds, epochs,
113                                     model_name)
114
115     #Appending the results
116     #results[ds] = {'train_loss': train_loss, 'val_acc': val_acc}
117     results[ds][model_name] = {'train_loss': train_loss, 'val_acc': val_acc}
118
119     #Evaluating (confusion matrix)
120     evaluate(model, val_loader, ds, model_name)
121
122     # Plotting the loss
123     plt.figure(figsize=(10, 5))
124     #plt.plot(range(1, epochs + 1), train_loss, label='Train Loss') # x-axis from 1
125     #to epochs
126     plt.plot(train_loss, label='Train Loss') # x-axis from 1 to epochs
127     #plt.plot(range(1, epochs + 1), val_acc, label='Validation Accuracy') #If
128     #you also want to plot validation accuracy
129     plt.xlabel('Epoch')
130     plt.ylabel('Loss')
131     plt.title('Model: ' + model_name + ' - Dataset: ' + ds + ' - Training Loss')
132     plt.legend()
133     plt.grid(True)
134     plt.xticks(range(1, epochs + 1)) # Set x-axis ticks to be integers
135     plt.show()
136
137     #If you also want to plot the validation accuracy:
138     plt.figure(figsize=(10, 5))
139     #plt.plot(range(1, epochs + 1), val_acc, label='Validation Accuracy', color='red
140     #')
141     plt.plot(val_acc, label='Validation Accuracy', color='red')
142     plt.xlabel('Epoch')
143     plt.ylabel('Validation Accuracy')
144     plt.title('Model: ' + model_name + ' - Dataset: ' + ds + ' - Validation
145     #accuracy')
146     #plt.title(ds+' with model: ' + model_name + ' - Validation Accuracy vs. Epoch')
147     plt.legend()
148     plt.grid(True)
149     plt.xticks(range(1, epochs + 1)) # Set x-axis ticks to be integers
150     plt.show()

```

### 3 Baseline HW4Net() Performance

The baseline network only contained 2 convolutional layers with limited number of filters (16, 32). The performance on such a CNN is limited and usually fails to learn the finer details that allows is to descricate between classes of objects. We anticipate that the other two (included in the bonus) network to perform much better.

In my case, the network performed best when it was presented with multiple instances of the same object, followed closely by single instance, and performed noticeably worse on the multi instances with different objects, which suggests that the more comlicated scenario where multiple different objects are present may need bigger network or more confolotional layers or differnet approach. And having multiple instances of the same object may assist the network to generalize as it see multiple instances with different conditions of the same class.

#### 3.1 Baseline HW4Net() Single Instance Loss

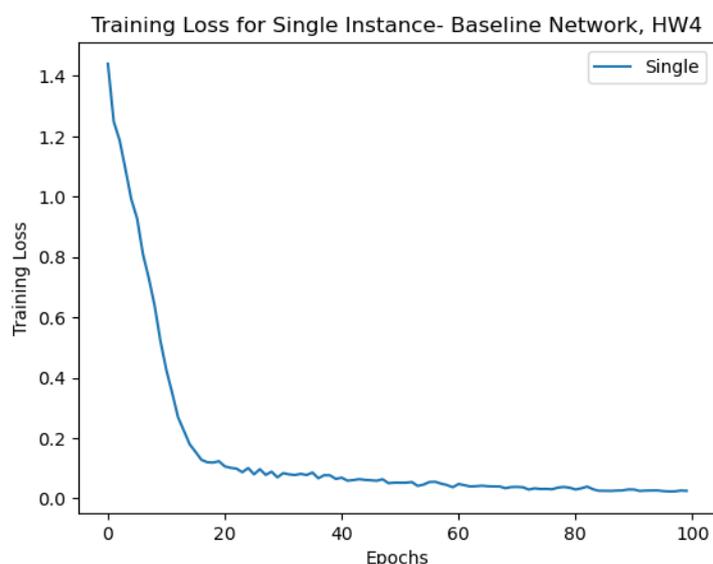


Figure 4: The baseline HW4Net() architecture with **Single Instance**

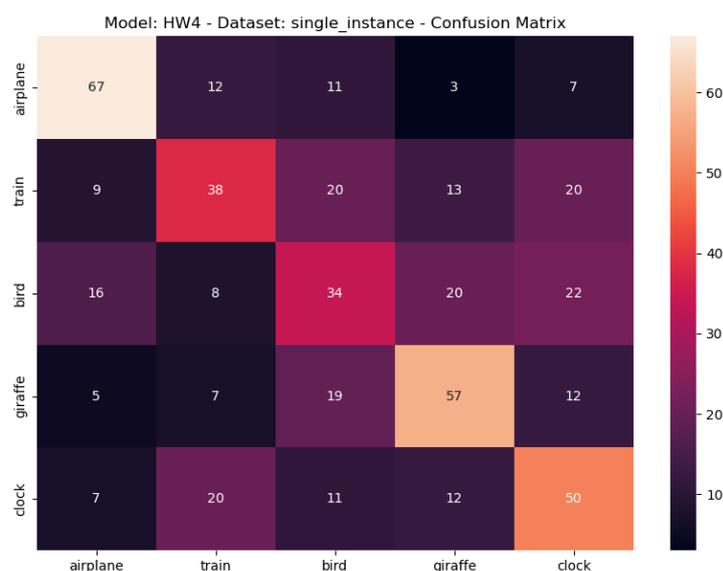


Figure 5: Confusion matrix for the single instance dataset on baseline HW4 network

### 3.2 Baseline HW4Net() Multi Instance Same Loss

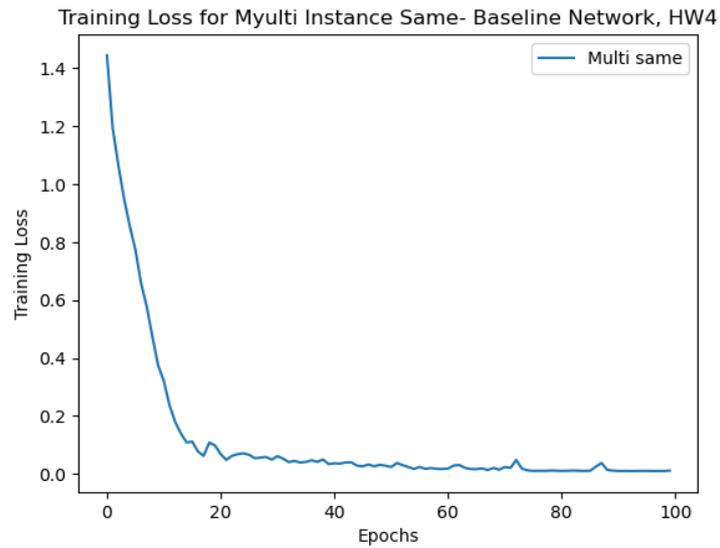


Figure 6: The baseline **HW4Net()** architecture with **Multi Instance Same**

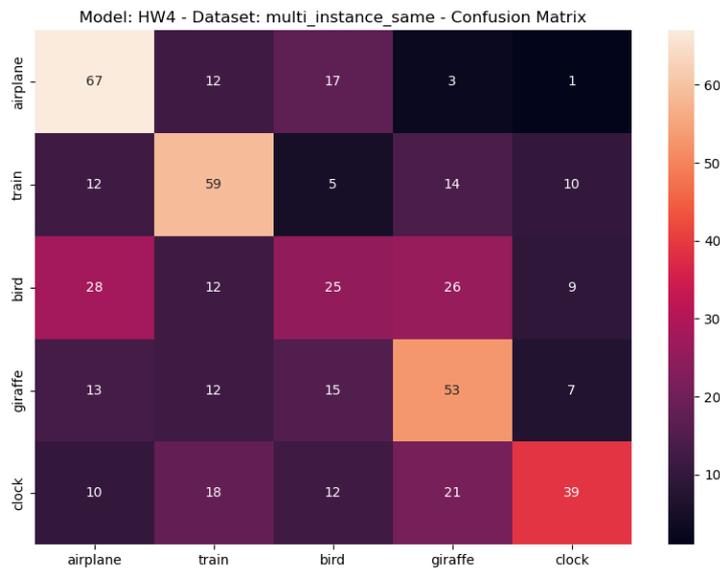


Figure 7: The baseline Confusion matrix for the multi instance same object dataset on baseline HW4 network

### 3.3 Baseline HW4Net() Multi Instance Different Loss

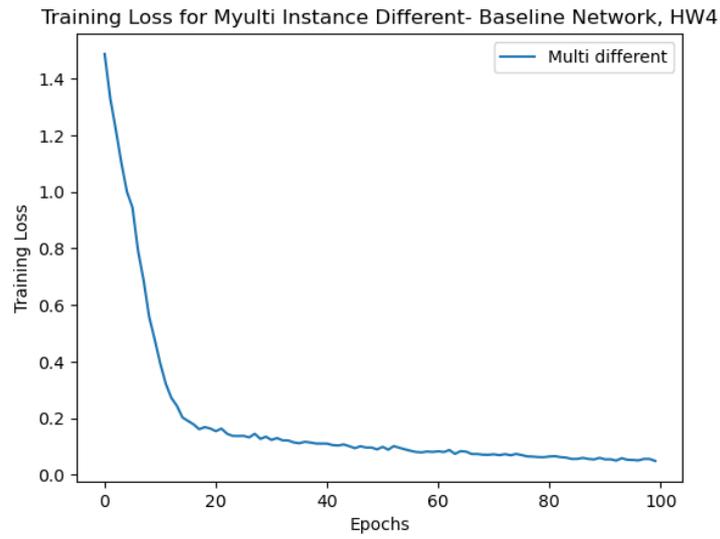


Figure 8: The baseline **HW4Net()** architecture with **Multi Instance Different**

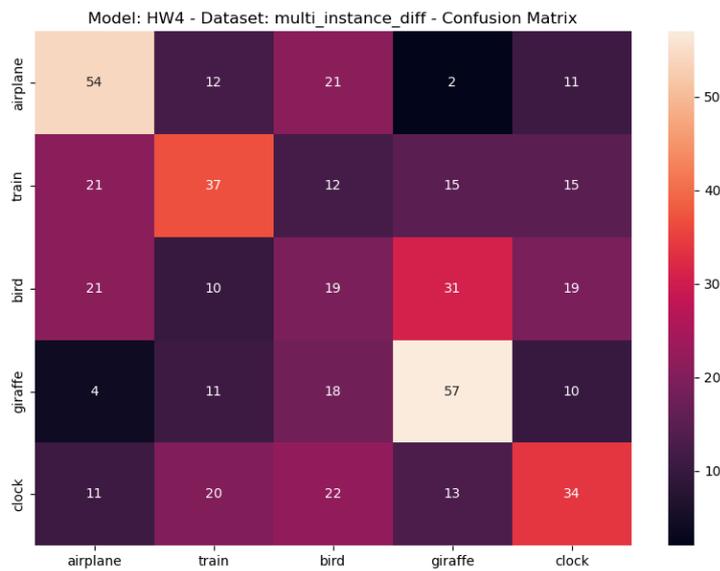


Figure 9: The baseline **HW4Net()** architecture confusion matrix for Multi Instance Same

### 3.4 Baseline HW4Net() Comparison for all Datasets.

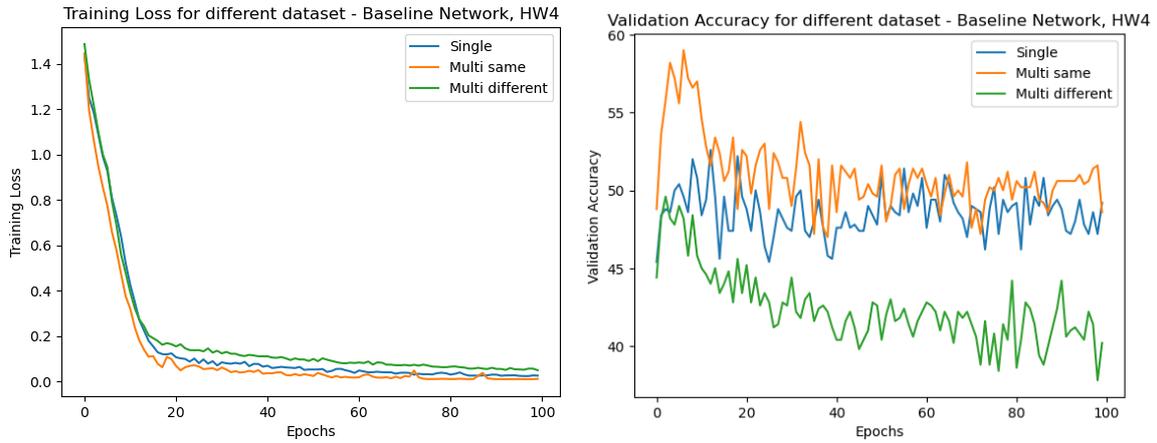


Figure 10: The baseline **HW4Net()** architecture. **Left**: loss value for the 3 different datasets as a function of epochs. **Right**: validation accuracy.

	Number of parameters	Accuracy at the end	Maximum Accuracy
Single Instance	406885	<b>49.200%</b>	52.600%
Multi Same	406885	48.600 %	<b>59.000%</b>
Multi Different	406885	40.200%	49.600%

Table 5: Baseline network **HW4Net()** number of parameters, and validation accuracy at the end final epoch and the maximum validation accuracy achieved for each dataset. The **boldface** value is the highest in each column.

## Misclassified Images

The following are examples of images that have been misclassified for each dataset. The images are normalized due to the transformation so they may appear darker than the original images. We'll try to hypothesize why might those images been misclassified.

**1. Single Instance** The images are of planes but have been categorized as: cloc, train, and train respectively. We can see the resemblance of a train car in the later two images (middle, and right image) which may look vaguely similar to the model. The clock however may be harder to interpret. However, we should bear in mind that the dataset include huge clocks that are installed in public spaces and plazas (think Big Ben in London). So those may share some edge features with the plane image on the left.



Figure 11: Misclassified images for the **Single instance** dataset

**2. Multiple Instance, same** The random selected images were again of planes, but this time they were categorized as train, giraffe, and train respectively. The image on the left have some repeated pattern which is often seen and associated with the repeated pattern of the train cars and rails, so the model might have mistaken it for a train accordingly. The middle image is harder to interpret, but if you look closely, there is a bulge (part of the plane) that is visible with a sky background. Usually giraffe have similar looking bulge (their faces) with the sky as a background, which may explain why the model thought it's a giraffe. The third image may again resemble some overall shape and color of some trains.



Figure 12: Misclassified images for the **Multi instance, same** dataset

**3. Multiple Instance, different** Those images might be the best to interpret compared with the previous datasets. The images again are for planes but have been misclassified as a bird, a clock, and a clock respectively. It is easy to see why the network could mistake the airplane for a bird, both of them are in the sky, and frankly look similar when they only occupy a small portion of the image. The middle image is more interesting, and even for humans, we could easily interpret this as a big clock (like those installed in big plazas), but with careful inspection, the propellers on the side tell us that it's a plane, something the model did not pick up during training. The third image is harder to interpret but it may have shared some features or edges or dials with a clock and hence confused the network.

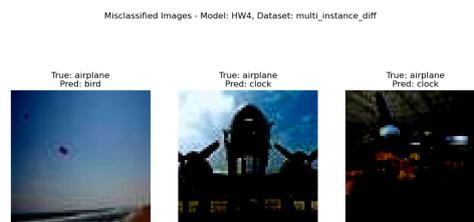


Figure 13: Misclassified images for the **Multi instance, different** dataset

**3.4.1 By observing your classification accuracies, which dataset think are more difficult to correctly differentiate and why?**

The **multi instances with different objects** was the most difficult to correctly classify. I would hypothesize that those objects are often appear cropped or at the corner of the image or only part of them shows up in the image but they are still counted as present. For example, often only part of the plane is showing in an image and it is still considered present, which presents a challenge for a small network like the baseline network we used with only 2 convolutional layers.

To account for those, we may need larger network, or convolutional layers with deeper filters (more number of filters) to extract finer descrimations between classes.

**3.4.2 By observing your confusion matrices, which class or classes do you think are more difficult to correctly differentiate and why?**

The **bird** class consistantly performed worst by a big margin compared to the other classes.

By examining the bird images, birds come in many different shapes, colors, and sizes. They are also not as deterministic in terms of pose, and location as images of planes or cars. They also appear in a wide variety of surrounding, from trees, to beaches, etc., all of which make them harder to classify.

**3.4.3 What is one thing that you propose to make the classification performance better?**

Make a bigger network with more convolutional layers and larger number of filters, and expand the augmentations to include random rotations, cropping, etc.

## 4 Bonus: Networks with more convolutional layers

As one may anticipate, using a CNN with more convolutional layers may hypothetically improve the performance. Here we implement two networks to test such hypothesis and provide the training loss, validation accuracy and some analysis. The two networks we implement are:

- **HWNet1()**: with 5 convolutional layers
- **HWNet2()**: with 10 convolutional layers

```
1 # Models for the bonus extra credit
2 class HWNet1(nn.Module):
3     def __init__(self):
4         super().__init__()
5         self.features = nn.Sequential(
6             nn.Conv2d(3, 16, 3, padding=1),
7             nn.ReLU(),
8             nn.MaxPool2d(2, 2),
9
10            nn.Conv2d(16, 32, 3, padding=1),
11            nn.ReLU(),
12            nn.Conv2d(32, 32, 3, padding=1),
13            nn.ReLU(),
14            nn.MaxPool2d(2, 2),
15
16            nn.Conv2d(32, 64, 3, padding=1),
17            nn.ReLU(),
18            nn.Conv2d(64, 64, 3, padding=1),
19            nn.ReLU(),
20            nn.MaxPool2d(2, 2),
21        )
22        self.classifier = nn.Sequential(
23            nn.Linear(64 * 8 * 8, 256),
24            nn.ReLU(),
25            nn.Linear(256, 5)
26        )
27
28    def forward(self, x):
29        x = self.features(x)
30        x = x.view(x.size(0), -1)
31        x = self.classifier(x)
32        return x
33
34 class HWNet2(nn.Module):
35     def __init__(self):
36         super().__init__()
37         self.features = nn.Sequential(
38             nn.Conv2d(3, 16, 3, padding=1),
39             nn.ReLU(),
40             nn.Conv2d(16, 16, 3, padding=1),
41             nn.ReLU(),
42             nn.MaxPool2d(2, 2),
43
44            nn.Conv2d(16, 32, 3, padding=1),
45            nn.ReLU(),
46            nn.Conv2d(32, 32, 3, padding=1),
47            nn.ReLU(),
48            nn.MaxPool2d(2, 2),
49
50            nn.Conv2d(32, 64, 3, padding=1),
51            nn.ReLU(),
52            nn.Conv2d(64, 64, 3, padding=1),
53            nn.ReLU(),
54            nn.MaxPool2d(2, 2),
55
56            nn.Conv2d(64, 128, 3, padding=1),
57            nn.ReLU(),
58            nn.Conv2d(128, 128, 3, padding=1),
59            nn.ReLU(),
60            nn.MaxPool2d(2, 2),
61
62            nn.Conv2d(128, 256, 3, padding=1),
```

```

63     nn.ReLU(),
64     nn.Conv2d(256, 256, 3, padding=1),
65     nn.ReLU(),
66     nn.MaxPool2d(2, 2),
67 )
68 self.classifier = nn.Sequential(
69     nn.Linear(256 * 2 * 2, 512),
70     nn.ReLU(),
71     nn.Linear(512, 5)
72 )
73
74 def forward(self, x):
75     x = self.features(x)
76     x = x.view(x.size(0), -1)
77     x = self.classifier(x)
78     return x

```

### Bonus: HWNet1() with 5 convolutional layers - Performance

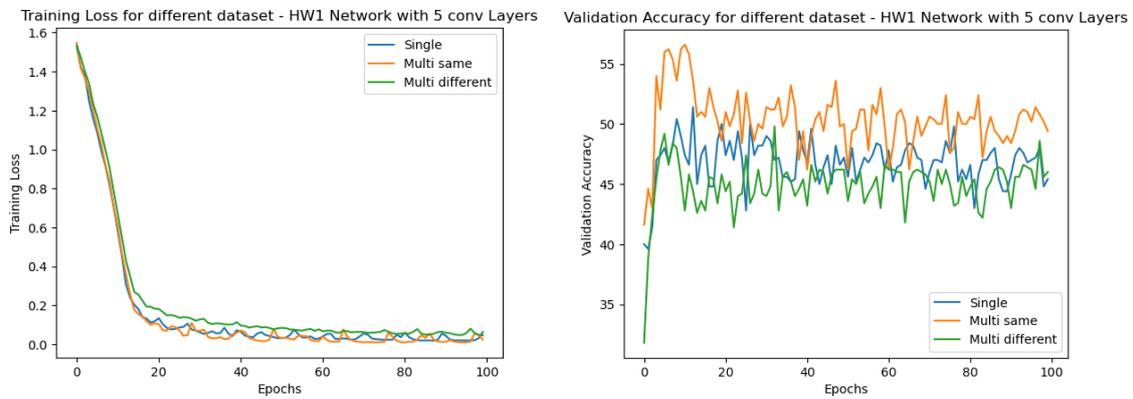


Figure 14: The bonus **HWNet1()** architecture with 5 convolutional layers. **Left**: loss value for the 3 different datasets as a function of epochs. **Right**: validation accuracy.

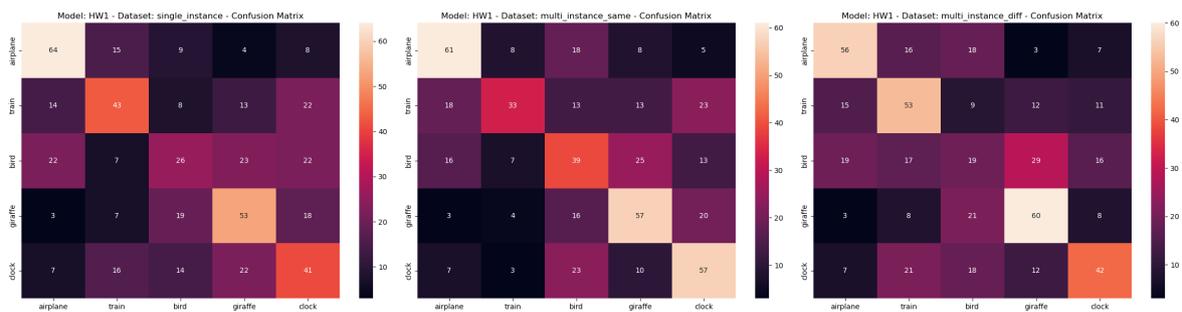


Figure 15: The bonus **HWNet1()** architecture with 5 convolutional layers. The three matrices from left to right are the confusion matrices for the single instance, multi instance same object, and multi instance different object.

## Bonus: HWNet2() with 10 convolutional layers - Performance

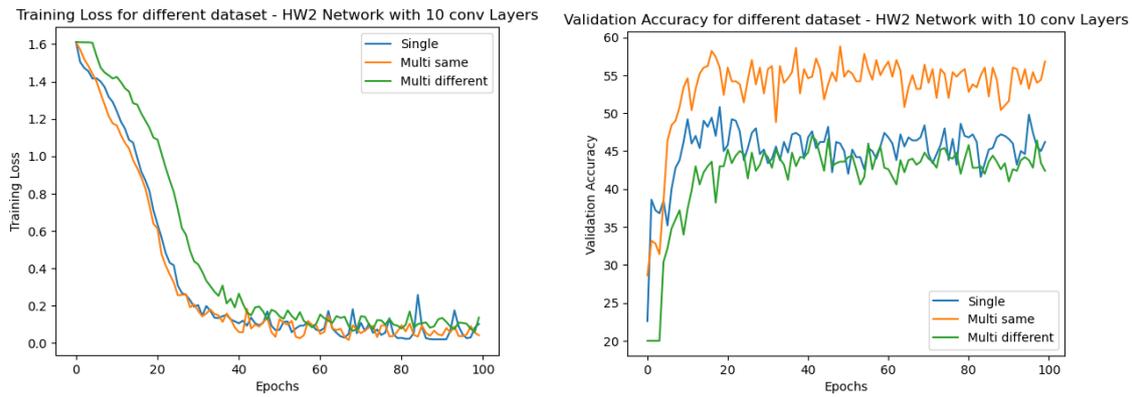


Figure 16: The bonus **HWNet2()** architecture with 10 convolutional layers. **Left**: loss value for the 3 different datasets as a function of epochs. **Right**: validation accuracy.

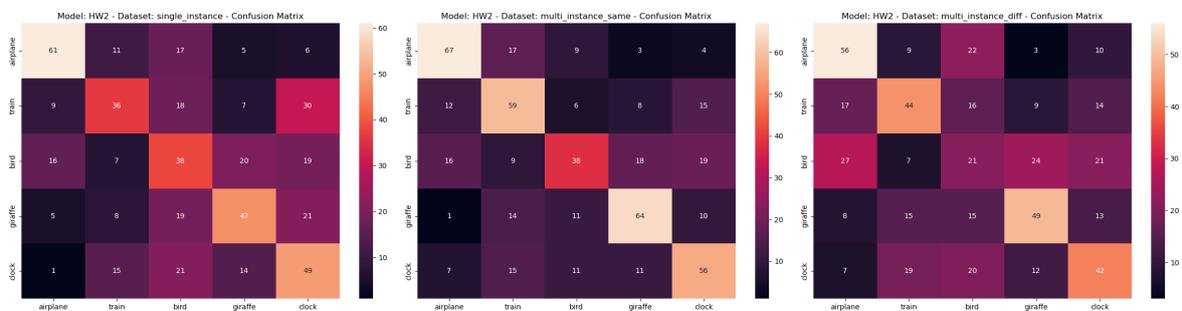


Figure 17: The bonus **HWNet2()** architecture with 10 convolutional layers. The three matrices from left to right are the confusion matrices for the single instance, multi instance same object, and multi instance different object.

### Analysis of the bonus models compared to the baseline

Network	Number of Parameters	Accuracy at the end
H1 (5 conv)	1119877	45.400%
H2 (10 conv)	1706133	46.200%
H4 Baseline	406885	<b>49.200%</b>

Table 6: Validation Accuracy for H1, H2, and H4 with for the **Single Instance** case. **Boldface** values are highest in each column.

	Number of Parameters	Accuracy at the end
H1 (5 conv)	1119877	49.400%
H2 (10 conv)	1706133	<b>56.800%</b>
H4 Baseline	406885	48.600%

Table 7: Validation Accuracy for H1, H2, and H4 with for the **Multi Instance Same** case. **Boldface** values are highest in each column.

	Number of Parameters	Accuracy at the end
H1 (5 conv)	1119877	<b>46.000%</b>
H2 (10 conv)	1706133	42.400%
H4 Baseline	406885	40.200%

Table 8: Validation Accuracy for H1, H2, and H4 with for the **Multi Instance Different** case. **Boldface** values are highest in each column.

## 4.1 Bonus: Analysis of the bonus networks

### Question: Does a deeper network always result in better performance?

Based on my experiment above with the architecture I used, the accuracy for the multi instance (same and different) both improved when using a larger network with more convolutional layers which suggests that if the task is more complicated, a more complex network might be needed to capture the underlying complexity of the data.

However, based on my limited experiment here, in the single instance case, the baseline network was sufficient and performed better than the other two. This could have been just a random coincidence, or an indication that larger network may not always be better, or they may need to train for longer to achieve comparable results on less complicated tasks.

One improvement that I could also suggest is fiddling with the number of filters, as this may result in learning finer features and improving the network classification accuracy.

## References

- [1] Avinash Kak. Bme 646 and ece 60146. deep learning. <https://engineering.purdue.edu/DeepLearn/>, 2025. Accessed: February 2025.
- [2] Avinash Kak. Dlstudio-2.5.1. <https://engineering.purdue.edu/kak/distDLS/>, 2025. Accessed: February 2025.
- [3] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014.