

DAT410 Assignment 8

Hongliang Zhong
zhongh@student.chalmers.se

Jianing Lei
jianingl@student.chalmers.se

March 21, 2025

1 Introduction

Remote sensing scene classification is a fundamental task in computer vision, aiming to automatically categorize aerial images into predefined scene classes. With the rapid development of satellite and drone-based imaging technologies, accurate classification of remote sensing images has become crucial for applications such as land use analysis, environmental monitoring, and urban planning. However, this task presents significant challenges due to the high intra-class variability (e.g., seasonal and weather changes) and low inter-class separability (e.g., visually similar scenes in different categories).

This project is directly related to **Module 5: Computer Vision and Remote Sensing** in our course. Similar to our coursework, we work with remote sensing data, but unlike Assignment 5, which focuses on a different task, our project primarily explores scene classification and evaluates the effectiveness of different deep learning models. In this work, we investigate and compare four different deep-learning architectures:

- Basic CNN – A simple convolutional neural network as a baseline.
- ResNet-50 – A widely used deep residual network with pre-trained ImageNet weights.
- Transformer-based models – Exploring vision transformers (ViT) for scene classification.
- Self-Supervised Learning (SSL) using SimCLR – A contrastive learning approach that learns representations from unlabeled data before fine-tuning on labeled samples.

To conduct this study, we use the **RSSCN7 dataset** [1], which consists of 2,800 remote sensing images from 7 typical scene categories: grassland, forest, farmland, parking lot, residential region, industrial region, and river & lake. The images are 400×400 pixels, covering diverse landscapes and environmental conditions. The variability in seasons, lighting conditions, and scales make RSSCN7 a challenging yet practical dataset for benchmarking remote sensing classification models.

By evaluating these models, we aim to analyze the impact of network depth, self-supervised learning, and transformer-based architectures on remote sensing image classification. Additionally, we explore whether self-supervised learning can reduce labeled data dependency and how Vision Transformers (ViTs) compare against traditional CNNs in this domain.

2 Methodology

2.1 Data preprocessing

We construct the dataset into a structured format with three columns: `image` (file path), `class_name` (scene label), and `class_num` (numerical label). To ensure fair and consistent evaluation, we randomly split the dataset into training, validation, and test sets with a ratio of **70%, 15%, and 15%**, respectively, while maintaining balanced class distributions. A fixed random seed was used to ensure reproducibility.

Data augmentation was applied to the training set to improve the model’s generalization ability and robustness to variations in lighting, orientation, and composition. Specifically, the training images were randomly rotated, horizontally flipped, and color jittered in terms of brightness, contrast, saturation, and hue. All images were then converted to tensors and normalized using ImageNet mean and standard deviation values. For the validation and test sets, only resizing and normalization were applied, without additional augmentations, to ensure fair evaluation.

The training process was configured with consistent hyperparameters across all experiments: the batch size was set to 32, the number of training epochs was 30, and the learning rate was 0.0001. To avoid overfitting and reduce unnecessary computation, early stopping was employed with patience of 5 epochs, monitoring the validation loss throughout training.

2.2 Basic CNN

As a baseline for comparison, we implemented a simple convolutional neural network (CNN) architecture to perform scene classification on the RSSCN7 dataset. The network includes three convolutional blocks, each comprising a convolutional layer with increasing channel depth, a ReLU activation function, and a max pooling operation to progressively reduce the spatial resolution while capturing hierarchical features. The output feature maps are then flattened and passed through two fully connected layers, with the final layer producing logits corresponding to the seven scene categories. This model serves as an essential reference point to evaluate the effectiveness of more advanced models.

2.3 ResNet-50

ResNet-50 is a widely used deep convolutional neural network known for its residual learning framework. It consists of 50 layers with identity shortcut connections that allow the network to mitigate the vanishing gradient problem during backpropagation. These residual connections enable the training of significantly deeper networks without degradation in performance, making ResNet-50 a powerful choice for image classification tasks.

In this project, we used the ImageNet-pretrained version of ResNet-50 and replaced its final fully connected (FC) layer to adapt it to our specific classification task.

Several key design choices influence the fine-tuning process and impact the model’s ability to generalize while maintaining efficient training. First, we freeze the initial 140 layers of ResNet-50 to preserve the fundamental visual features learned from ImageNet, such as edges and textures, while allowing only the final layers to adapt to the new dataset. This approach reduces the risk of overfitting and significantly decreases computational costs. Additionally, we replace the fully connected (FC) layer to match the number of target classes. A 256-unit hidden layer with ReLU activation is introduced to enhance learning capacity while maintaining efficiency. To further mitigate overfitting, we apply a dropout rate of 0.3, randomly deactivating neurons during training. The final output layer is adjusted to align with the number of classes in our classification task, ensuring compatibility and effective adaptation.

2.4 Transformer-based model

In addition, we explored a transformer-based model for our remote sensing scene classification task. We utilized the ViT-Base-Patch16-224 model introduced by Google, which applies the transformer architecture—originally developed for natural language processing—to image classification tasks.

The Vision Transformer (ViT) divides each input image into fixed-size patches (in our case, 16×16 pixels). Each patch is flattened and linearly projected into an embedding vector, forming a sequence that is fed into a standard transformer encoder. For the ViT-Base-Patch16-224 model, each image must be resized to 224×224 pixels before being split into patches. This is a notable distinction from the other models in our study, which implement directly on the original 400×400

resolution. Therefore, a key preprocessing step specific to the transformer model was resizing the original images to 224×224 prior to training and inference.

The model also combines positional encodings to retain spatial information and uses a learnable class token to aggregate global features for classification. Unlike convolutional networks, ViT relies on self-attention mechanisms to capture global context information. In our implementation, we fine-tuned the pre-trained model on the RSSCN7 dataset after replacing its classification head with a new fully connected layer suitable for seven classes.

2.5 Self-Supervised Learning (SSL)

2.5.1 SimCLR pretraining

Self-supervised learning has emerged as a powerful approach for representation learning, especially in scenarios with limited labeled data. In this project, we adopt SimCLR (Simple Framework for Contrastive Learning of Visual Representations) as our self-supervised method. SimCLR learns visual representations by maximizing the agreement between differently augmented views of the same image, without relying on labels during pretraining.

The core idea of SimCLR is to treat each augmented image pair as a positive pair and all other images in the batch as negatives. For each image, two augmented views are generated through a stochastic data augmentation pipeline. These are passed through a shared encoder (in our case, ResNet-50) and a projection head to obtain two latent representations. The model is trained using the Normalized Temperature-scaled Cross Entropy Loss (NT-Xent Loss), which encourages positive pairs to be close in the latent space while pushing apart the negative pairs [2]. We used a temperature parameter of 0.3 in the loss function to control the sharpness of similarity distribution.

In our implementation, the data augmentation for contrastive learning includes random resized cropping with a scale range of (0.2, 1.0), horizontal flipping, random application of color jittering (brightness, contrast, saturation, hue), and random grayscale conversion. These augmentations are useful for encouraging the model to learn invariant representations. All images are normalized using the standard ImageNet statistics.

2.5.2 Fine-tuning on RSSCN7

After pretraining, we evaluated the learned representations through supervised fine-tuning on the RSSCN7 dataset. We removed the projection head from the pre-trained SimCLR model and appended a new classification head consisting of a fully connected layer with 256 hidden units, a ReLU activation, dropout regularization, and a final output layer with seven units corresponding to the scene categories.

To investigate whether SimCLR pretraining can reduce the dependency on large amounts of labeled samples while still achieving competitive performance, we conducted two fine-tuning experiments. In the first setting, we used the entire training set (70% of the dataset), while in the second setting, we randomly sampled 50% of the training set. In both cases, the model was trained by using the labeled data.

3 Results and discussion

3.1 Performance Comparison Based on Test Accuracy

Table 1 presents the test accuracy of different models on the RSSCN7 dataset. Among the evaluated models, ResNet-50 achieves the highest accuracy (95.71%), followed closely by ViT-Base-Patch16-224 (94.05%). This suggests that deep convolutional networks and transformer-based architectures are highly effective for remote sensing image classification. The Basic CNN, serving as a baseline, performs significantly worse (76.19%), highlighting the benefits of deeper architectures and pretraining.

The SimCLR-based models, which employ self-supervised learning (SSL) followed by fine-tuning, exhibit lower accuracy than fully supervised models. The SimCLR model fine-tuned with 100% of the training data achieves 77.62%, slightly better than the Basic CNN but still well below ResNet-50 and ViT. When the labeled training data is reduced to 50%, the performance drops further to 68.33%, indicating that self-supervised pretraining alone does not fully compensate for the reduced supervision in this setting.

Model	Test Accuracy (%)
Basic CNN	76.19
ResNet-50	95.71
ViT-Base-Patch16-224	94.05
SimCLR + Fine-tuning (100% training data)	77.62
SimCLR + Fine-tuning (50% training data)	68.33

Table 1: Test accuracy comparison of different models

3.2 Training and Validation Loss/Accuracy Analysis

Examining the training and validation loss curves A reveals notable differences in convergence behavior across models. ResNet-50 demonstrates stable training, reaching near-optimal performance without significant overfitting. ViT, however, appears to require more epochs to converge, aligning with previous findings that transformer-based models often need extensive training. Additionally, ViT exhibits a larger gap between training and validation accuracy, suggesting a tendency to overfit.

The Basic CNN model aligns with expectations, as simpler networks with fewer parameters may struggle to generalize well without advanced architectural components like residual connections.

For SimCLR, the training curves indicate that pretraining learns feature representations, but fine-tuning still lags behind fully supervised models. The lower performance with 50% labeled data suggests that while SSL can reduce dependency on labels, it may require larger datasets or stronger augmentations to be more competitive in remote sensing tasks.

3.3 Confusion Matrix and Per-Class Performance

The confusion matrices B provide further insights into model performance. ResNet-50 and ViT show strong classification ability across all seven scene categories, though minor confusion occurs between visually similar classes, such as ViT model couldn’t distinguish industry and parking very well. The Basic CNN exhibits more pronounced misclassifications, especially in similar scenes like residential and industrial area, grass and field, indicating weaker feature extraction capabilities.

Interestingly, the SimCLR-based models show slightly different misclassification patterns, possibly due to the nature of the learned representations. The performance drop with reduced training data is more evident in certain classes, especially industry, suggesting that SSL pretraining does not generalize equally across all scene types.

3.4 Limitations and Future Work

While this study provides valuable insights, certain limitations should be acknowledged. First, the dataset size may not be sufficient for fully leveraging the advantages of transformer-based models and self-supervised learning. Future work could explore whether pretraining on larger remote sensing datasets improves performance.

Additionally, experimenting with alternative SSL approaches, such as MoCo or SwAV, may provide a better understanding of the role of contrastive learning in remote sensing. Further architectural improvements, such as hybrid CNN-ViT models, could also be considered to combine the strengths of both paradigms.

References

- [1] Google Drive Repository. *RSSCN7 Dataset for remote sensing scene classification*. Online. Accessed: 20-Mar-2025. 2025. URL: <https://drive.google.com/drive/folders/1A05g8Y0Nj2YZ7XdoJMA9p3rVsN40svCx>.
- [2] Ting Chen et al. “A Simple Framework for Contrastive Learning of Visual Representations”. In: *International conference on machine learning (ICML)*. 2020.

A Training and validation loss/accuracy curves

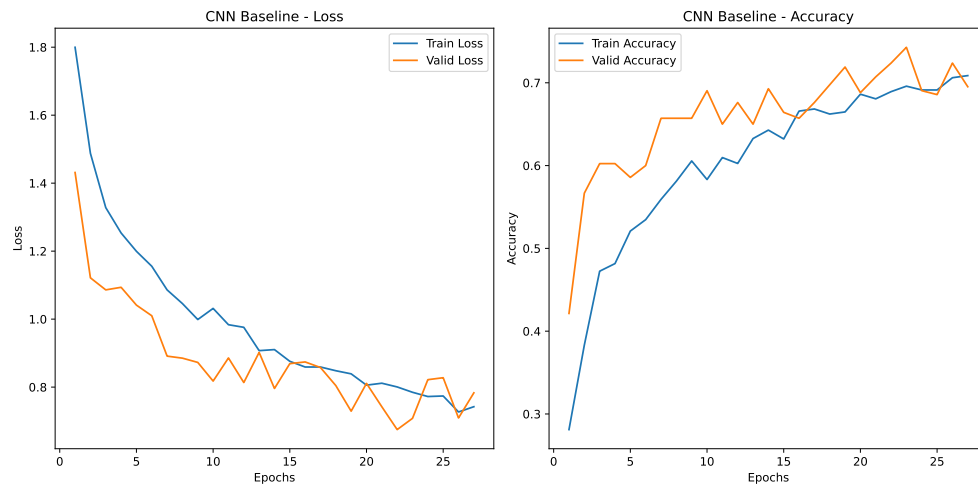


Figure 1: Training and validation loss/accuracy curves for the CNN baseline model.

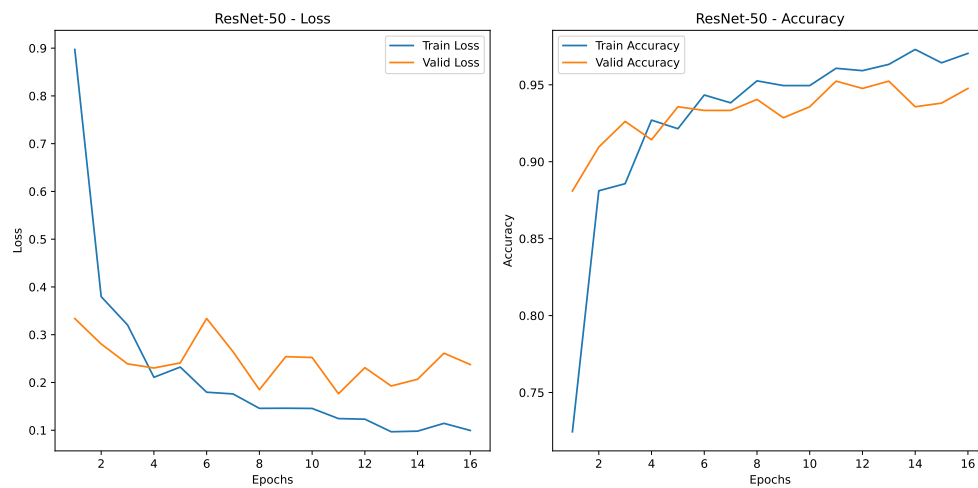


Figure 2: Training and validation loss/accuracy curves for the ResNet-50.

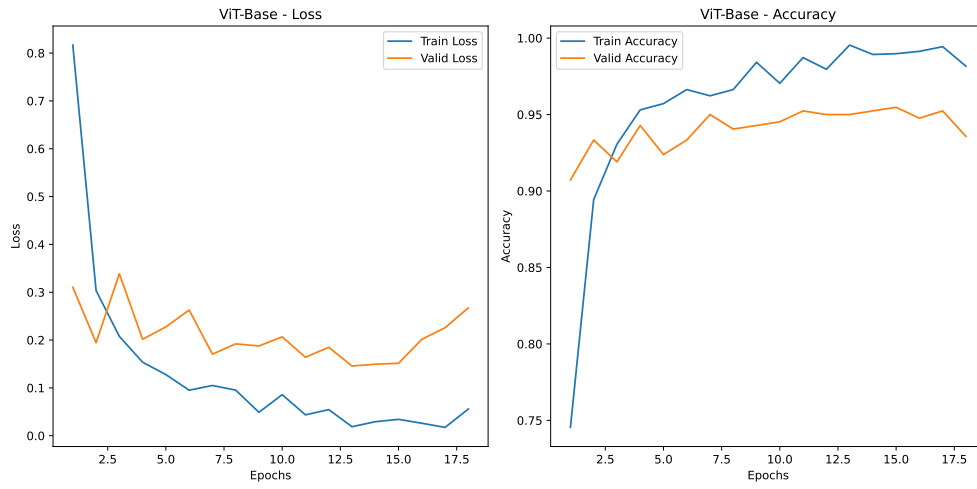


Figure 3: Training and validation loss/accuracy curves for the ViT-Base model.

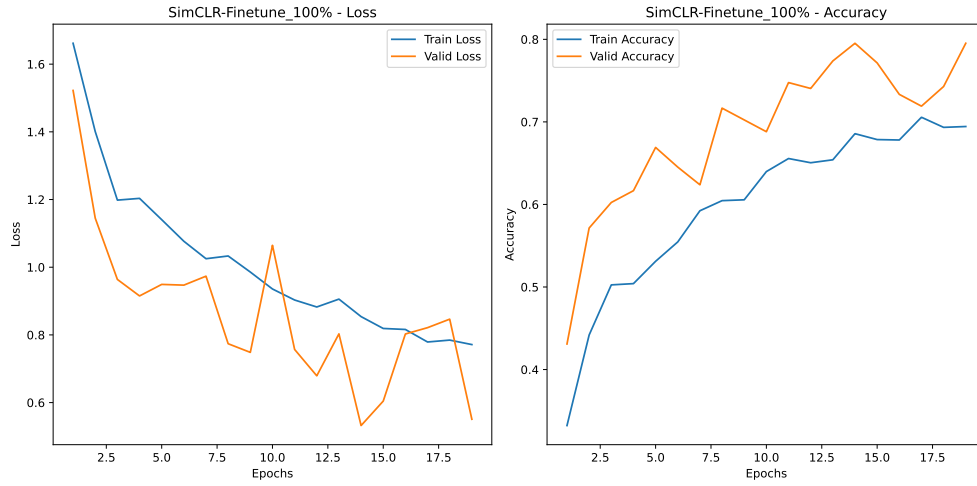


Figure 4: Training and validation loss/accuracy curves for the SimCLR-Finetune(Entry training set).

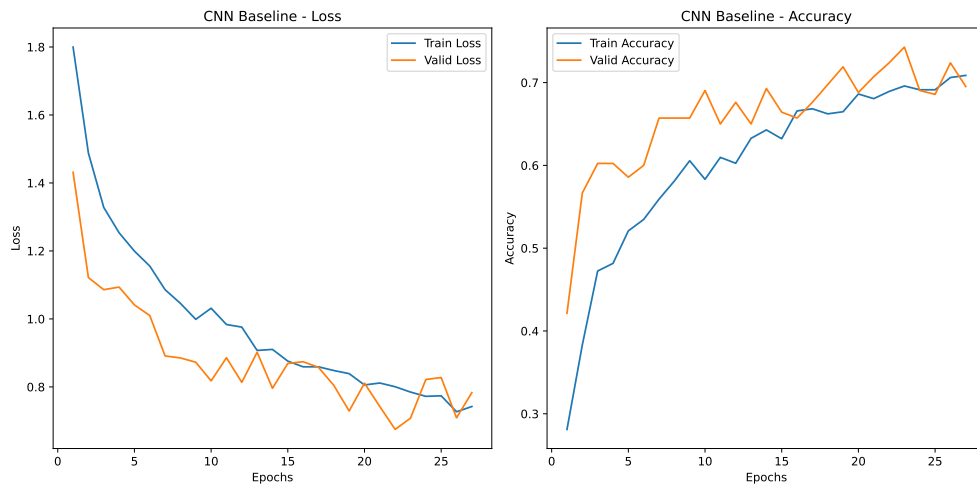


Figure 5: Training and validation loss/accuracy curves for the SimCLR-Finetune(Half of the training set).

B Confusion matrix

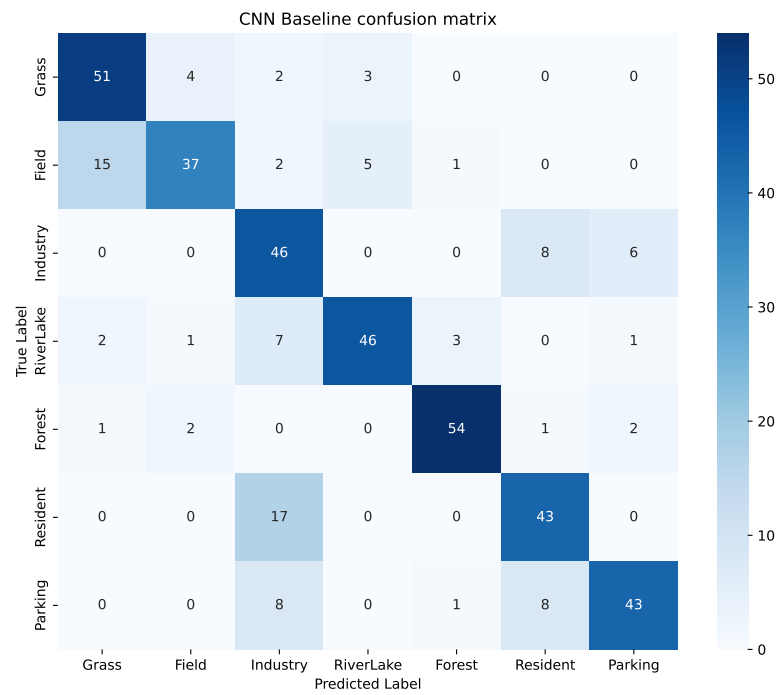


Figure 6: CNN Baseline confusion matrix

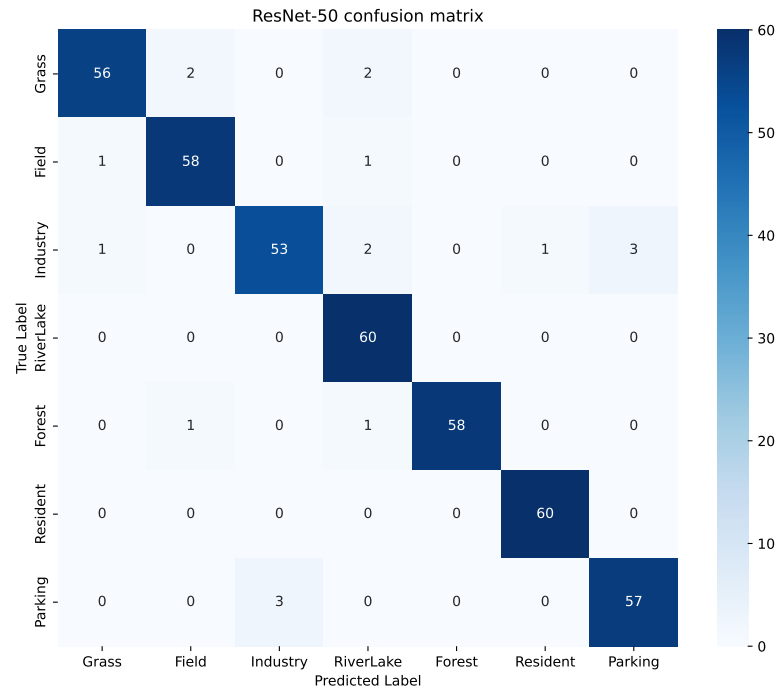


Figure 7: ResNet-50 confusion matrix

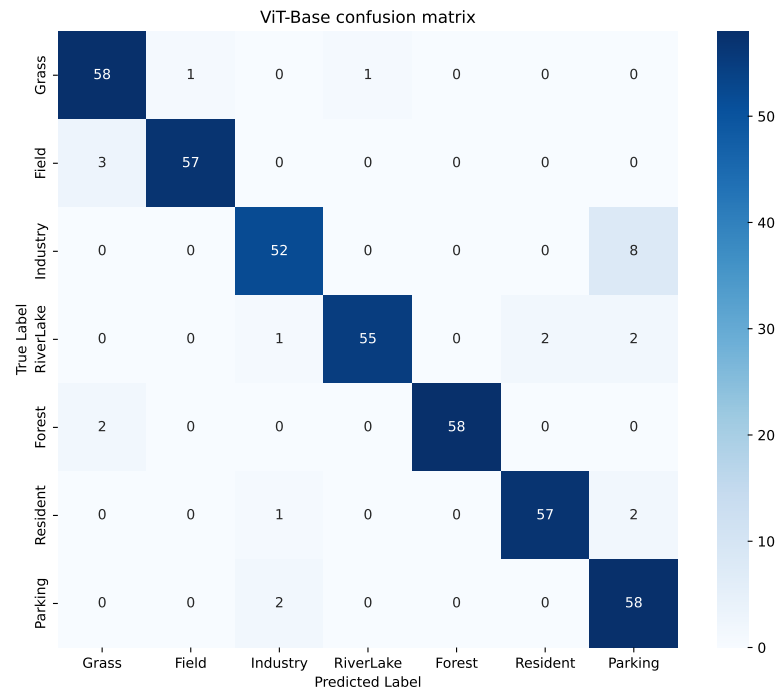


Figure 8: ViT-Base confusion matrix

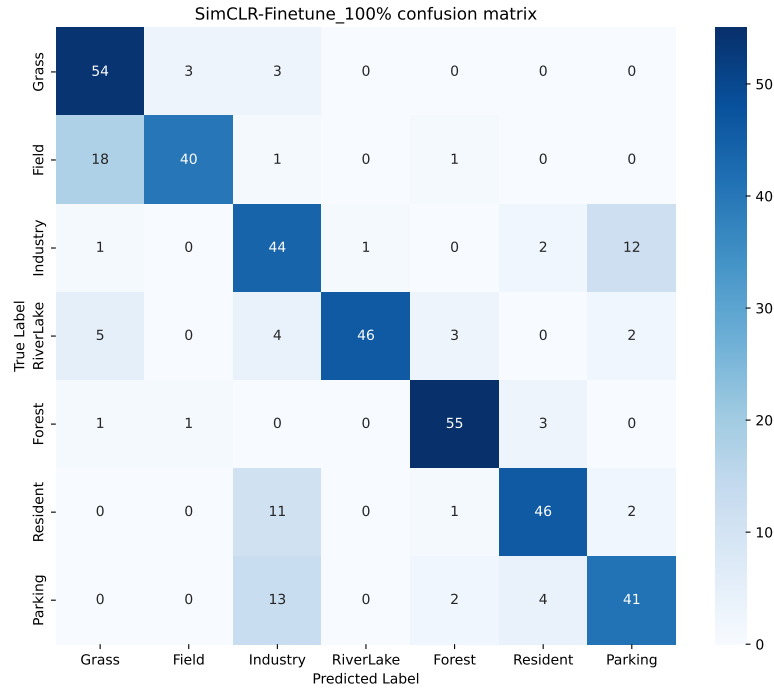


Figure 9: SimCLR-Finetune(Entrie training set) confusion matrix

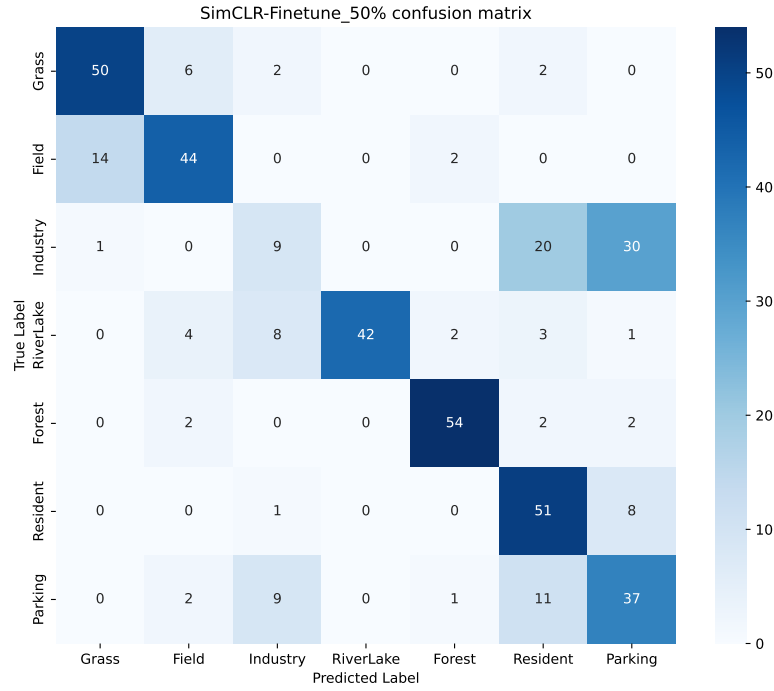


Figure 10: SimCLR-Finetune(Half of the training set) confusion matrix

C Python code

Colab code click [here](#).

C.1 Data preparation

```
1 import gdown
2
3 # Google Drive file ID
4 file_id = "1ExI4cImHyoFPL0anBBHePui5zhRWKQvj"
5 output_path = "RS_images_2800.rar"
6
7 # Use gdown to download
8 gdown.download(f"https://drive.google.com/uc?id={file_id}", output_path, quiet=False)
9
10 !apt-get install unrar
11 !unrar x RS_images_2800.rar -o+ # unzip
12
13
14 import os
15
16 dataset_path = "/content/RS_images_2800"
17 print(os.listdir(dataset_path))
18
19
20 import random
21 import pandas as pd
22
23
24 random.seed(410)
25
26 dataset_path = "./RS_images_2800"
27
28 # Split the dataset proportionally
29 train_ratio = 0.7 # Training set
30 valid_ratio = 0.15 # Validation set
31 test_ratio = 0.15 # Test set
32
33 # Read all categories
34 categories = sorted(os.listdir(dataset_path))
35 categories_cleaned = ["".join(c[1:]) for c in categories]
36 class_mapping = {c: i for i, c in enumerate(categories_cleaned)}
37
38 # Store
39 dataset = {"train": [], "valid": [], "test": []}
40
41 for category in categories:
42     category_path = os.path.join(dataset_path, category)
43     images = [os.path.join(category_path, img) for img in os.listdir(category_path)]
44
45     # Random shuffle
46     random.shuffle(images)
47
48     # Calculate partition index
49     total_count = len(images)
50     train_count = int(total_count * train_ratio)
51     valid_count = int(total_count * valid_ratio)
52
53     # Split
54     dataset["train"].extend([(img, category[1:], class_mapping[category[1:]]) for img
55                             in images[:train_count]])
```

```

55     dataset["valid"].extend([(img, category[1:], class_mapping[category[1:]]) for img
56                             in images[train_count:train_count + valid_count]])
57
58     dataset["test"].extend([(img, category[1:], class_mapping[category[1:]]) for img in
59                             images[train_count + valid_count:]]
60
61 # Convert to Pandas DataFrame
62 df_train = pd.DataFrame(dataset["train"], columns=["image", "class_name", "class_num"])
63 df_valid = pd.DataFrame(dataset["valid"], columns=["image", "class_name", "class_num"])
64 df_test = pd.DataFrame(dataset["test"], columns=["image", "class_name", "class_num"])
65
66 # Save
67 df_train.to_csv("train_data.csv", index=False)
68 df_valid.to_csv("valid_data.csv", index=False)
69 df_test.to_csv("test_data.csv", index=False)
70
71 import numpy as np
72 import torch
73 import torch.nn as nn
74 import torch.optim as optim
75 import torch.nn.functional as F
76 from torch.utils.data import Dataset, DataLoader
77 from torchvision import transforms
78 from PIL import Image
79 import matplotlib.pyplot as plt
80
81 # Set global random seed
82 seed = 410
83 random.seed(seed)
84 np.random.seed(seed)
85 torch.manual_seed(seed)
86 torch.cuda.manual_seed_all(seed)
87 torch.backends.cudnn.deterministic = True
88 torch.backends.cudnn.benchmark = False
89
90 # Use GPU
91 device = torch.device("cuda" if torch.cuda.is_available() else "CPU")
92
93 df_train = pd.read_csv("train_data.csv")
94 df_valid = pd.read_csv("valid_data.csv")
95 df_test = pd.read_csv("test_data.csv")
96
97 # Data enhancement
98 train_transform = transforms.Compose([
99     transforms.Resize((400, 400)),
100     transforms.RandomRotation(30), # Rotation
101     transforms.RandomHorizontalFlip(), # Horizontal flip
102     transforms.ColorJitter(brightness=0.5, contrast=0.5, saturation=0.5, hue=0.1), #
103     # Brightness change
104     transforms.ToTensor(),
105     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
106 ])
107
108 valid_test_transform = transforms.Compose([
109     transforms.Resize((400, 400)),
110     transforms.ToTensor(),
111     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
112 ])
113
114 class LoadDataset(Dataset):

```

```

113     def __init__(self, dataframe, transform=None):
114         self.dataframe = dataframe
115         self.transform = transform
116
117     def __len__(self):
118         return len(self.dataframe)
119
120     def __getitem__(self, idx):
121         img_path, class_name, class_num = self.dataframe.iloc[idx]
122         image = Image.open(img_path).convert("RGB")
123         if self.transform:
124             image = self.transform(image)
125         return image, class_num
126
127 # Create DataLoader
128 batch_size = 32
129
130 train_dataset = LoadDataset(df_train, transform=train_transform)
131 valid_dataset = LoadDataset(df_valid, transform=valid_test_transform)
132 test_dataset = LoadDataset(df_test, transform=valid_test_transform)
133
134 train_loader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True)
135 valid_loader = DataLoader(valid_dataset, batch_size=batch_size, shuffle=False)
136 test_loader = DataLoader(test_dataset, batch_size=batch_size, shuffle=False)

```

C.2 Basic CNN

```

1 class CNNBaseline(nn.Module):
2     def __init__(self, num_classes):
3         super(CNNBaseline, self).__init__()
4         self.conv1 = nn.Conv2d(3, 32, kernel_size=3, stride=1, padding=1)
5         self.conv2 = nn.Conv2d(32, 64, kernel_size=3, stride=1, padding=1)
6         self.conv3 = nn.Conv2d(64, 128, kernel_size=3, stride=1, padding=1)
7         self.pool = nn.MaxPool2d(kernel_size=2, stride=2, padding=0)
8         self.fc1 = nn.Linear(128 * 50 * 50, 256) #400x400 input, after 3 MaxPool
           iterations, becomes 50x50
9         self.fc2 = nn.Linear(256, num_classes)
10        self.dropout = nn.Dropout(0.5)
11
12    def forward(self, x):
13        x = self.pool(F.relu(self.conv1(x)))
14        x = self.pool(F.relu(self.conv2(x)))
15        x = self.pool(F.relu(self.conv3(x)))
16        x = torch.flatten(x, 1)
17        x = F.relu(self.fc1(x))
18        x = self.dropout(x)
19        x = self.fc2(x)
20        return x
21
22
23 num_classes = len(df_train["class_num"].unique())
24
25 model = CNNBaseline(num_classes).to(device)
26
27
28
29 import matplotlib.pyplot as plt
30
31 def plot_training_results(train_losses, valid_losses, train_accs, valid_accs,
    model_name="Model"):

```

```

32     epochs = range(1, len(train_losses) + 1)
33
34     plt.figure(figsize=(12, 6))
35
36     # Draw loss
37     plt.subplot(1, 2, 1)
38     plt.plot(epochs, train_losses, label="Train Loss")
39     plt.plot(epochs, valid_losses, label="Valid Loss")
40     plt.xlabel("Epochs")
41     plt.ylabel("Loss")
42     plt.title(f"{model_name} - Loss")
43     plt.legend()
44
45     # Draw accuracy
46     plt.subplot(1, 2, 2)
47     plt.plot(epochs, train_accs, label="Train Accuracy")
48     plt.plot(epochs, valid_accs, label="Valid Accuracy")
49     plt.xlabel("Epochs")
50     plt.ylabel("Accuracy")
51     plt.title(f"{model_name} - Accuracy")
52     plt.legend()
53
54     plt.tight_layout()
55     plt.savefig(f"{model_name}_training_results.pdf")
56     plt.show()
57
58
59
60 def train_model(model, train_loader, valid_loader, device, num_epochs=30,
61                 learning_rate=0.0001, early_stop_patience=5, model_name="Model"):
62     criterion = nn.CrossEntropyLoss()
63     optimizer = optim.Adam(model.parameters(), lr=learning_rate)
64
65     # Early stopping
66     best_val_loss = float("inf")
67     early_stop_counter = 0
68
69     # Record Loss and Accuracy
70     train_losses, valid_losses = [], []
71     train_accs, valid_accs = [], []
72
73     for epoch in range(num_epochs):
74         model.train()
75         running_loss = 0.0
76         correct = 0
77         total = 0
78
79         for images, labels in train_loader:
80             images, labels = images.to(device), labels.to(device)
81
82             optimizer.zero_grad()
83             outputs = model(images)
84             loss = criterion(outputs, labels)
85             loss.backward()
86             optimizer.step()
87
88             running_loss += loss.item()
89             _, predicted = torch.max(outputs, 1)
90             correct += (predicted == labels).sum().item()
91             total += labels.size(0)

```

```

92     train_loss = running_loss / len(train_loader)
93     train_acc = correct / total
94     train_losses.append(train_loss)
95     train_accs.append(train_acc)
96
97     # Validation
98     model.eval()
99     running_val_loss = 0.0
100     correct = 0
101     total = 0
102
103     with torch.no_grad():
104         for images, labels in valid_loader:
105             images, labels = images.to(device), labels.to(device)
106             outputs = model(images)
107             loss = criterion(outputs, labels)
108             running_val_loss += loss.item()
109
110             _, predicted = torch.max(outputs, 1)
111             correct += (predicted == labels).sum().item()
112             total += labels.size(0)
113
114     val_loss = running_val_loss / len(valid_loader)
115     val_acc = correct / total
116     valid_losses.append(val_loss)
117     valid_accs.append(val_acc)
118
119     print(f"Epoch {epoch+1}/{num_epochs}: Train Loss: {train_loss:.4f}, Train Acc:
120           {train_acc:.4f}, Val Loss: {val_loss:.4f}, Val Acc: {val_acc:.4f}")
121
122     # Early stopping & Save best model
123     if val_loss < best_val_loss:
124         best_val_loss = val_loss
125         early_stop_counter = 0
126         torch.save(model.state_dict(), f"{model_name}_best.pth")
127         print("Best model saved!")
128     else:
129         early_stop_counter += 1
130         if early_stop_counter >= early_stop_patience:
131             print("Early stopping triggered!")
132             break
133
134     # plot loss and acc
135     plot_training_results(train_losses, valid_losses, train_accs, valid_accs,
136                           model_name)
137
138 train_model(model, train_loader, valid_loader, device, num_epochs=30, model_name="CNN
139           Baseline")
140
141 def evaluate_model(model, test_loader, device):
142     model.eval()
143     correct = 0
144     total = 0
145     all_predictions = []
146     all_labels = []
147
148     with torch.no_grad():
149         for images, labels in test_loader:
150             images, labels = images.to(device), labels.to(device)
151             outputs = model(images)

```

```

150         _, predicted = torch.max(outputs, 1)
151
152         all_predictions.extend(predicted.cpu().numpy())
153         all_labels.extend(labels.cpu().numpy())
154
155         correct += (predicted == labels).sum().item()
156         total += labels.size(0)
157
158     accuracy = correct / total
159     print(f"Test Accuracy: {accuracy:.4f}")
160
161     return accuracy, all_predictions, all_labels
162
163
164 # Load model
165 best_model = CNNBaseline(num_classes).to(device)
166 best_model.load_state_dict(torch.load("best_model.pth"))
167
168 # Compute test accuracy
169 test_accuracy, predictions, labels = evaluate_model(best_model, test_loader, device)
170
171
172 import seaborn as sns
173 from sklearn.metrics import confusion_matrix
174
175 def plot_confusion_matrix(y_true, y_pred, class_names, model_name):
176     cm = confusion_matrix(y_true, y_pred)
177     plt.figure(figsize=(10, 8))
178     sns.heatmap(cm, annot=True, fmt="d", cmap="Blues", xticklabels=class_names,
179                yticklabels=class_names)
180     plt.xlabel("Predicted Label")
181     plt.ylabel("True Label")
182     plt.title(f"{model_name} confusion matrix")
183     plt.savefig(f"{model_name}_confusion_matrix.pdf")
184     plt.show()
185
186 # Draw confusion matrix
187 plot_confusion_matrix(labels, predictions, df_train["class_name"].unique(),
188                       model_name="CNN Baseline")

```

C.3 ResNet-50

```

1 from torchvision import models
2
3 # Load ResNet-50 pre trained model
4 model = models.resnet50(weights=models.ResNet50_Weights.IMAGENET1K_V1)
5
6 # Freeze the first 140 layers (only train the last few layers)
7 for param in list(model.parameters())[:140]:
8     param.requires_grad = False
9
10 # Replace the fully connected layer of ResNet-50
11 num_features = model.fc.in_features
12 model.fc = nn.Sequential(
13     nn.Linear(num_features, 256),
14     nn.ReLU(),
15     nn.Dropout(0.3),
16     nn.Linear(256, num_classes)
17 )
18

```



```

19 model = model.to(device)
20
21 train_model(model, train_loader, valid_loader, device, num_epochs=30,
    model_name="ResNet-50")
22
23 # Load best model
24 best_resnet_model = models.resnet50(weights=models.ResNet50_Weights.IMAGENET1K_V1)
25 best_resnet_model.fc = nn.Sequential(
26     nn.Linear(num_features, 256),
27     nn.ReLU(),
28     nn.Dropout(0.3),
29     nn.Linear(256, num_classes)
30 )
31 best_resnet_model.load_state_dict(torch.load("ResNet-50_best.pth"))
32 best_resnet_model.to(device)
33
34 # Compute test accuracy
35 test_accuracy, predictions, labels = evaluate_model(best_resnet_model, test_loader,
    device)
36
37 # Deaw confusion matrix
38 plot_confusion_matrix(labels, predictions, df_train["class_name"].unique(),
    model_name="ResNet-50")

```

C.4 ViT

```

1 from transformers import ViTModel, ViTFeatureExtractor
2
3 # Load ViT feature extractor
4 feature_extractor = ViTFeatureExtractor.from_pretrained("google/vit-base-patch16-224")
5 # Load ViT model (without classification head)
6 class ViTClassifier(nn.Module):
7     def __init__(self, num_classes=7):
8         super(ViTClassifier, self).__init__()
9         self.vit = ViTModel.from_pretrained("google/vit-base-patch16-224")
10        self.classifier = nn.Sequential(
11            nn.Linear(self.vit.config.hidden_size, 256),
12            nn.ReLU(),
13            nn.Dropout(0.3),
14            nn.Linear(256, num_classes)
15        )
16
17    def forward(self, x):
18        outputs = self.vit(x) # Extract features
19        cls_token = outputs.last_hidden_state[:, 0, :] # CLS token
20        return self.classifier(cls_token)
21
22 # Initialize model
23 model = ViTClassifier(num_classes=7).to(device)
24
25
26 train_transform_vit = transforms.Compose([
27     transforms.Resize((224, 224)), # Vit needs 224x224 input
28     transforms.RandomRotation(30),
29     transforms.RandomHorizontalFlip(),
30     transforms.ColorJitter(brightness=0.5, contrast=0.5, saturation=0.5, hue=0.1),
31     transforms.ToTensor(),
32     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
33 ])
34

```

```

35 valid_test_transform_vit = transforms.Compose([
36     transforms.Resize((224, 224)),
37     transforms.ToTensor(),
38     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
39 ])
40
41 class ViTDataset(Dataset):
42     def __init__(self, dataframe, transform):
43         self.dataframe = dataframe
44         self.transform = transform
45
46     def __len__(self):
47         return len(self.dataframe)
48
49     def __getitem__(self, idx):
50         img_path, class_num = self.dataframe.iloc[idx][["image", "class_num"]]
51
52         image = Image.open(img_path).convert("RGB")
53
54         image = self.transform(image)
55
56         # Ensure the shape is (3, 224, 224)
57         if image.shape != (3, 224, 224):
58             print(f"Shape mismatch: {image.shape} at index {idx}")
59
60         # Convert class_num to PyTorch tensor
61         class_num = torch.tensor(class_num, dtype=torch.long)
62
63         return image, class_num
64
65
66 train_dataset_vit = ViTDataset(df_train, transform=train_transform_vit)
67 valid_dataset_vit = ViTDataset(df_valid, transform=valid_test_transform_vit)
68 test_dataset_vit = ViTDataset(df_test, transform=valid_test_transform_vit)
69
70 train_loader_vit = DataLoader(train_dataset_vit, batch_size=batch_size, shuffle=True)
71 valid_loader_vit = DataLoader(valid_dataset_vit, batch_size=batch_size, shuffle=False)
72 test_loader_vit = DataLoader(test_dataset_vit, batch_size=batch_size, shuffle=False)
73
74
75 train_model(model, train_loader_vit, valid_loader_vit, device, num_epochs=30,
76             model_name="ViT-Base")
77
78 # Load best model
79 best_vit_model = ViTClassifier(num_classes=7).to(device)
80 best_vit_model.load_state_dict(torch.load("ViT-Base_best.pth"))
81
82 # Compute test accuracy
83 test_accuracy, predictions, labels = evaluate_model(best_vit_model, test_loader_vit,
84                                                     device)
85
86 # Draw confusion matrix
87 plot_confusion_matrix(labels, predictions, df_train["class_name"].unique(),
88                       model_name="ViT-Base")

```

C.5 SSL

```

1 contrastive_transforms = transforms.Compose([
2     transforms.RandomResizedCrop(400, scale=(0.2, 1.0)),

```

```

3     transforms.RandomHorizontalFlip(),
4     transforms.RandomApply([transforms.ColorJitter(0.4, 0.4, 0.4, 0.1)], p=0.8),
5     transforms.RandomGrayscale(p=0.2),
6     transforms.ToTensor(),
7     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
8 ])
9
10 def get_augmented_views(img_path):
11     img = Image.open(img_path).convert("RGB")
12     view1 = contrastive_transforms(img)
13     view2 = contrastive_transforms(img)
14     return view1, view2
15
16 class SimCLRDataset(Dataset):
17     def __init__(self, df):
18         self.image_paths = df["image"].values
19
20     def __len__(self):
21         return len(self.image_paths)
22
23     def __getitem__(self, idx):
24         img_path = self.image_paths[idx]
25         img1, img2 = get_augmented_views(img_path)
26         return img1, img2
27
28 train_dataset_sim = SimCLRDataset(df_train)
29 train_loader_sim = DataLoader(train_dataset, batch_size=32, shuffle=True, num_workers=2)
30
31
32 from torchvision.models import ResNet50_Weights
33
34 class SimCLR(nn.Module):
35     def __init__(self, feature_dim=128):
36         super(SimCLR, self).__init__()
37         base_model = models.resnet50(weights=ResNet50_Weights.DEFAULT)
38
39         # Remove ResNet-50 Classification layer
40         self.encoder = nn.Sequential(*list(base_model.children())[:-1])
41
42         # Projection Head
43         self.projector = nn.Sequential(
44             nn.Linear(2048, 512),
45             nn.ReLU(),
46             nn.Linear(512, feature_dim)
47         )
48
49     def forward(self, x):
50         h = self.encoder(x).squeeze()
51         z = self.projector(h)
52         return z
53
54
55 model = SimCLR().cuda()
56
57
58 def nt_xent_loss(z1, z2, temperature=0.3):
59     z1 = F.normalize(z1, dim=1)
60     z2 = F.normalize(z2, dim=1)
61     logits = torch.matmul(z1, z2.T) / temperature
62     labels = torch.arange(z1.shape[0]).cuda()
63     loss = F.cross_entropy(logits, labels)

```

```

64     return loss
65
66
67 import gc
68 import time
69 import json
70 from torch.amp import autocast, GradScaler
71
72 # Record time
73 start_time = time.time()
74
75 # Record training history
76 history = {"train_loss": [], "learning_rate": []}
77
78
79 scaler = GradScaler(device="cuda")
80 optimizer = torch.optim.Adam(model.parameters(), lr=1e-4)
81 scheduler = torch.optim.lr_scheduler.StepLR(optimizer, step_size=50, gamma=0.5)
82
83 best_train_loss = float("inf")
84 patience = 5
85 early_stop_counter = 0
86
87
88 for epoch in range(30):
89     total_loss = 0
90     model.train()
91
92     epoch_start = time.time()
93
94     for img1, img2 in train_loader:
95         img1, img2 = img1.cuda(), img2.cuda()
96
97         optimizer.zero_grad()
98
99         with autocast(device_type="cuda"):
100             z1, z2 = model(img1), model(img2)
101             loss = nt_xent_loss(z1, z2)
102
103         scaler.scale(loss).backward()
104         scaler.step(optimizer)
105         scaler.update()
106
107         total_loss += loss.item()
108
109         del img1, img2, z1, z2, loss
110         torch.cuda.empty_cache()
111         gc.collect()
112
113     avg_train_loss = total_loss / len(train_loader)
114     history["train_loss"].append(avg_train_loss)
115     history["learning_rate"].append(optimizer.param_groups[0]['lr'])
116
117     # update learning rate
118     scheduler.step()
119
120     # Compute training time
121     epoch_time = time.time() - epoch_start
122
123     print(f"Epoch [{epoch+1}/30] | Train Loss: {avg_train_loss:.4f} | LR:
          {optimizer.param_groups[0]['lr']:.6f} | Time: {epoch_time:.2f}s")

```

```

124
125     # Save best model
126     if avg_train_loss < best_train_loss:
127         best_train_loss = avg_train_loss
128         torch.save(model.encoder.state_dict(), "simclr_best_model.pth")
129         early_stop_counter = 0
130     else:
131         early_stop_counter += 1
132
133     # Early Stopping
134     if early_stop_counter >= patience:
135         break
136
137
138     =====
139     # Fine-tuning
140     =====
141     def balanced_sample(df, frac=0.5, random_state=410):
142         return df.groupby("class_name", group_keys=False).apply(lambda x:
143             x.sample(frac=frac, random_state=random_state)).reset_index(drop=True)
144
145     df_train_sample = balanced_sample(df_train)
146
147     # Check
148     print(df_train["class_name"].value_counts())
149     print(df_train_sample["class_name"].value_counts())
150
151     train_dataset_full = LoadDataset(df_train, transform=train_transform)
152     train_dataset_sample = LoadDataset(df_train_sample, transform=train_transform)
153     valid_dataset = LoadDataset(df_valid, transform=valid_test_transform)
154     test_dataset = LoadDataset(df_test, transform=valid_test_transform)
155
156     train_loader_full = DataLoader(train_dataset_full, batch_size=batch_size, shuffle=True)
157     train_loader_sample = DataLoader(train_dataset_sample, batch_size=batch_size,
158         shuffle=True)
159     valid_loader = DataLoader(valid_dataset, batch_size=batch_size, shuffle=False)
160     test_loader = DataLoader(test_dataset, batch_size=batch_size, shuffle=False)
161
162
163     # Load ResNet-50 pre trained with SimCLR
164     simclr_model = models.resnet50()
165     num_features = simclr_model.fc.in_features
166     simclr_model.fc = nn.Identity() # Remove the projection head
167     simclr_model.load_state_dict(torch.load("simclr_best_model.pth"), strict=False)
168     simclr_model = simclr_model.cuda()
169
170
171     # Add classification head
172     simclr_model.fc = nn.Sequential(
173         nn.Linear(num_features, 256),
174         nn.ReLU(),
175         nn.Dropout(0.3),
176         nn.Linear(256, num_classes)
177     )
178
179     simclr_model = simclr_model.to(device)
180
181     # Full training set

```

```

182 train_model(simclr_model, train_loader_full, valid_loader, device, num_epochs=30,
    model_name="SimCLR-Finetune_100%")
183
184
185 # Load the best model
186 best_simclr_model = models.resnet50()
187 best_simclr_model.fc = nn.Sequential(
188     nn.Linear(num_features, 256),
189     nn.ReLU(),
190     nn.Dropout(0.3),
191     nn.Linear(256, num_classes)
192 )
193 best_simclr_model.load_state_dict(torch.load("SimCLR-Finetune_100%_best.pth"))
194 best_simclr_model.to(device)
195
196 # Compute test accuracy
197 test_accuracy, predictions, labels = evaluate_model(best_simclr_model, test_loader,
    device)
198
199 # Draw confusion matrix
200 plot_confusion_matrix(labels, predictions, df_train["class_name"].unique(),
    model_name="SimCLR-Finetune_100%")
201
202 #-----
203 # Load ResNet-50 pre trained with SimCLR
204 simclr_model = models.resnet50()
205 num_features = simclr_model.fc.in_features
206 simclr_model.fc = nn.Identity() # Remove the projection head
207 simclr_model.load_state_dict(torch.load("simclr_best_model.pth"), strict=False)
208 simclr_model = simclr_model.cuda()
209
210 # Add classification head
211 simclr_model.fc = nn.Sequential(
212     nn.Linear(num_features, 256),
213     nn.ReLU(),
214     nn.Dropout(0.3),
215     nn.Linear(256, num_classes)
216 )
217
218 simclr_model = simclr_model.to(device)
219
220 # Half of training set
221 train_model(simclr_model, train_loader_sample, valid_loader, device, num_epochs=30,
    model_name="SimCLR-Finetune_50%")
222
223
224 # Load the best model
225 best_simclr_model = models.resnet50()
226 best_simclr_model.fc = nn.Sequential(
227     nn.Linear(num_features, 256),
228     nn.ReLU(),
229     nn.Dropout(0.3),
230     nn.Linear(256, num_classes)
231 )
232 best_simclr_model.load_state_dict(torch.load("SimCLR-Finetune_50%_best.pth"))
233 best_simclr_model.to(device)
234
235 # Compute test accuracy
236 test_accuracy, predictions, labels = evaluate_model(best_simclr_model, test_loader,
    device)
237

```

```
238 # Draw confusion matrix
239 plot_confusion_matrix(labels, predictions, df_train["class_name"].unique(),
                        model_name="SimCLR-Finetune_50%")
```