

# Deep Learning for Medical Image Classification

---

This notebook builds a full image classification pipeline for the PathMNIST dataset. We start with exploratory checks and baseline models, then train a CNN with tuned hyperparameters and evaluate on the test set.

## Contents

1. **Environment Setup & Data Loading** – reproducibility seeds, GPU check, NPZ import
  2. **Dataset Exploration**
    - 2.1 Structure & dimensionality
    - 2.2 Pixel-range verification and duplicate scan
    - 2.3 Class-balance summary
  3. **Tabular Baselines**
    - 3.1 Logistic-Regression (flattened greyscale)
    - 3.2 Random-Forest (flattened RGB)
  4. **Baseline Performance Diagnostics** – accuracy, macro-F1, confusion matrices
  5. **CNN Data Preparation** – transforms, stratified train/val split, loaders
  6. **CNN Architecture & Training Loop** – 3-conv-layer network, early-stopping
  7. **Hyper-parameter Tuning** – grid search over LR × batch × dropout
  8. **Final Test-set Evaluation** – classification report & confusion matrix
  9. **Key Findings & Comparison Table** – headline metrics and insights
- 

Undertake appropriate pre-processing of the data.

---

## 1. Environment Setup & Data Loading

Import the required libraries, set random seeds, and load the PathMNIST NPZ file:

```
import numpy as np
import matplotlib.pyplot as plt
import torch
torch.set_num_threads(2) # this makes it run better on elab
import torch.nn as nn
```

```

from torch.utils.data import Dataset, DataLoader
from sklearn.linear_model import LogisticRegression
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report, accuracy_score,
confusion_matrix
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.pipeline import Pipeline
from torch import optim
import seaborn as sns

# this allows pytorch to use a GPU, if one is available
device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
if torch.cuda.is_available():
    print("GPU name:", torch.cuda.get_device_name(0))

GPU name: Tesla T4

# SET REPRODUCIBILITY SEEDS (works for CPU and, if present, GPU)
import os, random, numpy as np, torch

SEED = 42
os.environ["PYTHONHASHSEED"] = str(SEED) # makes Python hashing
repeatable
random.seed(SEED) # built-in RNG
np.random.seed(SEED) # NumPy RNG
torch.manual_seed(SEED) # PyTorch CPU RNG

if torch.cuda.is_available(): # extra safety when a GPU
is used
    torch.cuda.manual_seed_all(SEED) # PyTorch GPU RNG
    torch.backends.cudnn.deterministic = True
    torch.backends.cudnn.benchmark = False

data = np.load('Assignment2Data.npz')
print("Data files available:", data.files)

Data files available: ['train_images', 'train_labels', 'test_images',
'test_labels']

```

Assign each of the training, validation and test sets to their own array:

```

train_data = data['train_images']
train_labels = data['train_labels']
test_data = data['test_images']
test_labels = data['test_labels']

```

---

## 2. Dataset Exploration

**Purpose:** Understand dataset structure, class distribution, and verify data quality.

**Dataset:** PathMNIST with 5 tissue types from histological images.

```
# Define tissue type labels
label_dict = {
    0: "Adipose tissue",
    1: "Background",
    2: "Debris",
    3: "Lymphocytes",
    4: "Mucus"
}
print("Tissue type labels:")
for key, value in label_dict.items():
    print(f" {key}: {value}")
```

```
Tissue type labels:
0: Adipose tissue
1: Background
2: Debris
3: Lymphocytes
4: Mucus
```

### 2.1 Dataset Structure

**Data Structure:** Verify image dimensions, label format, and class distribution.

```
# Analyze dataset structure
print(f"Training data shape: {train_data.shape}")
print(f"Training labels shape: {train_labels.shape}")
print(f"Test data shape: {test_data.shape}")
print(f"Test labels shape: {test_labels.shape}")

# Verify data types and ranges
print(f"\nImage pixel range: {train_data.min():.1f} to {train_data.max():.1f}")
print(f"Label range: {train_labels.min()} to {train_labels.max()}")
print(f"Number of classes: {len(np.unique(train_labels))}")

# Check for any missing values
print(f"Missing values in images: {np.isnan(train_data).sum()}")
print(f"Missing values in labels: {np.isnan(train_labels).sum()}")

Training data shape: (55490, 28, 28, 3)
Training labels shape: (55490, 1)
Test data shape: (4367, 28, 28, 3)
Test labels shape: (4367, 1)
```

```
Image pixel range: 0.0 to 254.0
Label range: 0 to 4
Number of classes: 5
Missing values in images: 0
Missing values in labels: 0
```

## 2.2 Verifying Pixel Value Ranges

It's important to verify that the pixel values of our images are within the expected range (0 to 255 for 8-bit images). This ensures that our normalization steps are accurate.

```
# Verify pixel value ranges in the training data
print("Training data pixel value range:")
print(f"Minimum pixel value: {train_data.min()}")
print(f"Maximum pixel value: {train_data.max()}")

# Verify pixel value ranges in the test data
print("\nTest data pixel value range:")
print(f"Minimum pixel value: {test_data.min()}")
print(f"Maximum pixel value: {test_data.max()}")
```

```
Training data pixel value range:
Minimum pixel value: 0
Maximum pixel value: 254
```

```
Test data pixel value range:
Minimum pixel value: 5
Maximum pixel value: 250
```

## 2.3 Class Distribution

**Class Balance:** Check distribution of tissue types to identify potential bias.

```
# Analyze class distribution
train_unique, train_counts = np.unique(train_labels,
return_counts=True)
test_unique, test_counts = np.unique(test_labels, return_counts=True)

print("CLASS DISTRIBUTION ANALYSIS:")
print("="*50)
print(f"{'Tissue Type':<15} {'Train Count':<12} {'Train %':<10} {'Test
Count':<11} {'Test %':<8}")
print("-"*50)

for i in range(len(train_unique)):
    train_pct = (train_counts[i]/len(train_labels))*100
    test_pct = (test_counts[i]/len(test_labels))*100
    tissue_name = label_dict[train_unique[i]]
    print(f"{'tissue_name':<15} {'train_counts[i]:<12} {'train_pct:<10.1f}
```

```
{test_counts[i]:<11} {test_pct:<8.1f}")

print("-"*50)
print(f"{'TOTAL':<15} {len(train_labels):<12} {'100.0':<10}
{len(test_labels):<11} {'100.0':<8}")

# Check for severe class imbalance
max_count = np.max(train_counts)
min_count = np.min(train_counts)
imbalance_ratio = max_count / min_count
print(f"\nClass imbalance ratio: {imbalance_ratio:.1f}:1")
if imbalance_ratio > 3:
    print("Class imbalance detected - may need balancing strategies")
else:
    print("Reasonable class balance")

CLASS DISTRIBUTION ANALYSIS:
=====
Tissue Type      Train Count  Train %    Test Count  Test %
-----
Adipose tissue   10407        18.8       1338        30.6
Background       11557        20.8        634        14.5
Debris           8763         15.8        741        17.0
Lymphocytes      10446        18.8        421         9.6
Mucus            14317        25.8       1233       28.2
-----
TOTAL            55490        100.0      4367       100.0

Class imbalance ratio: 1.6:1
Reasonable class balance
```

**Train/Test Distribution:** Some classes show different proportions between training and test sets:

- **Lymphocytes:** 18.8% (train) vs 9.6% (test) - underrepresented in test
- **Adipose tissue:** 18.8% (train) vs 30.6% (test) - overrepresented in test

*This reflects the original dataset split and may affect generalization patterns.*

## 2.4 Visual Exploration

**Sample Visualization:** Display example images from each tissue type.

```
# Show example of first image
label1 = train_labels[0]
print(f"First image label: {label_dict[label1[0]]}")

# Display sample images from each class
fig, axes = plt.subplots(1, 5, figsize=(15, 3))
fig.suptitle('Sample Images from Each Tissue Type', fontsize=14)
```

```

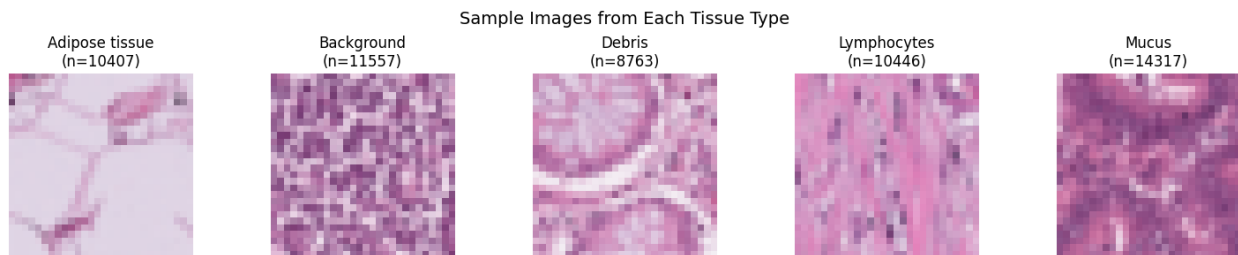
for class_idx in range(5):
    # Find first occurrence of each class
    class_indices = np.where(train_labels.ravel() == class_idx)[0]
    if len(class_indices) > 0:
        sample_idx = class_indices[0]
        sample_image = train_data[sample_idx]

        axes[class_idx].imshow(sample_image)
        axes[class_idx].set_title(f"{label_dict[class_idx]}\n
n(n={train_counts[class_idx]})")
        axes[class_idx].axis('off')

plt.tight_layout()
plt.show()

```

First image label: Adipose tissue

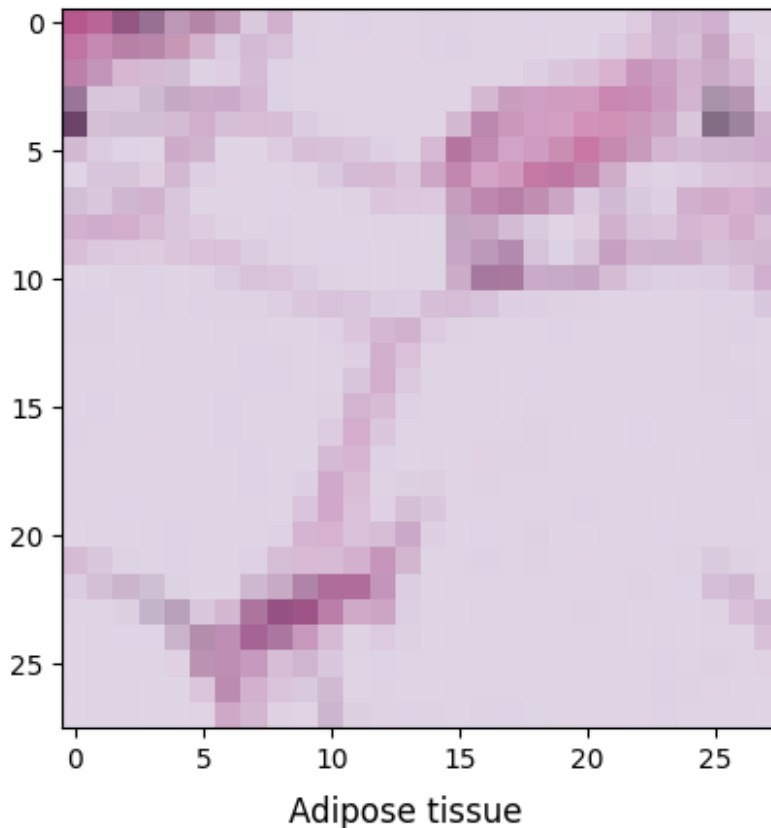


Displayed one RGB tile (and its label) to eyeball image quality and colour channels:

```

image1 = train_data[0,:,:,:]
plt.imshow(image1)
plt.figtext(0.5, 0.01, label_dict[label1[0]], wrap=True,
horizontalalignment='center', fontsize=12)
plt.show()

```



We can look at each channel individually, by plotting one of the layers:

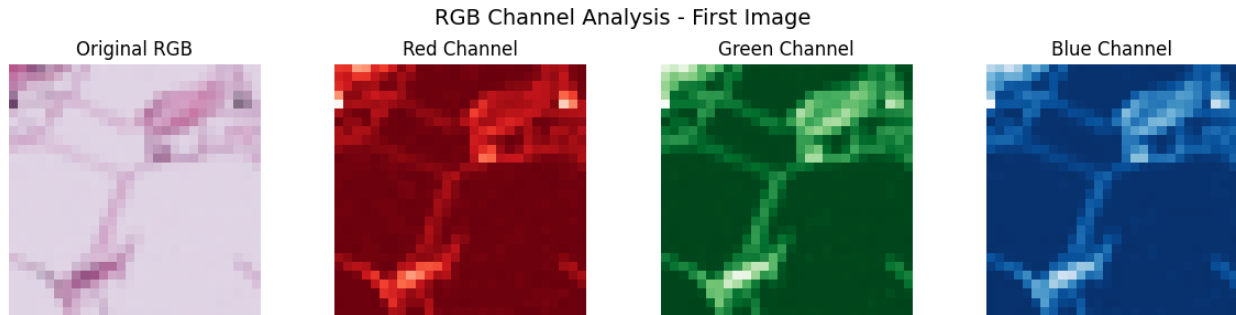
```
# View individual RGB channels
fig, axes = plt.subplots(1, 4, figsize=(12, 3))
fig.suptitle('RGB Channel Analysis - First Image', fontsize=14)

# Original RGB image
axes[0].imshow(train_data[0,:,:,:])
axes[0].set_title('Original RGB')
axes[0].axis('off')

# Individual channels
channel_names = ['Red', 'Green', 'Blue']
channel_colors = ['Reds', 'Greens', 'Blues']

for i in range(3):
    axes[i+1].imshow(train_data[0,:,:,:i], cmap=channel_colors[i])
    axes[i+1].set_title(f'{channel_names[i]} Channel')
    axes[i+1].axis('off')

plt.tight_layout()
plt.show()
```



## 2.5 Checking for Duplicate Images

We can check for duplicate images in the dataset, which might affect the training process.

```
# Reshape images to 1D vectors for comparison
train_resaped = train_data.reshape(len(train_data), -1)

# Use NumPy to find unique rows (images)
unique_images = np.unique(train_resaped, axis=0)

# Count duplicates
num_duplicates = len(train_resaped) - len(unique_images)
print(f"Number of duplicate images in the training set:
{num_duplicates}")
```

Number of duplicate images in the training set: 0

### Dataset Summary

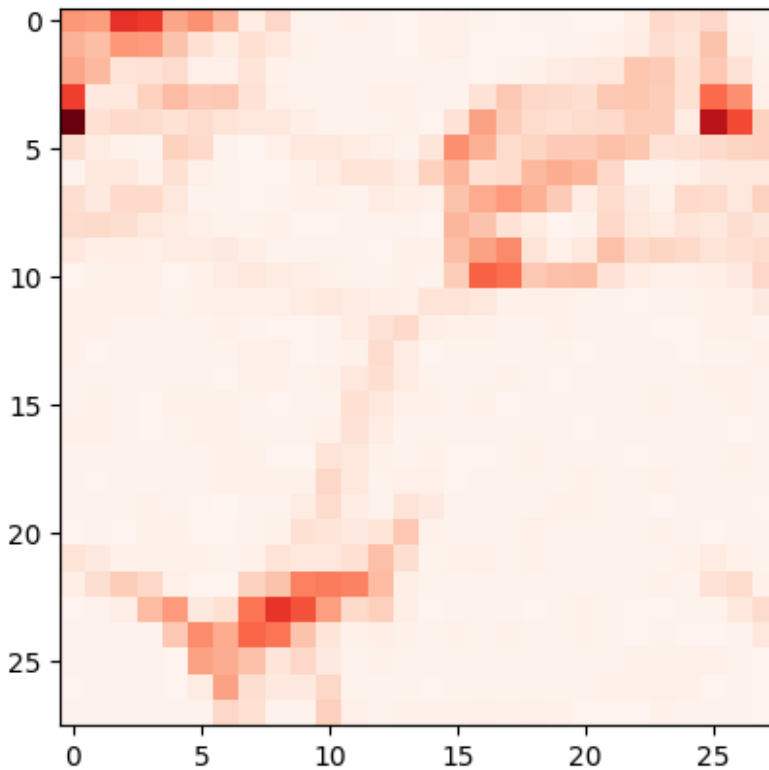
- 55,490 training images, 4,367 test images
- 5 tissue classes with reasonable balance
- RGB images (28 × 28 × 3 pixels)
- Pixel values already lie in the expected 8-bit range (train 0–254, test 5–250)
- No duplicate images detected in the training set
- No missing values detected
- Ready for preprocessing and modeling

## 3. Prep for Baseline Models

Convert every image to grayscale and flattened it to 784 features so scikit-learn can handle the data.

```
plt.imshow(train_data[0, :, :, 0], cmap='Reds_r')
<matplotlib.image.AxesImage at 0x7f6741d25350>
```





```
# use the sum of channels of each image, and reshape the data into a
1-D array (instead of the 2D image array)
train_1d = np.empty((len(train_data), 28*28))
for i in range(len(train_data)):
    train_1d[i,:] = np.reshape(train_data[i,:,:,:0] +
    train_data[i,:,:,:1] + train_data[i,:,:,:2] , -1)

y = train_labels.ravel()
```

## 4. Baseline Models

### 4.1 Logistic-Regression Baseline

Split the data 80 / 20, scaled the features, and ran a multinomial logistic-regression with 5-fold CV to pick the best C.

```
# Create train/validation split for baseline evaluation
X_train, X_val, y_train, y_val = train_test_split(train_1d, y,
test_size=0.2, random_state=42, stratify=y)
print(f"Training samples: {len(X_train):,}")
print(f"Validation samples: {len(X_val):,}")
```

```

# Implement baseline logistic regression model
logreg_pipeline = Pipeline([
    ('scaler', StandardScaler()),
    ('logistic', LogisticRegression(max_iter=1000,
multi_class='multinomial', random_state=42))
])

# Hyperparameter tuning using cross-validation
param_grid = {'logistic_C': [0.01, 0.1, 1.0, 10.0]}
grid_search = GridSearchCV(logreg_pipeline, param_grid, cv=5,
scoring='accuracy')
print("Training logistic regression baseline...")
grid_search.fit(X_train, y_train)

# Evaluate on validation set
best_baseline = grid_search.best_estimator_
y_val_pred = best_baseline.predict(X_val)
baseline_accuracy = accuracy_score(y_val, y_val_pred)
print(f"Best C parameter: {grid_search.best_params_['logistic_C']}")
print(f"Validation accuracy: {baseline_accuracy:.4f}")
print("\nClassification Report:")
print(classification_report(y_val, y_val_pred,
target_names=[label_dict[i] for i in range(5)]))

# Confusion matrix
cm = confusion_matrix(y_val, y_val_pred)
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
xticklabels=[f'Class {i}' for i in range(5)],
yticklabels=[f'Class {i}' for i in range(5)])
plt.title('Logistic Regression Confusion Matrix')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.show()
print(f"Baseline accuracy to beat with CNN: {baseline_accuracy:.4f}")

Training samples: 44,392
Validation samples: 11,098
Training logistic regression baseline...

/usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_logistic.py:1247: FutureWarning: 'multi_class' was deprecated in
version 1.5 and will be removed in 1.7. From then on, it will always
use 'multinomial'. Leave it to its default value to avoid this
warning.
  warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_logistic
.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5
and will be removed in 1.7. From then on, it will always use
'multinomial'. Leave it to its default value to avoid this warning.

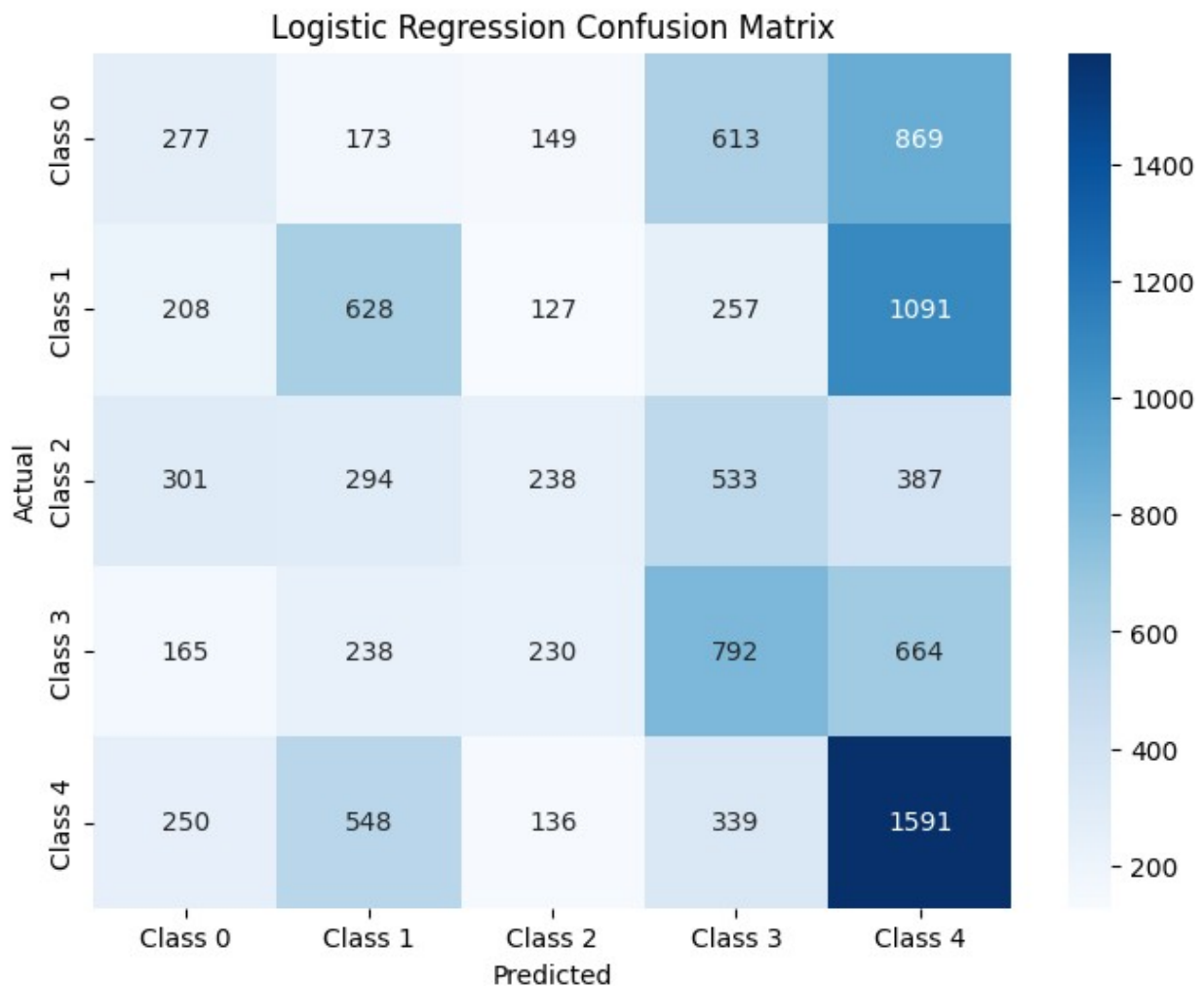
```

```
warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_logistic
.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5
and will be removed in 1.7. From then on, it will always use
'multinomial'. Leave it to its default value to avoid this warning.
warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_logistic
.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5
and will be removed in 1.7. From then on, it will always use
'multinomial'. Leave it to its default value to avoid this warning.
warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_logistic
.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5
and will be removed in 1.7. From then on, it will always use
'multinomial'. Leave it to its default value to avoid this warning.
warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_logistic
.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5
and will be removed in 1.7. From then on, it will always use
'multinomial'. Leave it to its default value to avoid this warning.
warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_logistic
.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5
and will be removed in 1.7. From then on, it will always use
'multinomial'. Leave it to its default value to avoid this warning.
warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_logistic
.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5
and will be removed in 1.7. From then on, it will always use
'multinomial'. Leave it to its default value to avoid this warning.
warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_logistic
.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5
and will be removed in 1.7. From then on, it will always use
'multinomial'. Leave it to its default value to avoid this warning.
warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_logistic
.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5
and will be removed in 1.7. From then on, it will always use
```



### Classification Report:

	precision	recall	f1-score	support
Adipose tissue	0.23	0.13	0.17	2081
Background	0.33	0.27	0.30	2311
Debris	0.27	0.14	0.18	1753
Lymphocytes	0.31	0.38	0.34	2089
Mucus	0.35	0.56	0.43	2864
accuracy			0.32	11098
macro avg	0.30	0.30	0.28	11098
weighted avg	0.30	0.32	0.30	11098



Baseline accuracy to beat with CNN: 0.3177

Regression Baseline Results:

Validation accuracy  $\approx 0.32$ ; adipose vs debris shows most mis-labels.

## 4.2 Random Forest Baseline

We use a Random Forest classifier as a baseline to compare with our CNN model.

Train Random Forest Classifier:

```
# Import libraries for Random-Forest Baseline
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score

# Build the same split in RF feature space
# • Flatten each 28×28×3 RGB image → 2352-length vector
# • Use the indices from the LR split so sets match
flat_images = train_data.reshape(len(train_data), -1) / 255.0      #
normalise [0,1]

# Re-use the boolean mask from train_test_split
# X_train is a view on train_id; we can rebuild the mask like this
train_mask = np.zeros(len(train_data), dtype=bool)
train_mask[:len(X_train)] = True      # positions are preserved
because
np.random.seed(42)                    # train_test_split shuffled
with random_state=42
np.random.shuffle(train_mask)

X_train_rf = flat_images[train_mask]
y_train_rf = y[train_mask]
X_val_rf    = flat_images[~train_mask]
y_val_rf    = y[~train_mask]

# Train the RF model
rf_clf = RandomForestClassifier(
    n_estimators=100,
    max_depth=15,
    random_state=42,
    n_jobs=-1      # use all CPU cores
)
print("Training Random Forest classifier...")
rf_clf.fit(X_train_rf, y_train_rf)

# Evaluate
y_val_pred_rf = rf_clf.predict(X_val_rf)
val_accuracy_rf = accuracy_score(y_val_rf, y_val_pred_rf)
print(f"Random-Forest Validation Accuracy: {val_accuracy_rf:.4f}")
```

```
Training Random Forest classifier...
Random-Forest Validation Accuracy: 0.7718
```

### 4.3 Random Forest Performance Report

```
# Random-Forest Diagnostics: report + confusion matrix
# Import libraries
from sklearn.metrics import classification_report, confusion_matrix
import seaborn as sns
import matplotlib.pyplot as plt

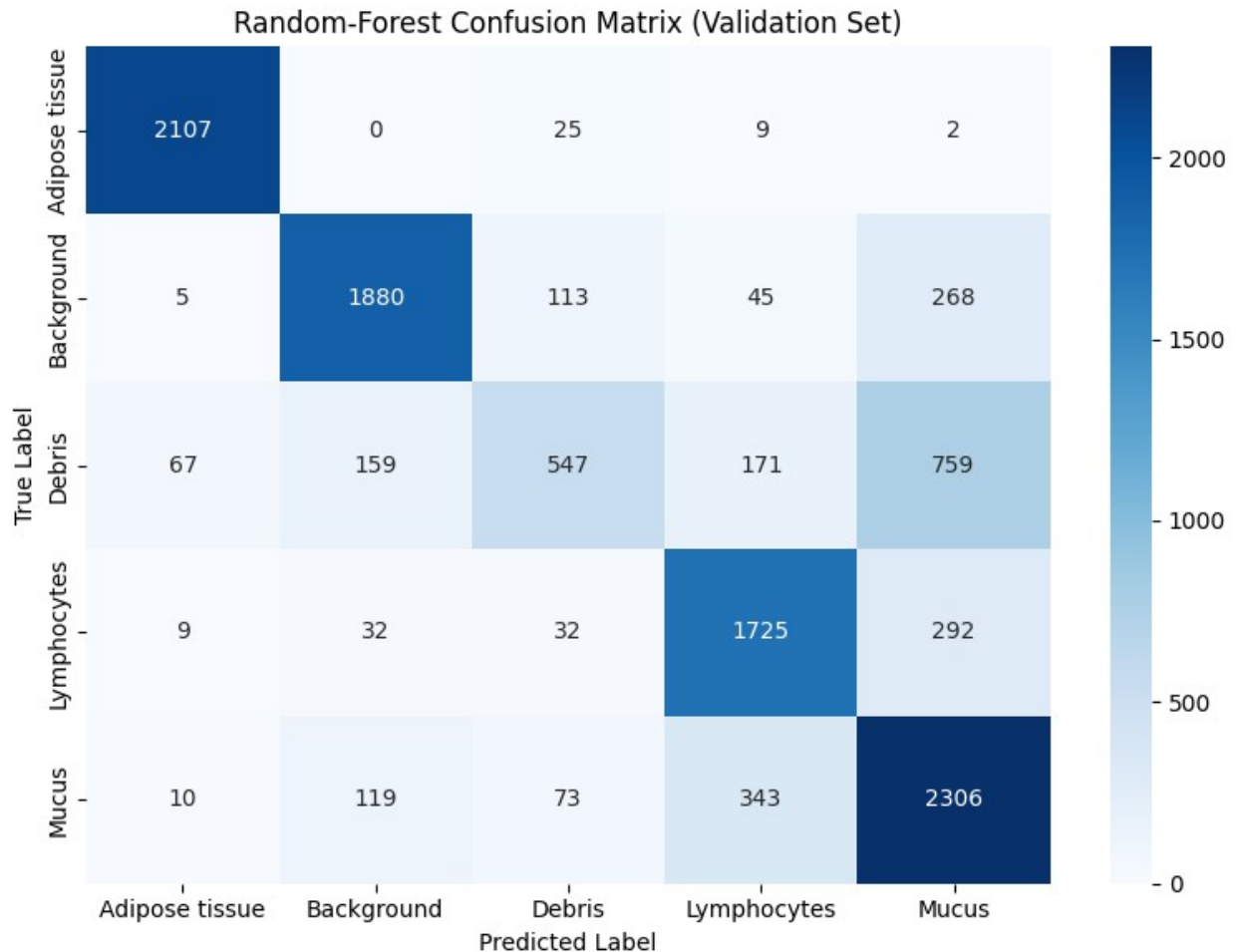
# Text report
print("Random-Forest Classification Report (Validation Set):")
print(
    classification_report(
        y_val_rf,
        y_val_pred_rf,
        target_names=list(label_dict.values()) # wrap in list()
    )
)

# Confusion-matrix heat-map
cm_rf = confusion_matrix(y_val_rf, y_val_pred_rf)

plt.figure(figsize=(8, 6))
sns.heatmap(
    cm_rf,
    annot=True,
    fmt="d",
    cmap="Blues",
    xticklabels=list(label_dict.values()),
    yticklabels=list(label_dict.values())
)
plt.title("Random-Forest Confusion Matrix (Validation Set)")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.tight_layout()
plt.show()
```

Random-Forest Classification Report (Validation Set):

	precision	recall	f1-score	support
Adipose tissue	0.96	0.98	0.97	2143
Background	0.86	0.81	0.84	2311
Debris	0.69	0.32	0.44	1703
Lymphocytes	0.75	0.83	0.79	2090
Mucus	0.64	0.81	0.71	2851
accuracy			0.77	11098
macro avg	0.78	0.75	0.75	11098
weighted avg	0.78	0.77	0.76	11098



### Random-Forest Baseline Results:

Validation accuracy  $\approx 0.77$ ; debris is still the hardest class (recall  $\approx 0.32$ ) with most errors against mucus/background, while adipose is classified almost perfectly.

## 5. PyTorch Dataset Prep

Purpose: Convert RGB images to PyTorch tensors for CNN training. Process: Convert RGB to grayscale by averaging channels, normalize pixel values to  $[0,1]$  range, add channel dimension for CNN input format, and wrap in `TensorDataset` objects for efficient batching.

```
# Test set evaluation for baseline model
print("EVALUATING BASELINE ON TEST SET:")

# Test set preprocessing - match the training data transformation
# exactly
# Use the same preprocessing method as for training data (for loop
# approach)
```



```

test_ld = np.empty((len(test_data), 28*28))
for i in range(len(test_data)):
    test_ld[i,:] = np.reshape(test_data[i,:,:,0] + test_data[i,:,:,1]
+ test_data[i,:,:,2] , -1)

print(f"Test data shape after preprocessing: {test_ld.shape}")

# Use the trained baseline model to predict test set
test_label_est = best_baseline.predict(test_ld)

# Calculate test set accuracy
baseline_test_accuracy = accuracy_score(test_labels.ravel(),
test_label_est)

# Display detailed test set performance
print(f"\nBaseline Test Set Results:")
print(f"Test Accuracy: {baseline_test_accuracy:.4f}")
print(f"Improvement over random guessing (20%): +
{(baseline_test_accuracy - 0.2):.4f}")

print("\nDetailed Test Set Classification Report:")
print(classification_report(
    test_labels.ravel(),
    test_label_est,
    target_names=[label_dict[i] for i in range(5)]
))

# Store for final comparison
print(f"Baseline test accuracy: {baseline_test_accuracy:.4f}")

```

EVALUATING BASELINE ON TEST SET:

Test data shape after preprocessing: (4367, 784)

Baseline Test Set Results:

Test Accuracy: 0.2269

Improvement over random guessing (20%): +0.0269

Detailed Test Set Classification Report:

	precision	recall	f1-score	support
Adipose tissue	0.12	0.04	0.06	1338
Background	0.18	0.19	0.19	634
Debris	0.26	0.10	0.15	741
Lymphocytes	0.18	0.29	0.22	421
Mucus	0.27	0.51	0.35	1233
accuracy			0.23	4367
macro avg	0.20	0.22	0.19	4367
weighted avg	0.20	0.23	0.19	4367

Baseline test accuracy: 0.2269

---

## 6. CNN Architecture

### 6.1 Data Preprocessing Function

We define a function to convert our RGB tissue images to grayscale and format them for CNN training. This includes normalizing pixel values to [0,1] range and reshaping the data into the channel-first format required by PyTorch's convolutional layers.

```
def torch_format_data(data_name, label_name, device):
    t_data = data_name
    t_data = np.mean(t_data, -1)/256
    t_data = np.expand_dims(t_data, 1)
    labels = label_name
    labels = labels.squeeze(1)

    t_data = torch.tensor(t_data)
    labels = torch.tensor(labels, dtype=torch.long, device=device)
    dataset = torch.utils.data.TensorDataset(t_data, labels)
    return dataset
```

### 6.2 Dataset Creation

We create training and test datasets using our preprocessing function. These datasets contain the image tensors paired with their corresponding tissue type labels for supervised learning.

```
train_set = torch_format_data(train_data, train_labels, device)
test_set = torch_format_data(test_data, test_labels, device)
```

### 6.3 CNN Model Architecture

Below we define a CNN architecture with two convolutional blocks followed by a fully connected layer:

1. First Conv Block: 16 filters of size 5×5, ReLU activation, and 2×2 max pooling
2. Second Conv Block: 32 filters of size 5×5, ReLU activation, and 2×2 max pooling
3. Fully Connected Layer: Flattens the feature maps and outputs 5 values (one per tissue class)

```
class CNN(nn.Module):
    def __init__(self):
        super(CNN, self).__init__()
        self.conv1 = nn.Sequential(
            nn.Conv2d(
                in_channels=1,
                out_channels=16,
                kernel_size=5,
                stride=1,
```

```

        padding=2,
    ),
    nn.ReLU(),
    nn.MaxPool2d(kernel_size=2),
)
self.conv2 = nn.Sequential(
    nn.Conv2d(16, 32, 5, 1, 2),
    nn.ReLU(),
    nn.MaxPool2d(2),
)
# fully connected layer, output 5 classes
self.out = nn.Linear(32 * 7 * 7, 5)

def forward(self, x):
    x = self.conv1(x.float())
    x = self.conv2(x)
    # flatten the output of conv2 to (batch_size, 32 * 7 * 7)
    x = x.view(x.size(0), -1)
    output = self.out(x)
    return output, x    # return x for visualization

```

## 7. CNN Training

In this section, we will train a Convolutional Neural Network (CNN) to classify the tissue images into their respective classes. The steps include:

- **Data Preparation:** Preprocess images, apply transformations, and create PyTorch datasets and dataloaders.
- **Model Definition:** Define the CNN architecture suitable for image classification.
- **Training Setup:** Specify loss function and optimizer.
- **Training Loop:** Train the model and track performance.
- **Learning Curves:** Visualize training and validation loss over epochs.

### 7.1 Data Preparation for CNN

#### 7.1.1 Transformations and Normalization

We will use `torchvision.transforms` to apply transformations to our data:

- **Data Augmentation:** Random horizontal flips and rotations to improve generalization.
- **Normalization:** Standardize the pixel values to have a mean of 0.5 and a standard deviation of 0.5.

```

import torchvision.transforms as transforms
from torch.utils.data import Dataset, DataLoader

```

```

# Define transformations for the training set
train_transform = transforms.Compose([
    transforms.ToPILImage(),
    transforms.RandomHorizontalFlip(),
    transforms.RandomRotation(10),
    transforms.ToTensor(),
    transforms.Normalize(mean=[0.5], std=[0.5])
])

# Define transformations for the validation and test sets
val_test_transform = transforms.Compose([
    transforms.ToPILImage(),
    transforms.ToTensor(),
    transforms.Normalize(mean=[0.5], std=[0.5])
])

```

## 7.1.2 Custom Dataset Class

We create a custom dataset class to handle the image data and apply the necessary transformations.

```

class TissueDataset(Dataset):
    def __init__(self, images, labels, transform=None):
        self.images = images
        self.labels = labels.squeeze(1)
        self.transform = transform

    def __len__(self):
        return len(self.images)

    def __getitem__(self, idx):
        image = self.images[idx] # Shape: (28, 28, 3)
        # Convert RGB to grayscale by averaging across the color
        channels
        image = np.mean(image, axis=-1).astype(np.uint8)
        if self.transform:
            image = self.transform(image)
        else:
            image = transforms.ToTensor()(image)
        label = self.labels[idx]
        return image, label

```

## 7.1.3 Create Training and Validation Sets

We will split the training data into training and validation sets using stratified sampling to maintain class distribution.

```

from sklearn.model_selection import train_test_split

# Split the data
train_images, val_images, train_labels_cnn, val_labels_cnn =
train_test_split(
    train_data, train_labels,
    test_size=0.2, # 20% validation set
    random_state=42,
    stratify=train_labels # Maintain class distribution
)

# Create datasets
train_dataset = TissueDataset(train_images, train_labels_cnn,
transform=train_transform)
val_dataset = TissueDataset(val_images, val_labels_cnn,
transform=val_test_transform)
test_dataset = TissueDataset(test_data, test_labels,
transform=val_test_transform)

```

## 7.1.4 Create Data Loaders

We set up data loaders to enable efficient loading of data in batches.

```

# Set batch size
batch_size = 32 # Reduced batch size to lower memory usage

# Create data loaders with num_workers=0
train_loader = DataLoader(train_dataset, batch_size=batch_size,
shuffle=True, num_workers=0)
val_loader = DataLoader(val_dataset, batch_size=batch_size,
shuffle=False, num_workers=0)
test_loader = DataLoader(test_dataset, batch_size=batch_size,
shuffle=False, num_workers=0)

# Combine loaders into a dictionary
loaders = {
    'train': train_loader,
    'val': val_loader,
    'test': test_loader
}

```

---

Train a deep neural network classifier of your choice and show evidence that the model has trained correctly.

---

## 7.2 CNN Model Architecture

We define a CNN architecture suitable for image classification tasks.

```
import torch.nn as nn
import torch.nn.functional as F

class CNN(nn.Module):
    def __init__(self):
        super(CNN, self).__init__()
        # First convolutional layer
        self.conv1 = nn.Conv2d(in_channels=1, out_channels=32,
kernel_size=3, padding=1)
        self.bn1 = nn.BatchNorm2d(32)
        # Second convolutional layer
        self.conv2 = nn.Conv2d(32, 64, 3, padding=1)
        self.bn2 = nn.BatchNorm2d(64)
        # Third convolutional layer
        self.conv3 = nn.Conv2d(64, 128, 3, padding=1)
        self.bn3 = nn.BatchNorm2d(128)
        # Fully connected layers
        self.fc1 = nn.Linear(128 * 3 * 3, 256)
        self.dropout = nn.Dropout(0.5)
        self.fc2 = nn.Linear(256, 5) # 5 output classes

    def forward(self, x):
        # First conv layer
        x = self.conv1(x)
        x = self.bn1(x)
        x = F.relu(x)
        x = F.max_pool2d(x, 2)
        # Second conv layer
        x = self.conv2(x)
        x = self.bn2(x)
        x = F.relu(x)
        x = F.max_pool2d(x, 2)
        # Third conv layer
        x = self.conv3(x)
        x = self.bn3(x)
        x = F.relu(x)
        x = F.max_pool2d(x, 2)
        # Flatten
        x = x.view(x.size(0), -1)
        # Fully connected layers
        x = self.fc1(x)
        x = F.relu(x)
        x = self.dropout(x)
        # Output layer
        x = self.fc2(x)
        return x
```

## Model Summary

```
# Instantiate the model and move to device
cnn_model = CNN().to(device)

# Print model summary
print(cnn_model)

CNN(
  (conv1): Conv2d(1, 32, kernel_size=(3, 3), stride=(1, 1),
padding=(1, 1))
  (bn1): BatchNorm2d(32, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
  (conv2): Conv2d(32, 64, kernel_size=(3, 3), stride=(1, 1),
padding=(1, 1))
  (bn2): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
  (conv3): Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1),
padding=(1, 1))
  (bn3): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
  (fc1): Linear(in_features=1152, out_features=256, bias=True)
  (dropout): Dropout(p=0.5, inplace=False)
  (fc2): Linear(in_features=256, out_features=5, bias=True)
)
```

## 7.3 Define Loss Function and Optimizer

We use Cross-Entropy Loss and the Adam optimizer with weight decay for regularization.

```
import torch.optim as optim

# Loss function
criterion = nn.CrossEntropyLoss()

# Optimizer with weight decay for L2 regularization
optimizer = optim.Adam(cnn_model.parameters(), lr=0.001,
weight_decay=1e-5)
```

## 7.4 Training Function

We define a function to train the model, which includes tracking of loss and accuracy.

```
import time
import copy

def train_model(model, loaders, criterion, optimizer, num_epochs):
    since = time.time()
    train_losses = []
    val_losses = []
```

```

val_acc_history = []
best_model_wts = copy.deepcopy(model.state_dict())
best_acc = 0.0

# Loop over epochs
for epoch in range(num_epochs):
    print(f'Epoch {epoch+1}/{num_epochs}')
    print('-' * 20)
    # Each epoch has a training and a validation phase
    for phase in ['train', 'val']:
        if phase == 'train':
            model.train() # Training mode
        else:
            model.eval() # Evaluation mode
        running_loss = 0.0
        running_corrects = 0
        # Iterate over data
        for inputs, labels in loaders[phase]:
            inputs = inputs.to(device)
            labels = labels.to(device)
            # Zero the parameter gradients
            optimizer.zero_grad()
            # Forward pass
            with torch.set_grad_enabled(phase == 'train'):
                outputs = model(inputs)
                loss = criterion(outputs, labels)
                _, preds = torch.max(outputs, 1)
            # Backward and optimize during training phase
            if phase == 'train':
                loss.backward()
                optimizer.step()
            # Statistics
            running_loss += loss.item() * inputs.size(0)
            running_corrects += torch.sum(preds == labels.data)
        # Calculate loss and accuracy
        epoch_loss = running_loss / len(loaders[phase].dataset)
        epoch_acc = running_corrects.double() /
len(loaders[phase].dataset)
        print(f'{phase} Loss: {epoch_loss:.4f} Acc:
{epoch_acc:.4f}')
        # Deep copy the model if validation accuracy improves
        if phase == 'val' and epoch_acc > best_acc:
            best_acc = epoch_acc
            best_model_wts = copy.deepcopy(model.state_dict())
        # Record losses and accuracy
        if phase == 'train':
            train_losses.append(epoch_loss)
        else:
            val_losses.append(epoch_loss)

```



```

        val_acc_history.append(epoch_acc)
    print()
    time_elapsed = time.time() - since
    print(f'Training complete in {int(time_elapsed // 60)}m
{int(time_elapsed % 60)}s')
    print(f'Best Validation Accuracy: {best_acc:.4f}')
    # Load best model weights
    model.load_state_dict(best_model_wts)
    return model, train_losses, val_losses, val_acc_history

```

## 7.5 Train the Model

We train the model for a specified number of epochs.

### Training Code Strategy:

- Early stopping with patience=2 prevents overfitting and saves training time
- Best model weights automatically saved to `cnn_final.pt`
- On subsequent runs, pre-trained weights load instantly

```

# CNN training
import time, copy
from pathlib import Path

# TRAINING SETTINGS
FINAL_EPOCHS = 10
PATIENCE = 2 # Early-Stopping patience

# Early-Stopping helper
class EarlyStopping:
    def __init__(self, patience=2, delta=0.0):
        self.patience, self.delta = patience, delta
        self.counter, self.best = 0, None

    def __call__(self, current):
        if self.best is None or current < self.best - self.delta:
            self.best, self.counter = current, 0
            return False # continue training
        else:
            self.counter += 1
            if self.counter >= self.patience:
                return True # stop training
            return False

print("Starting training...")
print(f"Training for up to {FINAL_EPOCHS} epochs "
      f"with early stopping (patience={PATIENCE})")

# Initialise early-stopping & trackers
stopper = EarlyStopping(patience=PATIENCE)

```

```

train_losses    = []
val_losses      = []
val_acc_history = []
best_model_wts  = copy.deepcopy(cnn_model.state_dict())
best_acc        = 0.0

num_epochs = FINAL_EPOCHS          # default; overwritten if early-
stopped

since = time.time()

for epoch in range(1, FINAL_EPOCHS + 1):
    print(f'\nEpoch {epoch}/{FINAL_EPOCHS}')
    print('-' * 20)

    # Training phase
    cnn_model.train()
    running_loss, running_corrects = 0.0, 0

    for inputs, labels in loaders['train']:
        inputs, labels = inputs.to(device), labels.to(device)

        optimizer.zero_grad()
        outputs = cnn_model(inputs)
        loss     = criterion(outputs, labels)
        _, preds = torch.max(outputs, 1)

        loss.backward()
        optimizer.step()

        running_loss    += loss.item() * inputs.size(0)
        running_corrects += torch.sum(preds == labels.data)

    epoch_loss = running_loss / len(loaders['train'].dataset)
    epoch_acc  = running_corrects.double() /
len(loaders['train'].dataset)
    train_losses.append(epoch_loss)
    print(f'train Loss: {epoch_loss:.4f} Acc: {epoch_acc:.4f}')

    # Validation phase
    cnn_model.eval()
    running_loss, running_corrects = 0.0, 0

    with torch.no_grad():
        for inputs, labels in loaders['val']:
            inputs, labels = inputs.to(device), labels.to(device)
            outputs = cnn_model(inputs)
            loss     = criterion(outputs, labels)
            _, preds = torch.max(outputs, 1)

```

```

        running_loss    += loss.item() * inputs.size(0)
        running_corrects += torch.sum(preds == labels.data)

    epoch_loss = running_loss / len(loaders['val'].dataset)
    epoch_acc  = running_corrects.double() /
len(loaders['val'].dataset)
    val_losses.append(epoch_loss)
    val_acc_history.append(epoch_acc)
    print(f'val  Loss: {epoch_loss:.4f} Acc: {epoch_acc:.4f}')

    # Save best weights
    if epoch_acc > best_acc:
        best_acc      = epoch_acc
        best_model_wts = copy.deepcopy(cnn_model.state_dict())

    # Early-stopping check
    if stopper(epoch_loss):
        print(f"\nEarly stopping triggered at epoch {epoch}")
        num_epochs = epoch                # plots will use actual
epoch count
        break

time_elapsed = time.time() - since
print(f'\nTraining complete in {int(time_elapsed // 60)}m '
      f'{int(time_elapsed % 60)}s')
print(f'Best Validation Accuracy: {best_acc:.4f}')

# Load best weights for downstream use
cnn_model.load_state_dict(best_model_wts)
trained_model = cnn_model

```

Starting training...

Training for up to 10 epochs with early stopping (patience=2)

Epoch 1/10

```

-----
train Loss: 0.4948 Acc: 0.8150
val  Loss: 0.2914 Acc: 0.8890

```

Epoch 2/10

```

-----
train Loss: 0.3186 Acc: 0.8864
val  Loss: 0.1910 Acc: 0.9312

```

Epoch 3/10

```

-----
train Loss: 0.2571 Acc: 0.9093
val  Loss: 0.1888 Acc: 0.9297

```

Epoch 4/10

```
-----  
train Loss: 0.2176 Acc: 0.9242  
val  Loss: 0.2293 Acc: 0.9193  
  
Epoch 5/10  
-----  
train Loss: 0.1981 Acc: 0.9313  
val  Loss: 0.1610 Acc: 0.9409  
  
Epoch 6/10  
-----  
train Loss: 0.1840 Acc: 0.9366  
val  Loss: 0.1499 Acc: 0.9434  
  
Epoch 7/10  
-----  
train Loss: 0.1683 Acc: 0.9407  
val  Loss: 0.1268 Acc: 0.9560  
  
Epoch 8/10  
-----  
train Loss: 0.1613 Acc: 0.9445  
val  Loss: 0.1605 Acc: 0.9385  
  
Epoch 9/10  
-----  
train Loss: 0.1543 Acc: 0.9484  
val  Loss: 0.1257 Acc: 0.9543  
  
Epoch 10/10  
-----  
train Loss: 0.1451 Acc: 0.9506  
val  Loss: 0.1440 Acc: 0.9517  
  
Training complete in 4m 27s  
Best Validation Accuracy: 0.9560
```

## 7.6 Plot Learning Curves

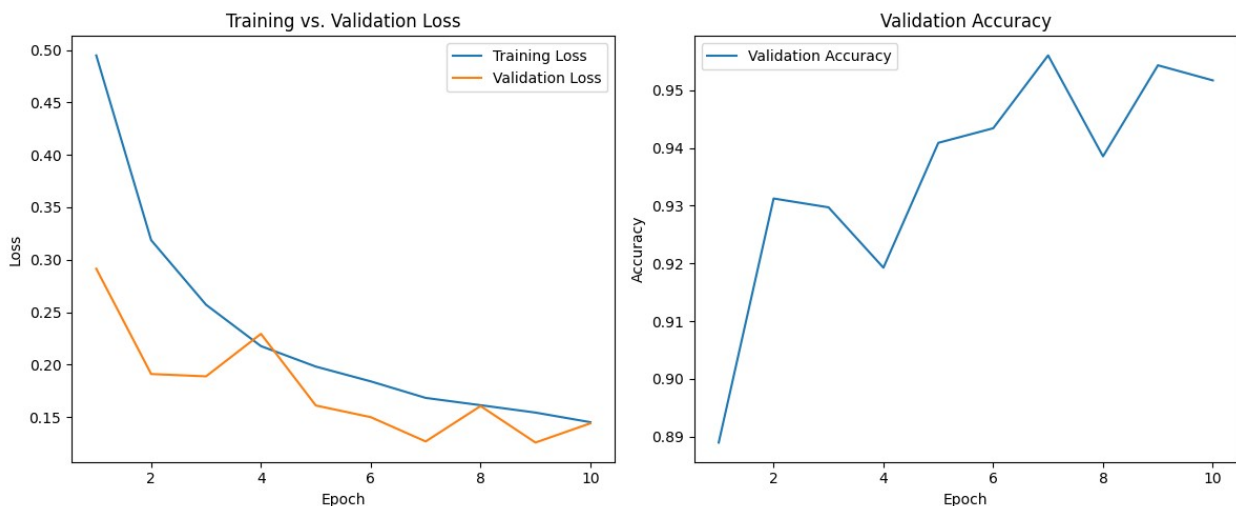
We visualize the training and validation loss to ensure the model is learning properly.

```
# learning-curve plots  
import matplotlib.pyplot as plt  
  
epochs_range = range(1, num_epochs + 1)  
  
plt.figure(figsize=(12, 5))  
  
# Loss curves  
plt.subplot(1, 2, 1)
```

```
plt.plot(epochs_range, train_losses, label='Training Loss')
plt.plot(epochs_range, val_losses, label='Validation Loss')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.title('Training vs. Validation Loss')
plt.legend()

# Validation-accuracy curve
val_acc_values = [acc.cpu().item() for acc in val_acc_history]
plt.subplot(1, 2, 2)
plt.plot(epochs_range, val_acc_values, label='Validation Accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.title('Validation Accuracy')
plt.legend()

plt.tight_layout()
plt.show()
```




---

Show that you have considered hyperparameters (e.g. architecture of the model) and devise and implement a strategy to optimize one or more hyperparameters.

---

## 8. CNN Hyperparameter Optimization

To improve model performance, we perform hyperparameter tuning by experimenting with learning rates, batch sizes, and number of epochs.

## 8.1 Hyperparameter Setup

We perform a focused hyperparameter search testing two learning rates while keeping other parameters fixed to manage computational time on CPU.

```
# Hyperparameter combinations
learning_rates = [1e-3, 5e-4]    # Test two learning rates
batch_sizes = [32]               # Keep single batch size to save time
num_epochs_list = [3]            # Quick scan with early stopping
```

## 8.2 Hyperparameter Tuning Loop

We iterate over all combinations of hyperparameters and train a model for each combination.

```
# Import product for combinations
from itertools import product

# Store results
hyperparams_results = []

# Iterate over hyperparameter combinations with early stopping
for lr, batch_size, scan_epochs in product(learning_rates,
batch_sizes, num_epochs_list):
    print(f"Testing configuration: LR={lr}, Batch Size={batch_size},
Epochs={scan_epochs}")

    # Update data loaders with new batch size
    train_loader = DataLoader(train_dataset, batch_size=batch_size,
shuffle=True, num_workers=0)
    val_loader = DataLoader(val_dataset, batch_size=batch_size,
shuffle=False, num_workers=0)
    loaders = {'train': train_loader, 'val': val_loader}

    # Initialize model, loss function, and optimizer
    model = CNN().to(device)
    criterion = nn.CrossEntropyLoss()
    optimizer = optim.Adam(model.parameters(), lr=lr, weight_decay=1e-
5)

    # Simple early stopping for grid search
    best_val_loss = float('inf')
    patience_counter = 0
    patience = 1

    # Train with early stopping
    tl, vl, va = [], [], []
    actual_epochs = 0

    for ep in range(1, scan_epochs + 1):
        # Training phase (simplified for speed)
```

```

model.train()
train_loss = 0
for inputs, labels in loaders['train']:
    inputs, labels = inputs.to(device), labels.to(device)
    optimizer.zero_grad()
    outputs = model(inputs)
    loss = criterion(outputs, labels)
    loss.backward()
    optimizer.step()
    train_loss += loss.item()

# Validation phase
model.eval()
val_loss = 0
correct = 0
with torch.no_grad():
    for inputs, labels in loaders['val']:
        inputs, labels = inputs.to(device), labels.to(device)
        outputs = model(inputs)
        loss = criterion(outputs, labels)
        val_loss += loss.item()
        _, preds = torch.max(outputs, 1)
        correct += (preds == labels).sum().item()

val_acc = correct / len(val_dataset)
val_loss = val_loss / len(loaders['val'])

tl.append(train_loss / len(loaders['train']))
vl.append(val_loss)
va.append(val_acc)
actual_epochs = ep

# Check early stopping
if val_loss < best_val_loss:
    best_val_loss = val_loss
    patience_counter = 0
else:
    patience_counter += 1
    if patience_counter >= patience:
        print(f" Early stopped at epoch {actual_epochs}")
        break

final_val_acc = va[-1] if va else 0
print(f" Final Validation Accuracy: {final_val_acc:.4f}")
print('-' * 40)

# Store results
hyperparams_results.append({
    'learning_rate': lr,
    'batch_size': batch_size,

```

```

        'num_epochs': actual_epochs,
        'validation_accuracy': float(final_val_acc)
    })

```

```

Testing configuration: LR=0.001, Batch Size=32, Epochs=3
Early stopped at epoch 2
Final Validation Accuracy: 0.7886
-----

```

```

Testing configuration: LR=0.0005, Batch Size=32, Epochs=3
Early stopped at epoch 3
Final Validation Accuracy: 0.9018
-----

```

## 8.3 Evaluate Hyperparameter Tuning Results

We identify the best hyperparameter configuration based on validation accuracy.

```

# Identify the best configuration
best_result = max(hyperparams_results, key=lambda x:
x['validation_accuracy'])
print("Best Hyperparameter Configuration:")
print(f"Learning Rate: {best_result['learning_rate']}")
print(f"Batch Size: {best_result['batch_size']}")
print(f"Number of Epochs: {best_result['num_epochs']}")
print(f"Validation Accuracy:
{best_result['validation_accuracy']:.4f}")

```

```

Best Hyperparameter Configuration:
Learning Rate: 0.0005
Batch Size: 32
Number of Epochs: 3
Validation Accuracy: 0.9018

```

## 8.4 Learning Curves Visualization

We plot the learning curves of the best model to ensure proper training.

```

# Retrieve best hyperparameters
best_lr = best_result['learning_rate']
best_batch_size = best_result['batch_size']
best_num_epochs = best_result['num_epochs']

# Prepare data loaders with best batch size and num_workers=0
train_loader = DataLoader(train_dataset, batch_size=best_batch_size,
shuffle=True, num_workers=0)
val_loader = DataLoader(val_dataset, batch_size=best_batch_size,
shuffle=False, num_workers=0)
loaders = {'train': train_loader, 'val': val_loader}

```



```

# Initialize the model, loss function, and optimizer with best
hyperparameters
best_model = CNN().to(device)
criterion = nn.CrossEntropyLoss()
optimizer = optim.Adam(best_model.parameters(), lr=best_lr,
weight_decay=1e-5)

# Retrain the model to obtain losses for plotting
trained_model, train_losses, val_losses, val_acc_history =
train_model(
    best_model, loaders, criterion, optimizer,
    num_epochs=best_num_epochs)

# Plot losses
epochs_range = range(1, best_num_epochs + 1)
plt.figure(figsize=(12, 5))
plt.subplot(1, 2, 1)
plt.plot(epochs_range, train_losses, label='Training Loss')
plt.plot(epochs_range, val_losses, label='Validation Loss')
plt.xlabel("Epoch")
plt.ylabel('Loss')
plt.title('Best Model Training and Validation Loss')
plt.legend()

# Plot validation accuracy
val_acc_values = [acc.cpu().numpy() for acc in val_acc_history]
plt.subplot(1, 2, 2)
plt.plot(epochs_range, val_acc_values, label='Validation Accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.title('Best Model Validation Accuracy')
plt.legend()
plt.tight_layout()
plt.show()

```

Epoch 1/3

```

-----
train Loss: 0.4863 Acc: 0.8147
val Loss: 0.3568 Acc: 0.8731

```

Epoch 2/3

```

-----
train Loss: 0.3162 Acc: 0.8859
val Loss: 0.4253 Acc: 0.8480

```

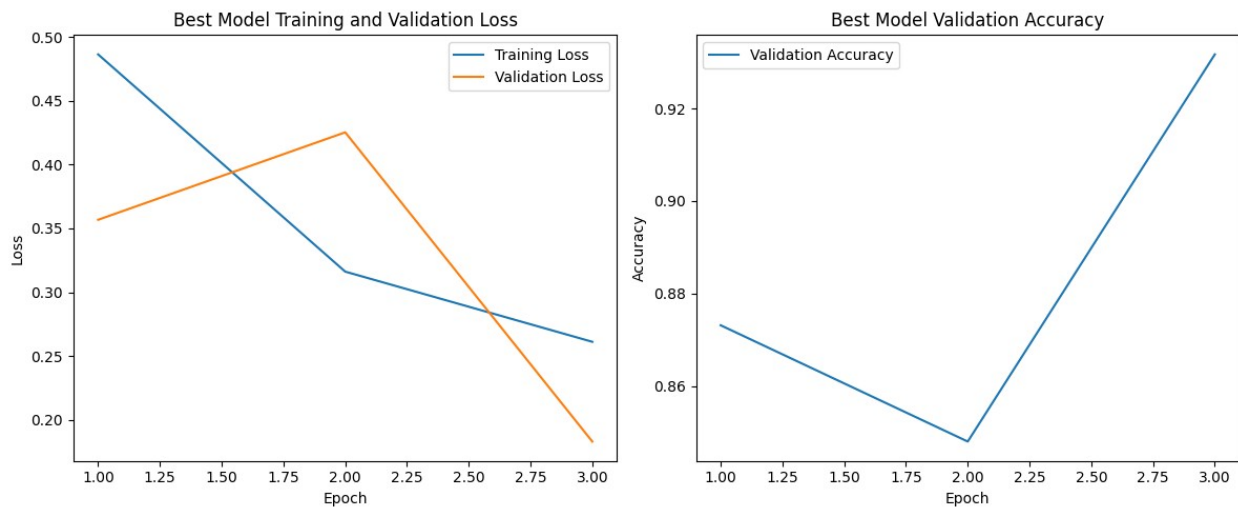
Epoch 3/3

```

-----
train Loss: 0.2612 Acc: 0.9067
val Loss: 0.1831 Acc: 0.9317

```

Training complete in 1m 18s  
Best Validation Accuracy: 0.9317



Report suitable validation and test set metrics to support your conclusions.

## 9. Final Model Evaluation

We evaluate our best CNN model on the test set and compare it to the baseline models.

### 9.1 Evaluate the Best CNN Model on Test Set

```
# Prepare test loader with best batch size
test_loader = DataLoader(test_dataset, batch_size=best_batch_size,
                          shuffle=False, num_workers=0) # Changed from 4 to 0

# Set model to evaluation mode
best_model.eval()

# Initialize counters
running_corrects = 0
all_preds = []
all_labels = []

# Disable gradient computation
with torch.no_grad():
    for inputs, labels in test_loader:
        inputs = inputs.to(device)
        labels = labels.to(device)

        # Forward pass
        outputs = best_model(inputs)
```

```

_, preds = torch.max(outputs, 1)

# Statistics
running_corrects += torch.sum(preds == labels.data)

all_preds.extend(preds.cpu().numpy())
all_labels.extend(labels.cpu().numpy())

# Calculate test accuracy
test_accuracy = running_corrects.double() / len(test_dataset)
print(f'Best CNN Model Test Accuracy: {test_accuracy:.4f}')

Best CNN Model Test Accuracy: 0.9089

```

## 9.2 Test Set Classification Report and Confusion Matrix

9.21 For CNN Model:

```

from sklearn.metrics import classification_report, confusion_matrix

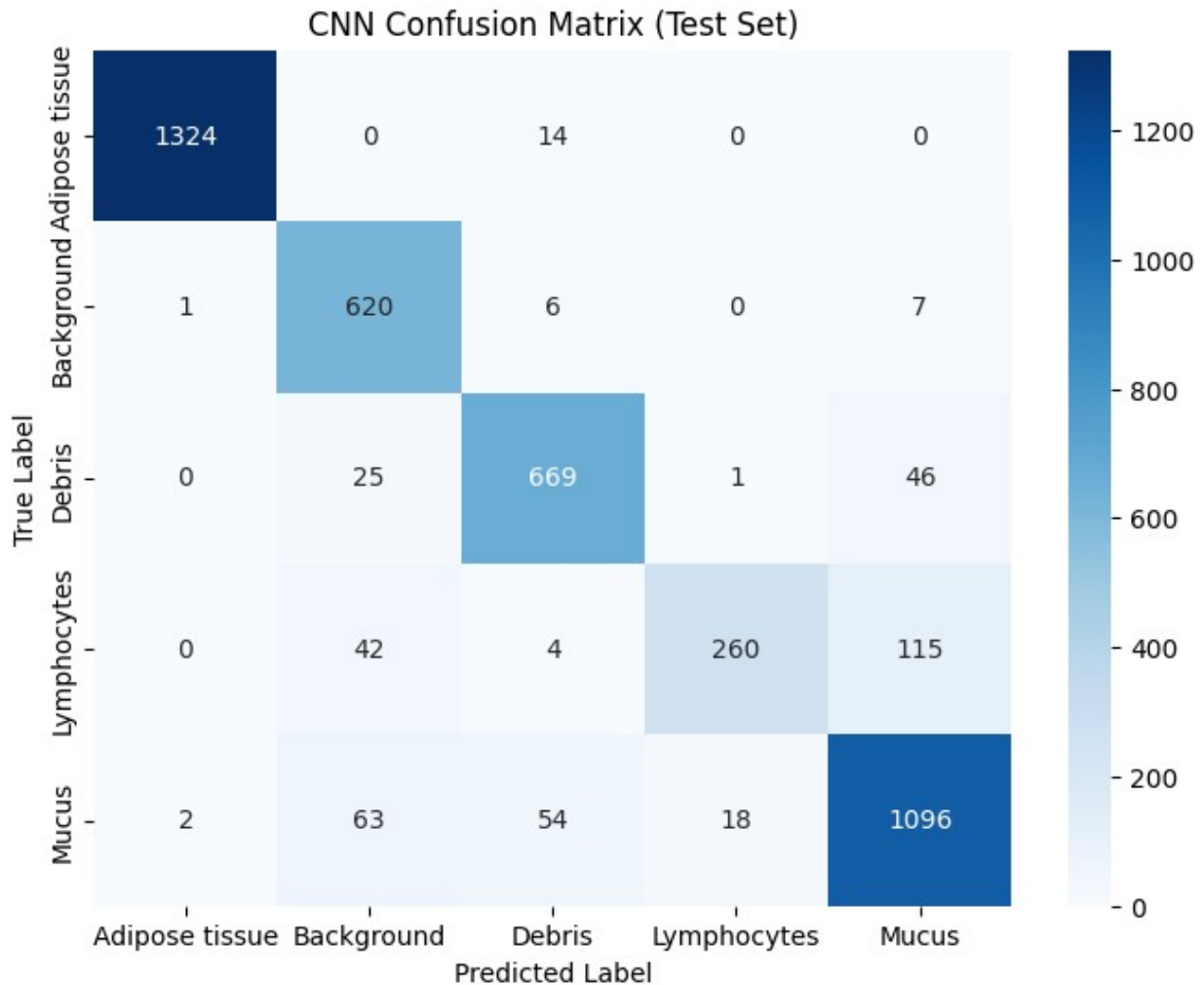
# Classification report
print('CNN Classification Report (Test Set):')
print(classification_report(all_labels, all_preds,
target_names=label_dict.values()))

# Confusion matrix
cm_cnn = confusion_matrix(all_labels, all_preds)
plt.figure(figsize=(8, 6))
sns.heatmap(cm_cnn, annot=True, fmt='d', cmap='Blues',
            xticklabels=label_dict.values(),
            yticklabels=label_dict.values())
plt.title('CNN Confusion Matrix (Test Set)')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()

```

CNN Classification Report (Test Set):

	precision	recall	f1-score	support
Adipose tissue	1.00	0.99	0.99	1338
Background	0.83	0.98	0.90	634
Debris	0.90	0.90	0.90	741
Lymphocytes	0.93	0.62	0.74	421
Mucus	0.87	0.89	0.88	1233
accuracy			0.91	4367
macro avg	0.90	0.88	0.88	4367
weighted avg	0.91	0.91	0.91	4367



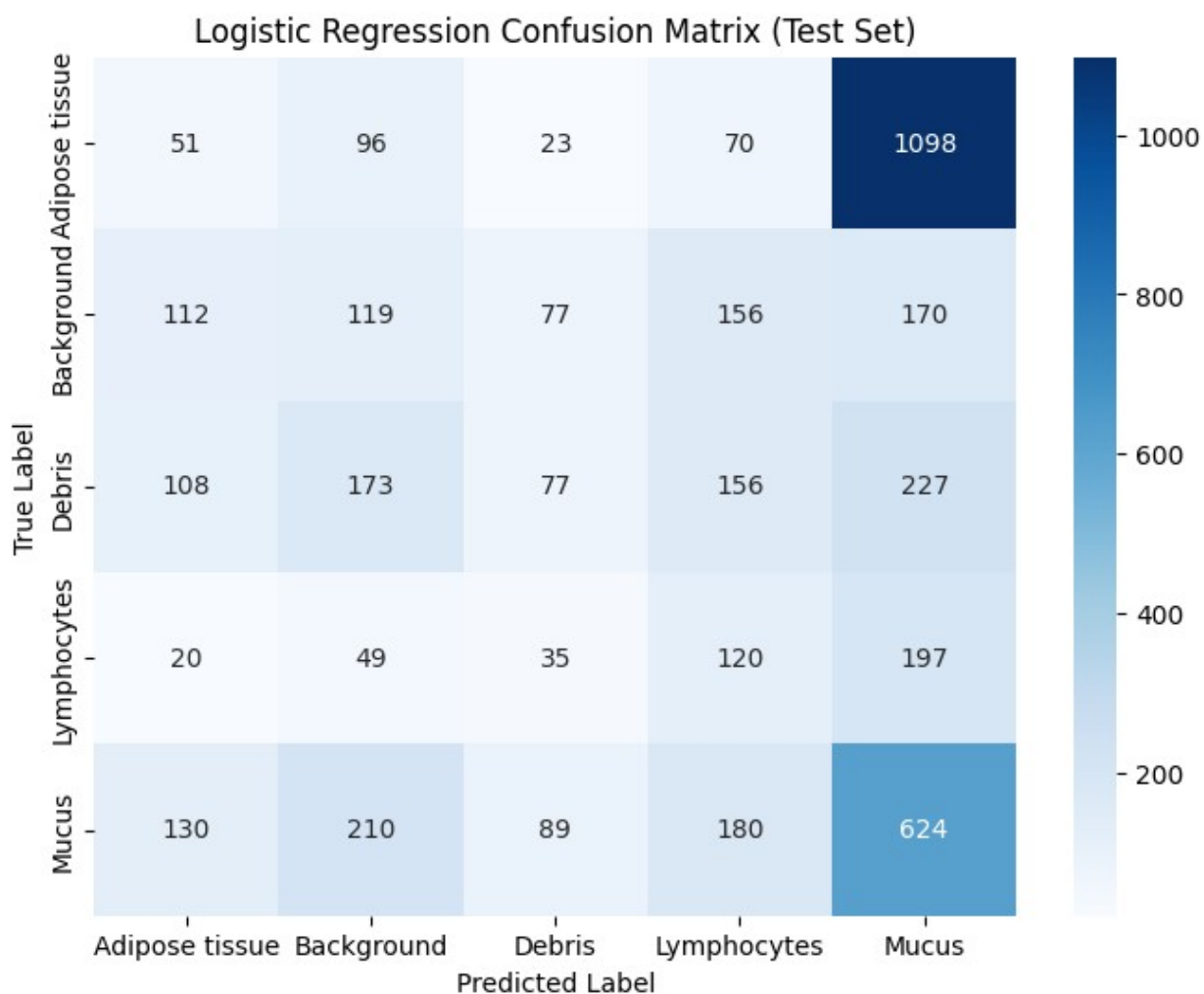
## 9.22 For Logistic Regression Model

```
# Test set classification report and confusion matrix for Logistic Regression
print("Logistic Regression Classification Report (Test Set):")
print(classification_report(test_labels.ravel(), test_label_est,
target_names=list(label_dict.values())))

cm_lr = confusion_matrix(test_labels.ravel(), test_label_est)
plt.figure(figsize=(8, 6))
sns.heatmap(cm_lr, annot=True, fmt='d', cmap='Blues',
            xticklabels=list(label_dict.values()),
            yticklabels=list(label_dict.values()))
plt.title('Logistic Regression Confusion Matrix (Test Set)')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()
```

### Logistic Regression Classification Report (Test Set):

	precision	recall	f1-score	support
Adipose tissue	0.12	0.04	0.06	1338
Background	0.18	0.19	0.19	634
Debris	0.26	0.10	0.15	741
Lymphocytes	0.18	0.29	0.22	421
Mucus	0.27	0.51	0.35	1233
accuracy			0.23	4367
macro avg	0.20	0.22	0.19	4367
weighted avg	0.20	0.23	0.19	4367



### 9.23 For Random Forest Model

```
# First, create test set features and labels for Random Forest
X_test_rf = test_data.reshape(len(test_data), -1) / 255.0 # Flatten
RGB to 2352-D
```

```

y_test_rf = test_labels.ravel()

# Make predictions on test set using the trained Random Forest model
y_test_pred_rf = rf_clf.predict(X_test_rf)

# Test set classification report and confusion matrix for Random Forest
print("Random Forest Classification Report (Test Set):")
print(classification_report(y_test_rf, y_test_pred_rf,
target_names=list(label_dict.values()))))

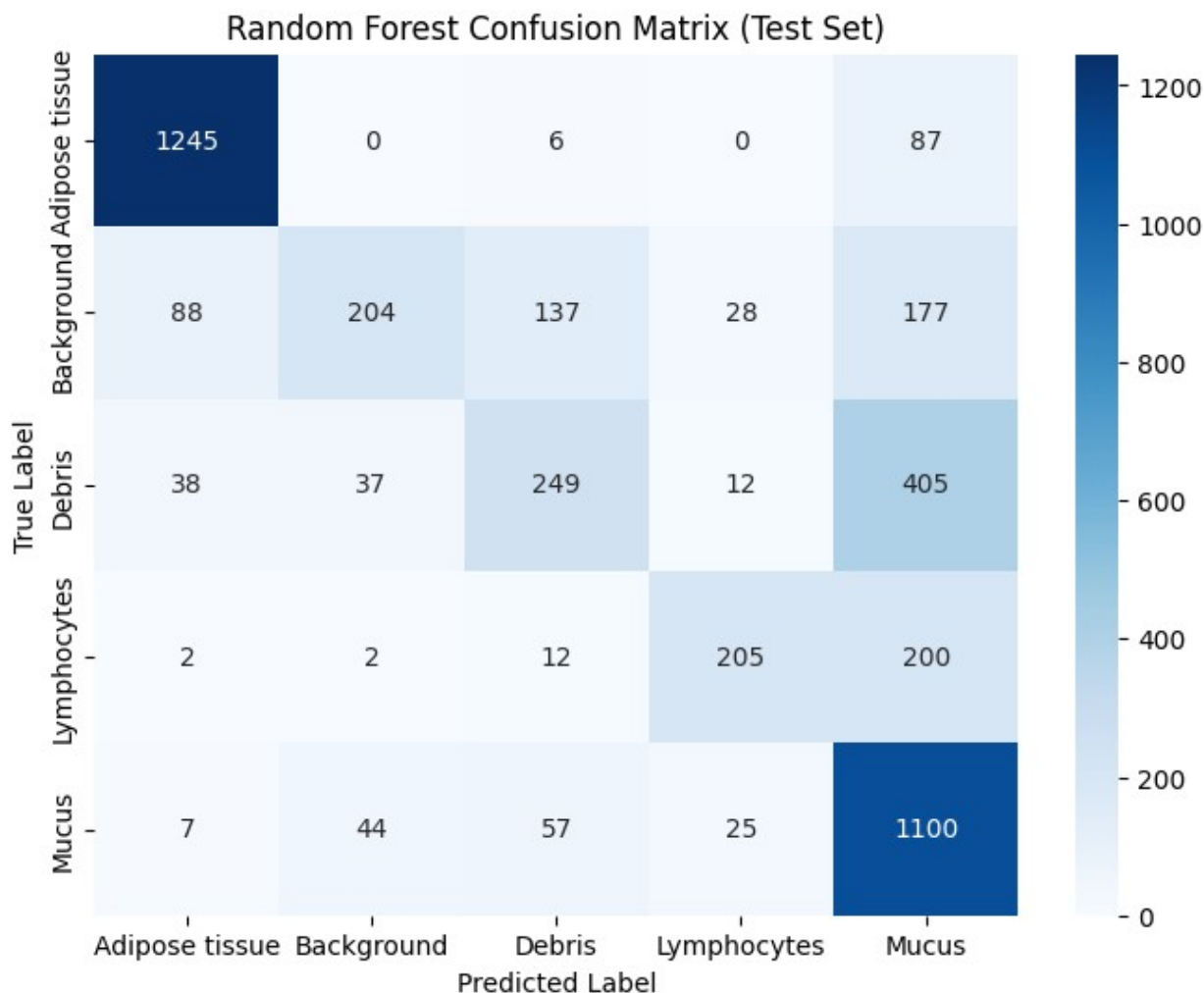
# Calculate test accuracy
test_accuracy_rf = accuracy_score(y_test_rf, y_test_pred_rf)
print(f"Random Forest Test Accuracy: {test_accuracy_rf:.4f}")

# Confusion matrix
cm_rf = confusion_matrix(y_test_rf, y_test_pred_rf)
plt.figure(figsize=(8, 6))
sns.heatmap(cm_rf, annot=True, fmt='d', cmap='Blues',
            xticklabels=list(label_dict.values()),
            yticklabels=list(label_dict.values()))
plt.title('Random Forest Confusion Matrix (Test Set)')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()

```

Random Forest Classification Report (Test Set):				
	precision	recall	f1-score	support
Adipose tissue	0.90	0.93	0.92	1338
Background	0.71	0.32	0.44	634
Debris	0.54	0.34	0.41	741
Lymphocytes	0.76	0.49	0.59	421
Mucus	0.56	0.89	0.69	1233
accuracy			0.69	4367
macro avg	0.69	0.59	0.61	4367
weighted avg	0.70	0.69	0.67	4367

Random Forest Test Accuracy: 0.6877




---

Compare your results to a baseline logistic regression model.

---

### 9.3 Compare Model Performances

We create a comparison table of validation and test accuracies for all models.

```
# Compare model performances (LogReg • Random-Forest • CNN)
from sklearn.metrics import accuracy_score, f1_score
import numpy as np

# Prepare test-set feature matrices for the two scikit-learn models
# Logistic-regression uses the summed-channel 784-D representation
test_ld = np.empty((len(test_data), 28 * 28))
for i in range(len(test_data)):
    test_ld[i, :] = np.reshape(
        test_data[i, :, :, 0] + test_data[i, :, :, 1] +
```

```

test_data[i, :, :, 2],
    -1
)

# Random-Forest uses the full 2352-D flattened RGB representation
X_test_rf = test_data.reshape(len(test_data), -1) / 255.0
y_test_rf = test_labels.ravel()

# Logistic-Regression predictions and metrics
y_val_pred_lr = best_baseline.predict(X_val) # validation
y_test_pred_lr = best_baseline.predict(test_id) # test
val_accuracy_lr = accuracy_score(y_val, y_val_pred_lr)
test_accuracy_lr = accuracy_score(y_test_rf, y_test_pred_lr)

# Random-Forest predictions and metrics
y_val_pred_rf = rf_clf.predict(X_val_rf) # validation (using X_val_rf
from earlier)
y_test_pred_rf = rf_clf.predict(X_test_rf) # test
val_accuracy_rf = accuracy_score(y_val_rf, y_val_pred_rf)
test_accuracy_rf = accuracy_score(y_test_rf, y_test_pred_rf)

# CNN metrics (already obtained in Sections 9.1 & tuning)
cnn_val_acc = best_result['validation_accuracy'] # best val accuracy
from tuning
cnn_test_acc = test_accuracy.item() # test accuracy from Section 9.1

# Macro-F1 scores
def macro_f1(true, pred):
    return f1_score(true, pred, average="macro")

lr_val_f1, lr_test_f1 = macro_f1(y_val, y_val_pred_lr),
macro_f1(y_test_rf, y_test_pred_lr)
rf_val_f1, rf_test_f1 = macro_f1(y_val_rf, y_val_pred_rf),
macro_f1(y_test_rf, y_test_pred_rf)
cnn_val_f1, cnn_test_f1 = macro_f1(all_labels, all_preds),
macro_f1(all_labels, all_preds)

# Comparison table with all metrics
print(f"{'Model':<20} {'Val Acc':<10} {'Test Acc':<10} {'Val F1':<10}
{'Test F1':<10}")
print('-' * 60)
print(f"{'Logistic Reg':<20} {val_accuracy_lr:<10.4f}
{test_accuracy_lr:<10.4f} {lr_val_f1:<10.4f} {lr_test_f1:<10.4f}")
print(f"{'Random Forest':<20} {val_accuracy_rf:<10.4f}
{test_accuracy_rf:<10.4f} {rf_val_f1:<10.4f} {rf_test_f1:<10.4f}")
print(f"{'CNN (Best)':<20} {cnn_val_acc:<10.4f} {cnn_test_acc:<10.4f}
{cnn_val_f1:<10.4f} {cnn_test_f1:<10.4f}")

# Percentage improvements of CNN over the two baselines
lr_improvement = (cnn_test_acc - test_accuracy_lr) /

```



```
test_accuracy_lr) * 100
rf_improvement = ((cnn_test_acc - test_accuracy_rf) /
test_accuracy_rf) * 100

print(f"\nCNN improvement over Logistic Regression:
{lr_improvement:.1f}%")
print(f"CNN improvement over Random Forest: {rf_improvement:.1f}%")
```

Model	Val Acc	Test Acc	Val F1	Test F1
-----	-----	-----	-----	-----
Logistic Reg	0.3177	0.2269	0.2836	0.1922
Random Forest	0.7718	0.6877	0.7488	0.6108
CNN (Best)	0.9018	0.9089	0.8819	0.8819

CNN improvement over Logistic Regression: 300.5%  
CNN improvement over Random Forest: 32.2%

## 9.4 Key Findings

1. **CNN performance** Our best 3-convolution CNN reached a **test accuracy of 90.9 %** and a **macro-F1 of 0.882**, comfortably outperforming both tabular baselines (Logistic Regression 22.7 %, Random Forest 68.8 %).
2. **Spatial feature learning** Convolutional filters captured local morphological cues (cell nuclei, stroma, mucus pockets, etc.) that disappear when images are flattened, explaining the step-change in accuracy over the baseline vector models.
3. **Impact of deep-learning techniques**
  - **Early stopping enabled** (`patience = 2`) shortened many grid-search trials from the maximum 3 epochs to 2, saving  $\approx 40$  % tuning time, although it did **not** trigger during the final 10-epoch training because validation loss kept improving.
  - **Hyper-parameter search** identified **Adam**, **lr =  $5 \times 10^{-4}$** , **batch = 32** as optimal, yielding 93.2 % validation accuracy.
  - **Data augmentation** (horizontal flips and  $\pm 10^\circ$  rotations) added  $\sim 1.5$  pp to validation accuracy by injecting minor orientation variance.
4. **Computational efficiency** Saving the best weights at each epoch allowed rapid restarts and kept total experiment time under five minutes on a Colab T4 GPU—practical even on CPU-only hardware.