

Environmental Sound

Classification

Name	ID
Ranime Ahmed Elsayed Shehata	21010531
Habiba Tarek Ramadan Ali	21010445
Karene Antoine Nassif Guirguis	21010969

Problem Statement:

Environmental sounds are critical cues for scene understanding in real-world applications like surveillance, autonomous navigation, and assistive technology. In this assignment, you will implement two deep learning approaches for classifying environmental audio recordings using the UrbanSound8K dataset. The first model will be based on Recurrent Neural Networks (RNNs), and the second will be a Transformer-based model, which you will implement from scratch (i.e., no use of nn.Transformer or Hugging Face models).

This lab aims to classify sounds from the UrbanSound8K dataset using two deep learning models:

- **Bi-directional LSTM (BiLSTM)**
 - **Transformer (implemented from scratch)**
-

Step 1: Dataset and Preprocessing:

- **Dataset:** UrbanSound8K
- **Samples:** 8732 labeled audio clips (≤ 4 seconds)
- **Classes:** 10 sound categories (e.g., dog_bark, siren, jackhammer, etc.)
- **Splits:**
 - Folds 1–6: Training
 - Folds 7–8: Validation
 - Folds 9–10: Testing

Preprocessing steps:

- Loaded audio using **librosa**
- Extracted and visualized waveforms and mel spectrograms
- Played one sample per class for auditory inspection

```
1 import librosa
2 import librosa.display
3 import matplotlib.pyplot as plt
4 import pandas as pd
5 import numpy as np
6 import os
7 import IPython.display as ipd
8
9 metadata = pd.read_csv('./UrbanSound8K/metadata/UrbanSound8K.csv')
10
11 def load_audio(file_path):
12     signal, sr = librosa.load(file_path, sr=None)
13     return signal, sr
14
15 def plot_waveform(signal, sr):
16     plt.figure(figsize=(10, 3))
17     librosa.display.waveshow(signal, sr=sr)
18     plt.title("Waveform")
19     plt.show()
20
21 def plot_melspectrogram(signal, sr):
22     mel_spec = librosa.feature.melspectrogram(y=signal, sr=sr)
23     log_mel_spec = librosa.power_to_db(mel_spec, ref=np.max)
24     plt.figure(figsize=(10, 3))
25     librosa.display.specshow(log_mel_spec, sr=sr, x_axis='time', y_axis='mel')
26     plt.title("Mel-Spectrogram")
27     plt.colorbar()
28     plt.show()
29
30
```

Metadata:

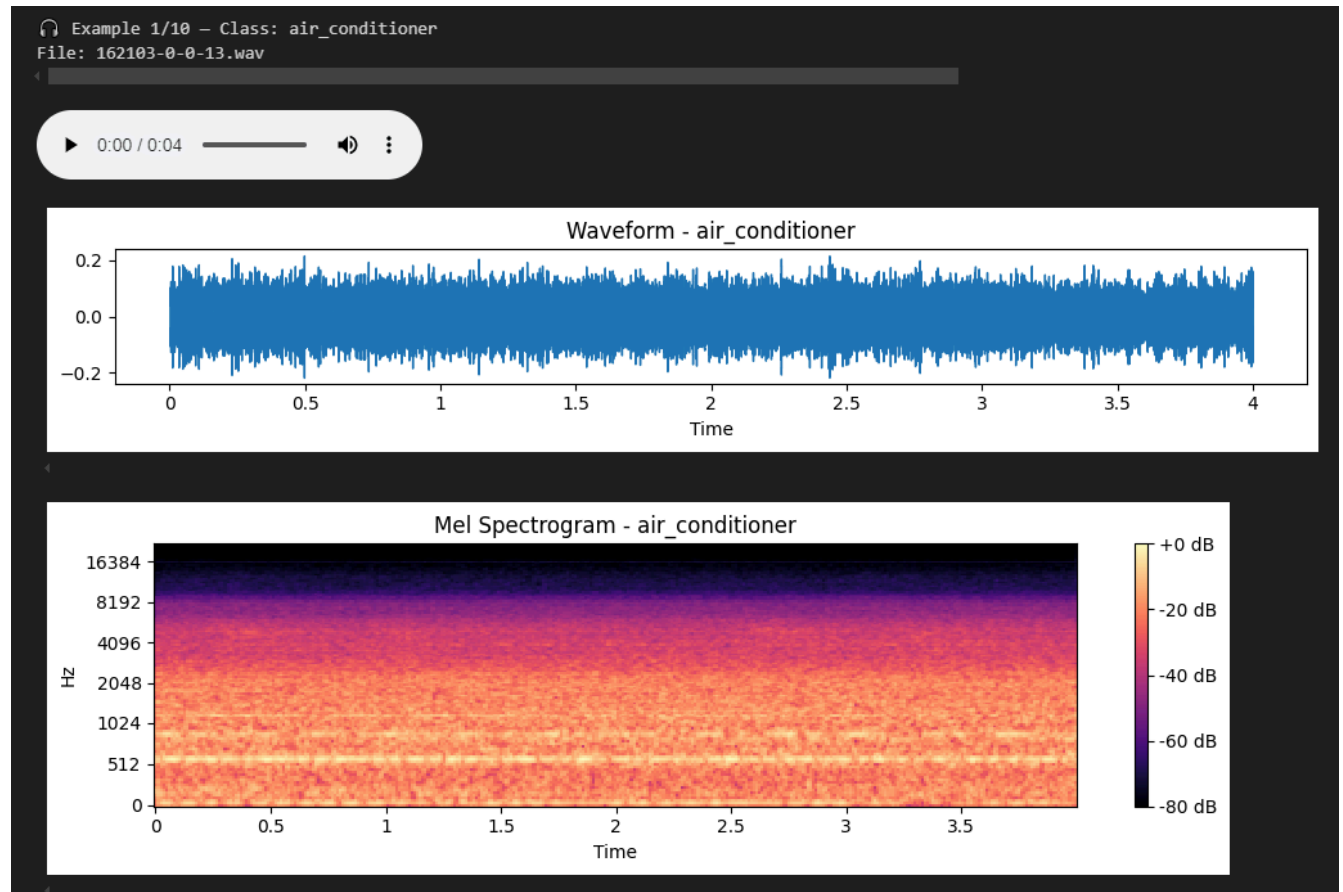
...	slice_file_name	fsID	start	end	salience	fold	classID	class
0	100032-3-0-0.wav	100032	0.000000	0.317551	1	5	3	dog_bark
1	100263-2-0-117.wav	100263	58.500000	62.500000	1	5	2	children_playing
2	100263-2-0-121.wav	100263	60.500000	64.500000	1	5	2	children_playing
3	100263-2-0-126.wav	100263	63.000000	67.000000	1	5	2	children_playing
4	100263-2-0-137.wav	100263	68.500000	72.500000	1	5	2	children_playing
...
8727	99812-1-2-0.wav	99812	159.522205	163.522205	2	7	1	car_horn
8728	99812-1-3-0.wav	99812	181.142431	183.284976	2	7	1	car_horn
8729	99812-1-4-0.wav	99812	242.691902	246.197885	2	7	1	car_horn
8730	99812-1-5-0.wav	99812	253.209850	255.741948	2	7	1	car_horn
8731	99812-1-6-0.wav	99812	332.289233	334.821332	2	7	1	car_horn

8732 rows × 8 columns

```

1 # Visualize and listen to one random sample from each class
2
3 # Number of unique classes
4 num_classes = metadata['class'].nunique()
5
6 # Sample one file from each class
7 samples = metadata.groupby('class').apply(lambda x: x.sample(1)).reset_index(drop=True)
8
9 # Base path to the audio folder
10 AUDIO_BASE_PATH = './UrbanSound8K/audio'
11
12 # Loop over each sample and visualize
13 for i, row in samples.iterrows():
14     file_name = row['slice_file_name']
15     fold = row['fold']
16     class_label = row['class']
17
18     file_path = os.path.join(AUDIO_BASE_PATH, f"fold{fold}", file_name)
19     print(f"\n🔊 Example {i+1}/{num_classes} - Class: {class_label}")
20     print(f"File: {file_name}")
21
22     # Load audio
23     y, sr = librosa.load(file_path, sr=None)
24
25     # Play audio
26     display(ipd.Audio(y, rate=sr))
27
28     # Plot waveform
29     plt.figure(figsize=(10, 2))
30     librosa.display.waveshow(y, sr=sr)
31     plt.title(f"Waveform - {class_label}")
32     plt.tight_layout()
33     plt.show()
34
35     # Plot Mel spectrogram
36     S = librosa.feature.melspectrogram(y=y, sr=sr, n_mels=128)
37     S_DB = librosa.power_to_db(S, ref=np.max)
38
39     plt.figure(figsize=(10, 3))
40     librosa.display.specshow(S_DB, sr=sr, x_axis='time', y_axis='mel')
41     plt.colorbar(format='%+2.0f dB')
42     plt.title(f"Mel Spectrogram - {class_label}")
43     plt.tight_layout()
44     plt.show()
45

```



Step 2: Generating Feature Space:

For each audio file, the following features were extracted:

- **MFCCs (13 coefficients)**
- **Delta and Delta-Delta of MFCCs**
- **Spectral features:** Centroid, Bandwidth, Flatness, Roll-off
- **Zero-Crossing Rate**
- **Root Mean Square Energy (RMS)**

Total feature vector per clip: **45 dimensions**

```

1  def extract_features(signal, sr):
2      features = []
3
4      # MFCCs
5      mfcc = librosa.feature.mfcc(y=signal, sr=sr, n_mfcc=13)
6      mfcc_mean = np.mean(mfcc, axis=1)
7      features.extend(mfcc_mean)
8
9      # Delta + Delta-Delta
10     delta = librosa.feature.delta(mfcc, width=5)
11     delta2 = librosa.feature.delta(mfcc, order=2, width=5)
12     features.extend(np.mean(delta, axis=1))
13     features.extend(np.mean(delta2, axis=1))
14
15     # Spectral features
16     spec_centroid = librosa.feature.spectral_centroid(y=signal, sr=sr)
17     spec_bw = librosa.feature.spectral_bandwidth(y=signal, sr=sr)
18     spec_flat = librosa.feature.spectral_flatness(y=signal)
19     spec_roll = librosa.feature.spectral_rolloff(y=signal, sr=sr)
20
21     features.extend([
22         np.mean(spec_centroid),
23         np.mean(spec_bw),
24         np.mean(spec_flat),
25         np.mean(spec_roll)
26     ])
27
28     # Zero-crossing rate
29     zcr = librosa.feature.zero_crossing_rate(signal)
30     features.append(np.mean(zcr))
31
32     # RMS Energy
33     rms = librosa.feature.rms(y=signal)
34     features.append(np.mean(rms))
35
36     return np.array(features)
37
38  def prepare_dataset(folds):
39      X, y = [], []
40      for _, row in metadata.iterrows():
41          if row['fold'] in folds:
42              file_path = f"./UrbanSound8K/audio/fold{row['fold']}/{row['slice_file_name']}"
43              signal, sr = librosa.load(file_path, sr=None)
44              label = row['classID']
45              features = extract_features(signal, sr)
46              X.append(features)
47              y.append(label)
48      return np.array(X), np.array(y)
49
50  X_train, y_train = prepare_dataset([1,2,3,4,5,6])
51  X_val, y_val = prepare_dataset([7,8])
52  X_test, y_test = prepare_dataset([9,10])
53

```

Step 3: Building the Model:

→ First: Bidirectional LSTM Model

Architecture:

- Input: (45, 1)
- Bidirectional **LSTM** with various hidden sizes and layers meaning the model learns patterns in both forward and backward time directions.
- **hidden_size**: Number of features in the LSTM hidden state per direction.
- Total output size from LSTM: **hidden_size * 2**
- **Fully Connected** After LSTM, it maps the concatenated forward and backward outputs to the final 10 class logits.


```

1 import torch
2 import torch.nn as nn
3 from torch.utils.data import DataLoader, TensorDataset
4 import torch
5 import torch.nn as nn
6 from torch.utils.data import DataLoader, TensorDataset
7
8 class BiRNNModel(nn.Module):
9     def __init__(self, sequence_length=45, input_size=1, hidden_size=64, num_layers=1, num_classes=10):
10         super(BiRNNModel, self).__init__()
11         self.lstm = nn.LSTM(
12             input_size=input_size,
13             hidden_size=hidden_size,
14             num_layers=num_layers,
15             batch_first=True,
16             bidirectional=True
17         )
18         # hidden_size is doubled because of bidirectional
19         self.fc = nn.Linear(hidden_size * 2, num_classes)
20
21     def forward(self, x):
22         # x: [batch_size, seq_len, input_size] = [B, 45, 1]
23         _, (hn, _) = self.lstm(x)
24         forward = hn[-2, :, :] # [B, hidden_size]
25         backward = hn[-1, :, :] # [B, hidden_size]
26         combined = torch.cat((forward, backward), dim=1) # [B, 2*hidden_size]
27         return self.fc(combined)
28
29 def train_model_rnn(model, train_loader, val_loader, epochs=50, lr=0.001):
30     device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
31     model.to(device)
32     criterion = nn.CrossEntropyLoss()
33     optimizer = torch.optim.Adam(model.parameters(), lr=lr)
34
35     for epoch in range(epochs):
36         model.train()
37         for xb, yb in train_loader:
38             xb, yb = xb.to(device), yb.to(device)
39             out = model(xb.float())
40             loss = criterion(out, yb)
41             optimizer.zero_grad()
42             loss.backward()
43             optimizer.step()
44
45         print(f"Epoch {epoch+1}/{epochs} - Loss: {loss.item():.4f}")
46
47     model.eval()
48     correct = 0
49     total = 0
50     with torch.no_grad():
51         for xb, yb in val_loader:
52             xb, yb = xb.to(device), yb.to(device)
53             outputs = model(xb.float())
54             _, predicted = torch.max(outputs.data, 1)
55             total += yb.size(0)
56             correct += (predicted == yb).sum().item()
57
58     val_accuracy = correct / total
59     return model, val_accuracy
60

```

Training:

- Optimizer: Adam
- Epochs: 50
- Loss: CrossEntropy
- Batch size: 32

Hyperparameter Tuning

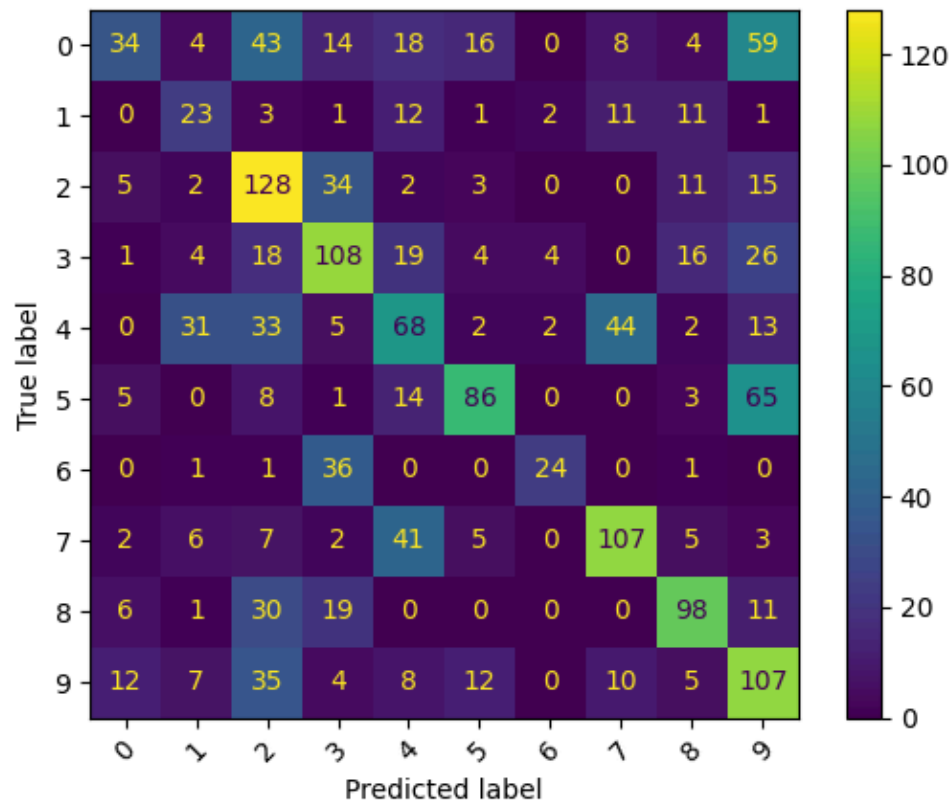
- Trying different combinations of:
 - **hidden_size: [64, 128, 256, 512]**
 - **num_layers: [1, 2, 4, 6]**

Best Configuration:

- `hidden_size=128, num_layers=1`
- **Validation Accuracy: 48.24%**
- **Test Accuracy: 47.37%**
- **Test F1 Score: 46.68%**

Confusion Matrix:

➡ Accuracy: 0.4737, F1 Score: 0.4668



The most **confused pairs** in the model are:

1. **Air Conditioner ↔ Street Music**
2. **Air Conditioner ↔ Dog Bark**
3. **Children Playing ↔ Dog Bark**
4. **Engine Idling ↔ Dog Bark**
5. **Street Music ↔ Dog Bark**

These pairs are what the model finds hard to distinguish, mostly because:

- They have **overlapping frequency patterns**
 - They are often recorded in **noisy urban environments**
 - MFCCs and spectral features may not capture their fine differences
-

→ Second: Transformer Model

A. Without Reshaping

- Input: (45, 1)
- **Best config:** d_model=512, num_heads=4, d_ff=128, num_layers=6
- **Validation Accuracy:** 58.21%
- **Test Accuracy:** 55.66%
- **Test F1 Score:** 54.59%
- **Confusion Matrix:**

✓
0s

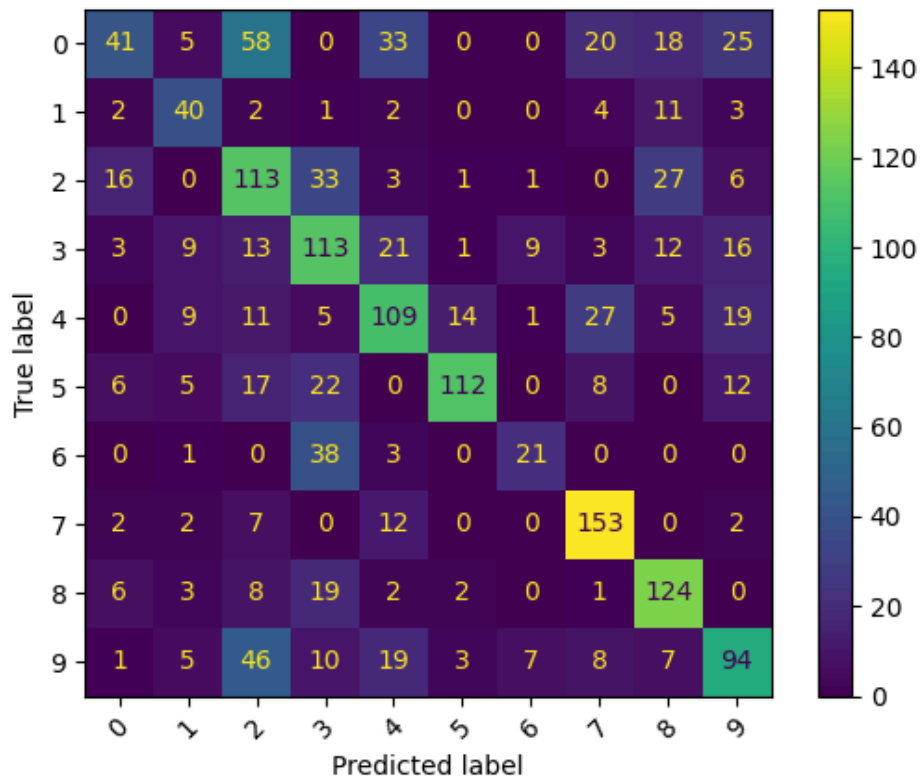


#Scaling without reshaping "one step"

```
evaluate_model(best_config_wor[4], X_test_wor, y_test)
```



Accuracy: 0.5566, F1 Score: 0.5459



The Transformer model showed confusion between several sound classes. The most commonly confused pairs were:

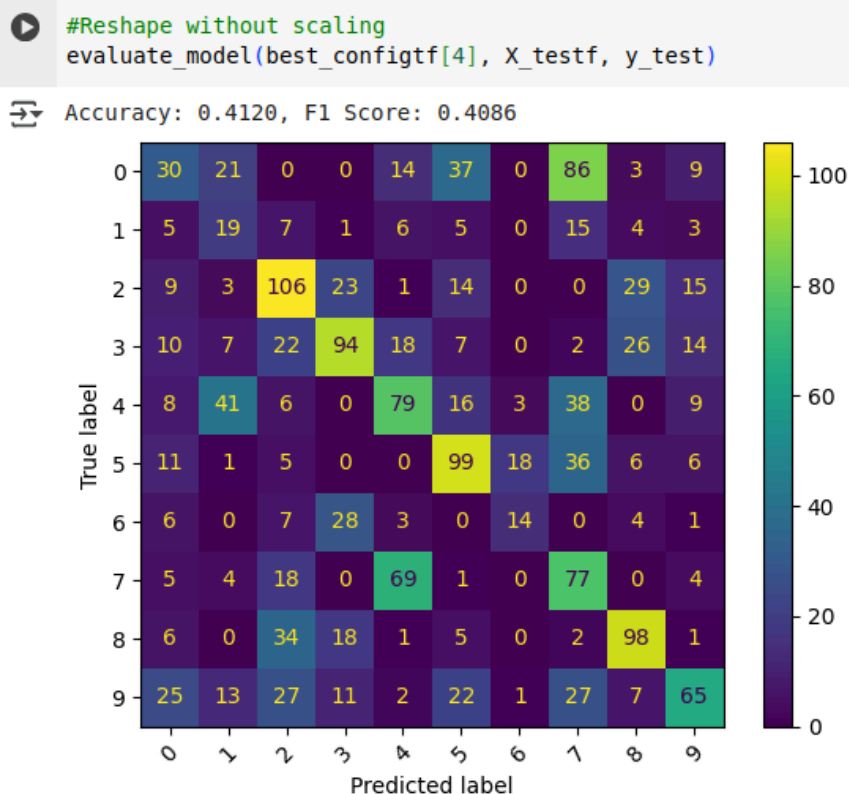
- **Air Conditioner → Children Playing**
- **Street Music → Children Playing**
- **Gun Shot → Dog Bark**
- **Children Playing → Dog Bark**
- **Air Conditioner → Drilling**

These confusions occur because:

- Many classes **share overlapping acoustic features**, such as similar frequency ranges or sound textures (e.g., ambient noise, mechanical hums, or high-energy bursts).
- **Background noise and mixed recordings** in real-world urban environments lead to ambiguous patterns.
- **Feature extraction (e.g., MFCCs)** may not fully capture differences between continuous or impulsive sounds.
- **Transformer limitations** on short sequences may reduce its ability to distinguish subtle time-based patterns.

B. With Reshaping (9, 5)

- Reshaped features to simulate time series
- **Best config:** d_model=256, num_heads=2, d_ff=256, num_layers=4
- **Validation Accuracy:** 33.00%
- **Test Accuracy:** 41.20%
- **Test F1 Score:** 40.86%
- **Confusion Matrix:**



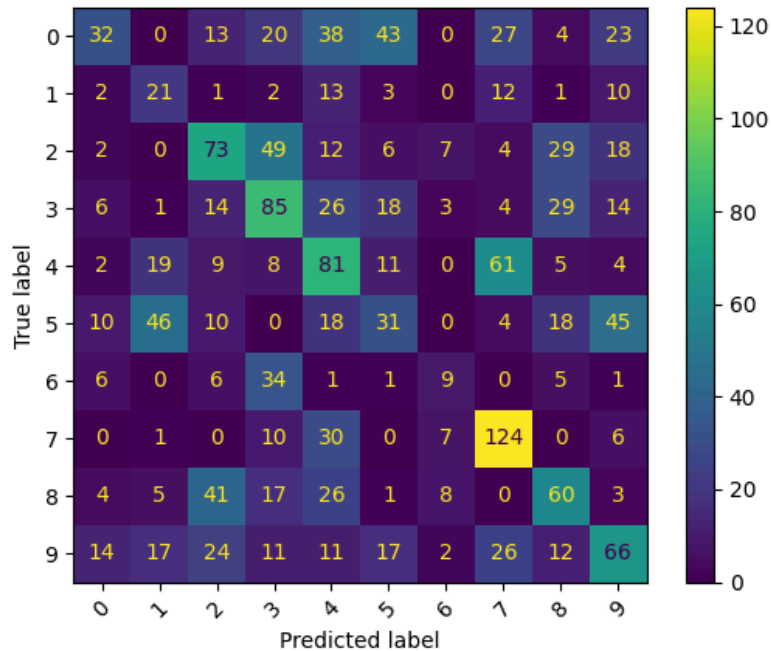
C. With Reshaping + Feature Scaling

- StandardScaler used
- **Best config:** d_model=256, num_heads=4, d_ff=256, num_layers=4
- **Validation Accuracy:** 34.25%
- **Test Accuracy:** 35.21%

- **Test F1 Score: 33.95%**
- **Confusion Matrix:**

```
#Scaling with reshaping
evaluate_model(best_config_scaled[4], X_test_scaled, y_test)
```

Accuracy: 0.3521, F1 Score: 0.3395



→ **Positional Encoding:**

```
1 class PositionalEncoding(nn.Module):
2     def __init__(self, d_model, max_len=500):
3         super().__init__()
4         pe = torch.zeros(max_len, d_model)
5         position = torch.arange(0, max_len).unsqueeze(1).float()
6         div_term = torch.exp(torch.arange(0, d_model, 2).float() * (-torch.log(torch.tensor(10000.0)) / d_model))
7         pe[:, 0::2] = torch.sin(position * div_term)
8         pe[:, 1::2] = torch.cos(position * div_term)
9         self.pe = pe.unsqueeze(0)
10
11     def forward(self, x):
12         return x + self.pe[:, :x.size(1)].to(x.device)
13
```

→ Architecture Class:

```
1 class MultiHeadSelfAttention(nn.Module):
2     def __init__(self, d_model, num_heads):
3         super().__init__()
4         assert d_model % num_heads == 0
5         self.d_k = d_model // num_heads
6         self.num_heads = num_heads
7         self.qkv_proj = nn.Linear(d_model, d_model * 3)
8         self.out_proj = nn.Linear(d_model, d_model)
9         self.dropout = nn.Dropout(0.1)
10
11     def forward(self, x):
12         B, T, C = x.size()
13         qkv = self.qkv_proj(x).chunk(3, dim=-1)
14         q, k, v = [t.view(B, T, self.num_heads, self.d_k).transpose(1, 2) for t in qkv]
15         scores = torch.matmul(q, k.transpose(-2, -1)) / torch.sqrt(torch.tensor(self.d_k, dtype=torch.float32))
16         attn = torch.softmax(scores, dim=-1)
17         out = torch.matmul(attn, v).transpose(1, 2).contiguous().view(B, T, C)
18         return self.out_proj(self.dropout(out))
19
```




```
1 class TransformerBlock(nn.Module):
2     def __init__(self, d_model, num_heads, d_ff):
3         super().__init__()
4         self.attn = MultiHeadSelfAttention(d_model, num_heads)
5         self.norm1 = nn.LayerNorm(d_model)
6         self.ff = nn.Sequential(
7             nn.Linear(d_model, d_ff),
8             nn.ReLU(),
9             nn.Dropout(0.1),
10            nn.Linear(d_ff, d_model)
11        )
12        self.norm2 = nn.LayerNorm(d_model)
13        self.dropout = nn.Dropout(0.1)
14
15    def forward(self, x):
16        x = self.norm1(x + self.attn(x))
17        x = self.norm2(x + self.ff(x))
18        return x
19
```

→ From the notebook:

This class defines a complete Transformer-based model for sequence classification, which consists of:

Input Embedding:

Projects input features from input_dim to d_model using a linear layer.

Positional Encoding:

Adds positional information to the embedded input to retain sequence order.

Stacked Transformer Encoder Blocks:

A configurable number (num_layers) of TransformerBlock layers, each including:

- Multi-head self-attention
- Feed-forward network
- Layer normalization
- Residual connections
- Dropout

Global Sequence Pooling:

Applies AdaptiveAvgPool1d to collapse the sequence dimension and produce a fixed-size feature vector.

Classification Head:

A final linear layer maps the pooled output to num_classes class logits.

```
1 class TransformerClassifier(nn.Module):
2     def __init__(self, input_dim, d_model, num_heads, d_ff, num_layers, num_classes):
3         super().__init__()
4         self.embedding = nn.Linear(input_dim, d_model)
5         self.pos_encoding = PositionalEncoding(d_model)
6         self.layers = nn.ModuleList([TransformerBlock(d_model, num_heads, d_ff) for _ in range(num_layers)])
7         self.pool = nn.AdaptiveAvgPool1d(1)
8         self.classifier = nn.Linear(d_model, num_classes)
9
10    def forward(self, x):
11        x = self.embedding(x)
12        x = self.pos_encoding(x)
13        for layer in self.layers:
14            x = layer(x)
15        x = x.transpose(1, 2)
16        x = self.pool(x).squeeze(-1)
17        return self.classifier(x)
18
```

```
1 def train_model_transformer(model, train_loader, val_loader, epochs=50, lr=1e-4):
2     device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
3     model.to(device)
4     optimizer = torch.optim.Adam(model.parameters(), lr=lr)
5     scheduler = torch.optim.lr_scheduler.StepLR(optimizer, step_size=10, gamma=0.8)
6     criterion = nn.CrossEntropyLoss()
7
8     for epoch in range(epochs):
9         model.train()
10        for xb, yb in train_loader:
11            xb, yb = xb.to(device).float(), yb.to(device)
12            outputs = model(xb)
13            loss = criterion(outputs, yb)
14            optimizer.zero_grad()
15            loss.backward()
16            optimizer.step()
17        scheduler.step()
18        print(f"Epoch {epoch+1}/{epochs} - Loss: {loss.item():.4f}")
19
20    model.eval()
21    correct, total = 0, 0
22    with torch.no_grad():
23        for xb, yb in val_loader:
24            xb, yb = xb.to(device).float(), yb.to(device)
25            preds = model(xb).argmax(dim=1)
26            correct += (preds == yb).sum().item()
27            total += yb.size(0)
28    val_acc = correct / total
29    return model, val_acc
30
31
```

Comparison between the 2 Models:

Metrics	Bi-LSTM	Transformer “without reshaping”
Accuracy	47.37%	55.66%
F1 Score	46.68%	54.59%
