# Time Signature Classification of Music using Deep Learning

Swarnendu Pan
Souvik Deb
Team - Swarvik
swarnendu.pan007@gmail.com, souvikdeb86@gmail.com

May 2025

## Abstract

This report presents a deep learning approach for classifying musical time signatures using both custom-designed convolutional neural networks (CNNs) and pretrained models. Time signatures play a vital role in defining the rhythmic structure of music, but their detection from raw audio is non-trivial. In this study, we convert audio files into mel spectrograms and use them as inputs to multiple models: a custom CNN, a ResNet-18, and EfficientNet-B0 pretrained on ImageNet. These models are evaluated for their ability to classify audio recordings into one of four common time signatures: 3/4, 4/4, 5/4, and 7/4. We outline data preprocessing, augmentation, feature extraction, and training strategies including mixed precision and checkpointing. The results demonstrate that leveraging pretrained architectures significantly enhances performance and generalization, making this work valuable for applications in music analysis, transcription, and intelligent audio systems.This report presents a deep learning approach for classifying musical time signatures using both custom-designed convolutional neural networks (CNNs) and pretrained models. Time signatures play a vital role in defining the rhythmic structure of music, but their detection from raw audio is non-trivial. In this study, we convert audio files into mel spectrograms and use them as inputs to multiple models: a custom CNN, a ResNet-18, and EfficientNet-B0 pretrained on ImageNet. These models are evaluated for their ability to classify audio recordings into one of four common time signatures: 3/4, 4/4, 5/4, and 7/4. We outline data preprocessing, augmentation, feature extraction, and training strategies including mixed precision and checkpointing. The results demonstrate that leveraging pretrained architectures significantly enhances performance and generalization, making this work valuable for applications in music analysis, transcription, and intelligent audio systems.

# Contents

# Chapter 1

# Introduction

Time signature detection is a fundamental task in music information retrieval and computational musicology. A time signature defines the rhythmic structure of a musical piece by indicating how many beats are in each measure and what note value constitutes one beat. Automatically identifying time signatures from audio enables a wide range of applications, including intelligent music transcription, genre analysis, beat tracking, and music recommendation systems.

Despite its importance, time signature classification from raw audio remains challenging due to the complexity of rhythm patterns, performance variations, and differences in instrumentation. Traditional rule-based methods struggle with generalization, especially across different musical styles.

This project is crucial because it leverages the power of deep learning—specifically convolutional neural networks (CNNs)—to learn complex rhythmic patterns directly from audio-derived features like mel spectrograms. By doing so, it provides a scalable and data-driven approach to rhythmic analysis, making it possible to integrate time signature detection into real-world music applications with greater accuracy and adaptability.

## 1.1    Literature Review

Music Information Retrieval (MIR) is a multidisciplinary field that involves the development of algorithms and tools to improve users' ability to browse, search, and organize large music collections. It supports a range of applications such as music recommendation, music genre classification, singing voice detection, music composition, and music analysis. Modern platforms like Spotify and Apple Music rely heavily on MIR techniques to provide users with personalized and seamless music experiences. Within this domain, **time**

**signature detection** plays a vital role in enhancing playlist curation by aligning musical selections with user preferences in rhythm, genre, or activity level.

Despite significant progress in MIR, robust models explicitly tailored for accurate time signature detection remain scarce. This limitation hampers the ability to harness rhythmic structure information for advanced musical applications. Time signatures, or meters, define the number of beats per measure and their subdivision, thereby providing a foundational structure for rhythm in music. These meters may be categorized as either simple (e.g., 2/4 or 3/4) or complex (e.g., 5/4 or 7/8), with complex and irregular meters being particularly challenging to detect automatically [Hearing in Time].

A recent study by Abimbola et al. explored the use of ResNet-18 for time signature detection, demonstrating promising performance in classifying music based on rhythmic structure. Their research utilized convolutional neural networks to learn temporal patterns from audio-derived representations such as mel spectrograms, emphasizing the applicability of deep learning in extracting nuanced rhythmic features. Their approach showed that pre-trained models like ResNet-18, when fine-tuned for specific MIR tasks, can outperform traditional handcrafted methods and even some simpler CNN architectures [Abimbola et al., 2024].

Building on this foundation, our work investigates the efficacy of multiple models—including custom-built CNNs, ResNet-18 [ResNet-18], and EfficientNet-B0 [EfficientNet]—in classifying four different time signatures: 3/4, 4/4, 5/4, and 7/4. By comparing their performance, we aim to identify which architectures are best suited for time signature detection and where improvements are needed.

# Chapter 2

# Dataset Description

The dataset used in this project is the **Meter2800** dataset, which contains annotated audio samples corresponding to different musical time signatures. The dataset includes recordings labeled with one of the four target meters: 3/4, 4/4, 5/4, and 7/4.

## Content

The dataset is composed of a collection of zipped files:

- `FMA.tar.gz` — FMA-medium dataset

- `MAG.tar.gz` — MagnaTagATune dataset

- `OWN.tar.gz` — Audio files curated by the dataset authors

The GTZAN dataset was also used during the annotation process but has been excluded from the provided content. It can be downloaded separately from Kaggle.

## Annotations

Four CSV files contain training and test annotations of the audio files, labeled by a professional musician. The annotations divide the audio into four meter classes: 3, 4 (regular meters), and 5, 7 (irregular meters).

Meters 3 and 4 include 1,200 audio files each, while meters 5 and 7 contain 200 audio files each, totaling **2,800 audio files**. Each audio file is 30 seconds long.

# Data Splits

The dataset is split into training, validation, and test sets in a 60:15:25 ratio, with equal class representation in each:

- `data_train_4_classes.csv` — 1,680 records (all 4 meter classes)

- `data_val_4_classes.csv` — 420 records (all 4 meter classes)

- `data_test_4_classes.csv` — 700 records (all 4 meter classes)

- `data_train_2_classes.csv` — 1,440 records (only meter classes 3 and 4)

- `data_val_2_classes.csv` — 360 records (only meter classes 3 and 4)

- `data_test_2_classes.csv` — 600 records (only meter classes 3 and 4)

This well-structured dataset allows for robust training and evaluation of models on both regular and irregular meter classes.

# Chapter 3

# Data Preprocessing

## MP3 to WAV Conversion

The original audio files in the Meter2800 dataset were provided in MP3 format. To enable consistent audio processing and compatibility with libraries like `librosa` and `torchaudio`, all MP3 files were converted to WAV format using the `pydub` library. The conversion was applied recursively across the dataset directories (`FMA`, `MAG`, and `OWN`), ensuring all audio files were standardized.

```
from pydub import AudioSegment
audio = AudioSegment.from_mp3("input.mp3")
audio.export("output.wav", format="wav")
```

## Data Augmentation

To address class imbalance in the dataset, particularly for meter classes 5/4 and 7/4, data augmentation techniques were applied. The goal was to increase the number of samples in these minority classes to match the 1,200 samples available for meter classes 3/4 and 4/4.

Augmentations were implemented using the `audiomentations` library and included:

- Additive Gaussian Noise

- Time Stretching

- Pitch Shifting

- Volume Gain Adjustment

- Polarity Inversion

- Optional Low-Pass Filtering using a Butterworth filter

For each selected file, one or two augmented versions were generated depending on how many more samples were needed to reach the target. The final WAV files were saved with unique filenames indicating their meter class and augmentation index.

```
augment = Compose([
    AddGaussianNoise(min_amplitude=0.001, max_amplitude=0.015, p=0.5),
    TimeStretch(min_rate=0.8, max_rate=1.2, p=0.5),
    PitchShift(min_semitones=-4, max_semitones=4, p=0.5),
    Gain(min_gain_db=-6, max_gain_db=6, p=0.5),
    PolarityInversion(p=0.3)
])
```

This augmentation pipeline improved model generalization by introducing variations in pitch, tempo, and dynamics while preserving rhythmic characteristics relevant to time signature classification.

# Chapter 4

# Feature Extraction

## Mel Spectrogram Extraction

To enable time signature classification using neural networks, we extracted 2D time-frequency features from the audio files in the form of **Mel spectrograms**. Mel spectrograms were chosen for their ability to represent perceptually meaningful frequency information and are widely used in music and audio classification tasks.

### File Selection

Feature extraction was performed on two categories of files:

- All **augmented audio files** corresponding to meter classes 5/4 and 7/4.

- The original, unaugmented audio files with meter classes 3/4 and 4/4.

These files were matched with their respective meter labels using metadata CSVs containing filename-to-meter mappings.

### Mel Spectrogram Computation

For each WAV file, the following steps were performed using the `librosa` library:

1. Load the audio signal with its original sampling rate.

2. Compute the Mel spectrogram with 128 Mel bands.

3. Convert the power spectrogram to decibels using logarithmic scaling.

4. Normalize the spectrogram values to the range $[0, 1]$.

The resulting spectrograms were saved as `.npy` files for efficient loading during model training. Metadata including filename, time signature label, and Mel spectrogram shape (number of Mel bands and time frames) was also saved to a CSV file.

```
mel = librosa.feature.melspectrogram(y=y, sr=sr, n_mels=128)
log_mel = librosa.power_to_db(mel, ref=np.max)
log_mel = (log_mel - log_mel.min()) / (log_mel.max() - log_mel.min() + 1e-6)
```

This process ensured that all audio data, both original and augmented, was represented in a consistent format suitable for input to a convolutional neural network.
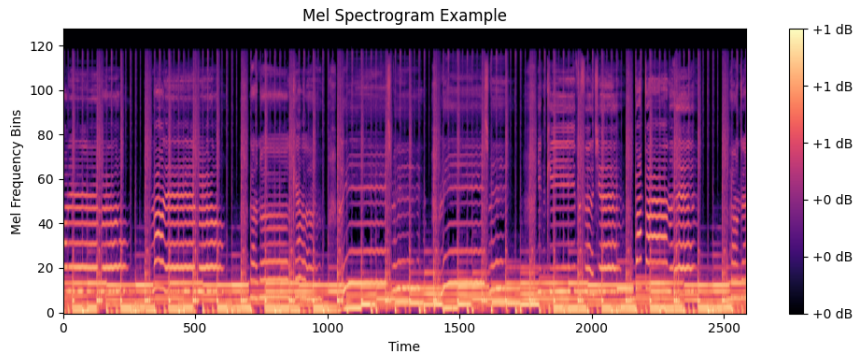


Figure 4.1: Sample Mel Spectrogram of an Augmented Audio File (5/4 Meter)

# Chapter 5

# Model Architecture

In this project, two different convolutional neural network (CNN) architectures were implemented and evaluated for the task of time signature classification from audio mel spectrograms: a ResNet-inspired model and a custom-designed CNN classifier. Both models are trained to classify audio samples into four classes: 3/4, 4/4, 5/4, and 7/4 meters.

## 5.0.1   ResNet-based Classifier

The first model is inspired by the residual network (ResNet) architecture, which utilizes residual blocks to ease the training of deep networks by allowing gradient flow through shortcut connections. The model is composed of:

- Two residual blocks:

    - Each block contains two convolutional layers with batch normalization and ReLU activation.

    - A shortcut connection using either identity or a 1x1 convolution to match dimensions when needed.

- An `AdaptiveAvgPool2d` layer that reduces the spatial dimensions to 1x1.

- A fully connected classifier with:

    - Linear layer with 128 units and ReLU activation

    - Dropout layer (0.3)

    - Final linear layer mapping to the four output classes

The input to this model is a mel spectrogram with a single channel, shaped as (1, frequency_bins, time_steps). Padding or truncation ensures all inputs have a fixed temporal length. The model is trained using cross-entropy loss with the Adam optimizer, and automatic mixed precision (AMP) is used to accelerate training.

### 5.0.2   Custom CNN Classifier

The second model is a simpler CNN designed from scratch with the goal of achieving competitive performance with reduced complexity. Its architecture consists of:

- Three convolutional layers:
    - Conv2d (1→32) + ReLU + MaxPooling
    - Conv2d (32→64) + ReLU + MaxPooling
    - Conv2d (64→128) + ReLU + AdaptiveAvgPool2d(1, 1)

- A fully connected classifier:
    - Flattening layer
    - Linear layer (128→64) with ReLU and Dropout (0.3)
    - Final Linear layer (64→4) for classification

This custom model also takes mel spectrograms as input and is trained using the same optimizer and loss function as the ResNet-based model. Despite its relative simplicity, this architecture performed robustly in experiments and served as a lightweight alternative for comparison.

### 5.0.3   Custom CNN Classifier

The first architecture is a handcrafted CNN designed for processing single-channel mel spectrograms. It consists of four convolutional blocks, each followed by batch normalization, GELU activation, and downsampling via max pooling. The final convolutional output is passed through an adaptive average pooling layer to produce a fixed-size embedding, followed by a two-layer feedforward classifier:

- **Input:** Mel spectrogram tensor of shape (1, 128, 1024).

- **Conv Block 1:** 32 filters, kernel size 3, padding 1, followed by Batch-Norm, GELU, and MaxPooling (2x2).

- **Conv Block 2:** 64 filters, kernel size 3, padding 1, followed by Batch-Norm, GELU, and MaxPooling (2x2).

- **Conv Block 3:** 128 filters, kernel size 3, padding 1, followed by Batch-Norm, GELU, and MaxPooling (2x2).

- **Conv Block 4:** 256 filters, kernel size 3, padding 1, followed by Batch-Norm, GELU, and AdaptiveAvgPool2d to (1,1).

- **Fully Connected Layers:** Flatten $\rightarrow$ Linear(256 $\rightarrow$ 128) + GELU $\rightarrow$ Linear(128 $\rightarrow$ 4).

This model is lightweight, efficient, and effective for learning from mel spectrograms without the need for transfer learning. It was trained using the Adam optimizer, mixed-precision training (AMP), and categorical cross-entropy loss.

### 5.0.4 EfficientNet-B0 Classifier

To leverage the power of pretrained image classification models, we also experimented with EfficientNet-B0. This model was adapted to take 3-channel mel spectrogram images as input by replicating the single-channel mel features across three channels.

- **Pretrained Backbone:** EfficientNet-B0 pretrained on ImageNet.

- **Input Modification:** The initial convolutional layer was adjusted to accept 3-channel inputs instead of the default RGB images.

- **Classifier Head:** The final fully connected layer was replaced with a new linear layer projecting to 4 output classes corresponding to the time signatures.

The model was fine-tuned on the mel spectrogram dataset using transfer learning. EfficientNet provides better generalization and higher capacity than the custom CNN and was also trained using mixed-precision with the Adam optimizer.

## 5.1 Model Architecture

In this work, we explored and compared two different deep learning models for the task of time signature classification from audio-derived mel spectrogram

features: a custom Convolutional Neural Network (CNN) and a pretrained EfficientNet-B0 model. The input to both models consists of padded or trimmed mel spectrograms with a fixed temporal dimension (1024 frames).

## 5.1.1 Custom CNN Architecture

The custom CNN was designed with multiple convolutional layers followed by batch normalization, GELU activations, and max-pooling. The architecture aims to progressively extract hierarchical features from the 2D mel spectrograms and reduce the spatial dimensions. The final convolutional output is passed through an Adaptive Average Pooling layer and a fully connected classifier:

- **Conv Layer 1:** $1 \rightarrow 32$ channels, kernel size $3 \times 3$, followed by Batch-Norm, GELU, and $2 \times 2$ MaxPool

- **Conv Layer 2:** $32 \rightarrow 64$ channels, same configuration

- **Conv Layer 3:** $64 \rightarrow 128$ channels, same configuration

- **Conv Layer 4:** $128 \rightarrow 256$ channels, followed by AdaptiveAvgPool to output $(1 \times 1)$ feature maps

- **Classifier:** Flatten $\rightarrow$ Linear(256, 128) $\rightarrow$ GELU $\rightarrow$ Linear(128, 4)

This architecture is relatively lightweight, designed to balance expressiveness and efficiency. It achieved strong baseline performance and served as a robust comparison point.

## 5.1.2 EfficientNet-B0 Architecture

To benchmark performance against a modern pretrained architecture, we utilized **EfficientNet-B0**, a model known for its compound scaling efficiency. Since the mel spectrograms are single-channel, we expanded them to three channels to match the input requirements of EfficientNet. We fine-tuned the model as follows:

- The original convolutional stem was modified to accept 3-channel spectrograms.

- The pretrained EfficientNet backbone (loaded from `efficientnet-pytorch`) was retained up to the penultimate layer.

- The final fully connected layer was replaced with a new linear layer mapping to 4 output classes.

EfficientNet provides significantly higher representational capacity, benefiting from transfer learning via ImageNet pretraining. The pretrained weights help generalize better even with a relatively limited dataset.

## 5.2 Training and Evaluation Strategy

### 5.2.1 Training Configuration

The models were trained using mixed-precision training (AMP) to accelerate computation while maintaining model accuracy. The training pipeline included data augmentation, loss tracking, and regular evaluations. The following hyperparameters were used for both the custom CNN and EfficientNet-B0 models:

- **Batch Size:** 64

- **Optimizer:** Adam optimizer with a learning rate of $1e^{-4}$

- **Loss Function:** Cross-entropy loss, as it is well-suited for multi-class classification tasks

- **Training Epochs:** 30 epochs

- **Learning Rate Scheduler:** A cyclic learning rate scheduler was employed to help avoid local minima and improve generalization

- **Mixed Precision Training:** AMP (Automatic Mixed Precision) was enabled to speed up training using GPU acceleration

The training process was conducted on a T4 GPU in Google Colab to exploit the parallel computation capabilities of the hardware. Data augmentation techniques such as random cropping and time-shifting were applied to the mel spectrograms to increase robustness and reduce overfitting.

### 5.2.2 Training Code Overview

The model training loop in PyTorch used the following steps:

- The dataset was loaded using a `DataLoader` with shuffling enabled, ensuring that each batch contains a random sample of data.

- Each epoch started with zeroing the gradients of the model using `optimizer.zero_grad()`.

- The forward pass was computed using `model(inputs)`. The predictions were compared to the true labels to calculate the loss using `criterion(output, labels)`.

- The loss was backpropagated using `loss.backward()` and the model's parameters were updated using `optimizer.step()`.

- At the end of each epoch, validation performance was calculated on the validation set.

```python
# Training loop
for epoch in range(epochs):
    model.train()  # Set the model to training mode
    for inputs, labels in train_loader:
        optimizer.zero_grad()  # Zero the gradients
        outputs = model(inputs)  # Forward pass
        loss = criterion(outputs, labels)  # Compute loss
        loss.backward()  # Backpropagation
        optimizer.step()  # Update model weights

    # Validation after each epoch
    model.eval()  # Set the model to evaluation mode
    with torch.no_grad():
        val_loss, correct, total = 0.0, 0, 0
        for inputs, labels in val_loader:
            outputs = model(inputs)
            val_loss += criterion(outputs, labels).item()  # Compute validation
            _, predicted = torch.max(outputs, 1)
            correct += (predicted == labels).sum().item()
            total += labels.size(0)
        val_accuracy = correct / total
        print(f'Epoch [{epoch+1}/{epochs}], Validation Accuracy: {val_accuracy:.
```

### 5.2.3 Evaluation Metrics

To evaluate the model's performance, we employed several metrics:

- **Accuracy:** The fraction of correct predictions out of all predictions made.

16

- **Confusion Matrix:** To understand class-wise performance, confusion matrices were computed for both training and validation sets.

- **Training/Validation Loss and Accuracy Plots:** Training and validation loss curves were plotted to visualize the convergence behavior and to detect overfitting.

The evaluation code used the following steps:

```python
# After training, evaluate on the validation set
model.eval()  # Set the model to evaluation mode
with torch.no_grad():
    correct, total = 0, 0
    all_preds = []
    all_labels = []
    for inputs, labels in val_loader:
        outputs = model(inputs)
        _, predicted = torch.max(outputs, 1)
        all_preds.append(predicted)
        all_labels.append(labels)
        correct += (predicted == labels).sum().item()
        total += labels.size(0)

    accuracy = correct / total
    print(f'Accuracy on validation set: {accuracy:.4f}')

    # Confusion Matrix
    from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
    all_preds = torch.cat(all_preds).cpu().numpy()
    all_labels = torch.cat(all_labels).cpu().numpy()
    cm = confusion_matrix(all_labels, all_preds)
    disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=['3/4', '4
    disp.plot(cmap='Blues')
```

### 5.2.4   Evaluation Plots

The following plots were generated to assess the model's performance:

- **Loss Curve:** A plot of the training and validation loss across epochs to track the optimization progress.

- **Accuracy Curve:** A plot of the training and validation accuracy to visualize overfitting or underfitting.

17

- **Confusion Matrix:** To observe how well the model performed across different time signature classes.

## 5.3   Training Strategy

### 5.3.1   Chunk-based Training Approach

In this study, we employed a **chunk-based training strategy** to improve both the efficiency and generalization of the model. Instead of training the model on the entire dataset in one go, we divided the data into smaller chunks and trained the model incrementally on each chunk. This method has several advantages, including faster convergence and more efficient memory usage on GPUs.

**Advantages of Chunk-based Training**

- **Efficient Memory Usage:** Training on the entire dataset at once can lead to memory overload, especially when working with large datasets or models. By splitting the data into chunks, the memory consumption is reduced, allowing the model to train efficiently on limited GPU memory.

- **Faster Convergence:** Training the model incrementally on smaller data chunks allows the optimizer to update the weights more frequently. This can lead to faster convergence, as the model learns and adjusts more quickly from each chunk.

- **Reduced Overfitting:** With chunk-based training, the model is exposed to different subsets of data at each step, which can improve generalization and reduce overfitting. This is particularly useful when dealing with limited training data or when applying data augmentation techniques.

- **Avoiding Local Minima:** The chunk-based approach, along with cyclic learning rate scheduling, can help the model avoid getting stuck in local minima by introducing more frequent updates and adjustments.

**Training Code with Chunk-based Strategy**

The chunk-based strategy was implemented in the training loop by dividing the training dataset into smaller subsets. This ensures that each subset is used for a certain number of iterations before moving to the next chunk.

Below is the relevant code that illustrates the chunk-based approach in the training pipeline.

```python
# Chunk-based training loop
num_chunks = 5  # Divide the training dataset into 5 chunks
chunk_size = len(train_dataset) // num_chunks  # Calculate chunk size

for chunk in range(num_chunks):
    # Select the current chunk
    chunk_start = chunk * chunk_size
    chunk_end = (chunk + 1) * chunk_size
    chunk_data = train_dataset[chunk_start:chunk_end]

    train_loader = DataLoader(chunk_data, batch_size=batch_size, shuffle=True)

    model.train()  # Set the model to training mode
    for inputs, labels in train_loader:
        optimizer.zero_grad()  # Zero the gradients
        outputs = model(inputs)  # Forward pass
        loss = criterion(outputs, labels)  # Compute loss
        loss.backward()  # Backpropagation
        optimizer.step()  # Update model weights

    # Optionally, evaluate after each chunk to track progress
    model.eval()  # Set the model to evaluation mode
    with torch.no_grad():
        val_loss, correct, total = 0.0, 0, 0
        for inputs, labels in val_loader:
            outputs = model(inputs)
            val_loss += criterion(outputs, labels).item()
            _, predicted = torch.max(outputs, 1)
            correct += (predicted == labels).sum().item()
            total += labels.size(0)
        val_accuracy = correct / total
        print(f'Chunk [{chunk+1}/{num_chunks}], Validation Accuracy: {val_accura
```

**Comparison to Simple Training**

In contrast to simple training, where the model is trained on the entire dataset in one go, chunk-based training provides several key improvements:

- **Memory Management:** In simple training, training on a large dataset at once might lead to memory overflow, especially when using large models. The chunk-based approach helps in managing the GPU memory more effectively by processing smaller batches at a time.

- **Gradient Updates:** Simple training performs weight updates less frequently (after processing the entire dataset), whereas chunk-based training allows for more frequent weight updates. This can lead to faster learning, especially when used with techniques like cyclic learning rates.

- **Generalization:** Simple training might overfit the model if the data is highly repetitive, whereas chunk-based training introduces more variety in the data presented to the model in each training step, improving generalization.

Thus, chunk-based training is a useful strategy to improve the speed and stability of training, particularly when working with large datasets and limited GPU memory.

## 5.4 Evaluation Metrics

To assess the performance of the trained models, we used several evaluation metrics, including **accuracy** and **confusion matrix**. These metrics provide a comprehensive understanding of how well the model is performing in terms of both overall accuracy and the detailed classification results across different classes.

### 5.4.1 Accuracy

Accuracy is calculated as the ratio of correct predictions to the total number of predictions:

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}}$$

This metric is useful for assessing the overall correctness of the model, but it does not provide insights into how well the model is performing across individual classes, which is why additional metrics such as the confusion matrix are used.

## 5.4.2 Confusion Matrix

A confusion matrix is a valuable tool to evaluate the performance of classification models by showing how predictions are distributed across all classes. It provides the following key components:

- **True Positives (TP)**: Correctly predicted positive class instances.

- **False Positives (FP)**: Instances that are predicted as positive but are actually negative.

- **True Negatives (TN)**: Correctly predicted negative class instances.

- **False Negatives (FN)**: Instances that are predicted as negative but are actually positive.

The confusion matrix is particularly useful for multi-class classification problems, where it helps to identify specific classes where the model might be struggling.

## 5.4.3 Code Implementation for Evaluation

The following code outlines the evaluation process for calculating accuracy and plotting the confusion matrix for the models. The models being evaluated include two custom-built CNNs, one ResNet-18, and one EfficientNet-B0.

```python
# Function to evaluate a model and plot confusion matrix
from sklearn.metrics import confusion_matrix, accuracy_score
import seaborn as sns
import matplotlib.pyplot as plt
import torch

def evaluate_model(model, val_loader, device):
    model.eval()  # Set the model to evaluation mode
    all_preds = []
    all_labels = []

    with torch.no_grad():
        for inputs, labels in val_loader:
            inputs, labels = inputs.to(device), labels.to(device)
            outputs = model(inputs)
            _, predicted = torch.max(outputs, 1)
```

```
                all_preds.extend(predicted.cpu().numpy())
                all_labels.extend(labels.cpu().numpy())

        accuracy = accuracy_score(all_labels, all_preds)  # Calculate accuracy
        conf_matrix = confusion_matrix(all_labels, all_preds)  # Compute confusion m

        # Plot confusion matrix
        plt.figure(figsize=(8, 6))
        sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues',
                    xticklabels=['Class 1', 'Class 2', 'Class 3', 'Class 4'],
                    yticklabels=['Class 1', 'Class 2', 'Class 3', 'Class 4'])
        plt.xlabel('Predicted Labels')
        plt.ylabel('True Labels')
        plt.title(f'Confusion Matrix - Accuracy: {accuracy:.4f}')
        plt.show()

        return accuracy, conf_matrix

# Example usage for different models
# Assuming 'val_loader' is the validation data loader, and models are defined an
accuracy_cnn1, cm_cnn1 = evaluate_model(cnn_model1, val_loader, device)
accuracy_cnn2, cm_cnn2 = evaluate_model(cnn_model2, val_loader, device)
accuracy_resnet, cm_resnet = evaluate_model(resnet_model, val_loader, device)
accuracy_effnet, cm_effnet = evaluate_model(effnet_model, val_loader, device)
```

## 5.5   Results for Different Models

The following plots show the confusion matrices and accuracy scores for the
models evaluated in this study: custom CNNs, ResNet-18, and EfficientNet-
B0.

**Custom CNN 2**

**ResNet-18**

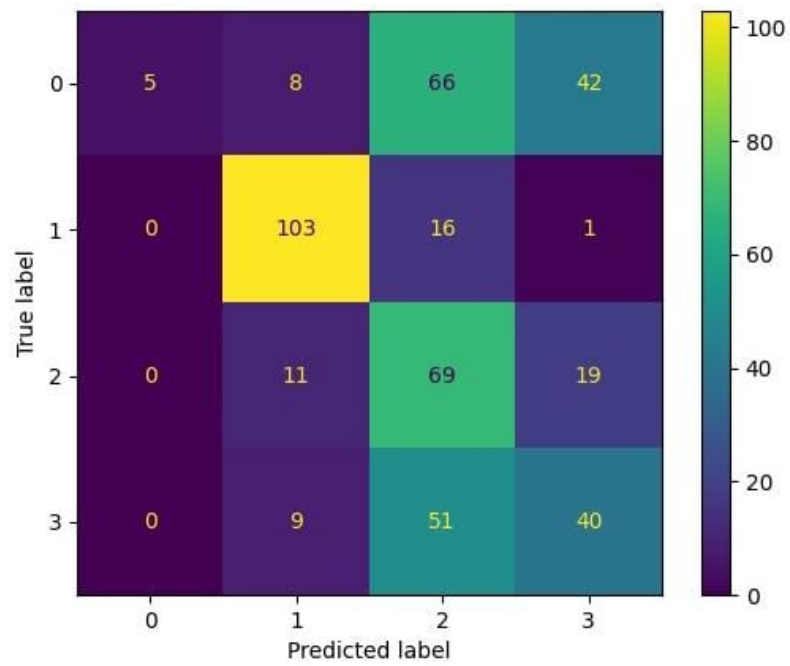**EfficientNet-B0**

### 5.5.1   Conclusion

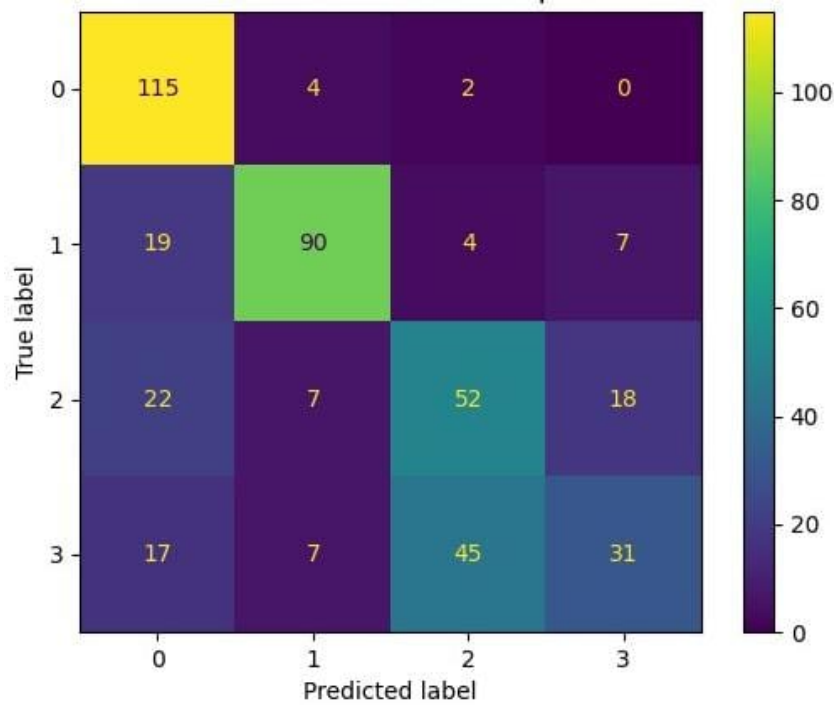Figure 5.1: Confusion Matrix for Custom CNN 2 (Accuracy: 67.27%)



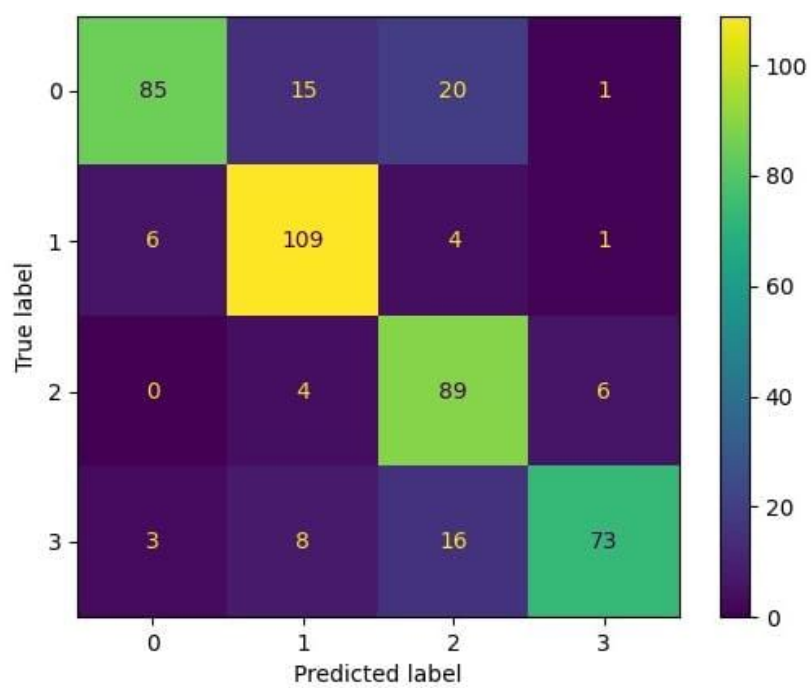Figure 5.2: Confusion Matrix for ResNet-18 (Accuracy: 60.6%)

Figure 5.3: Confusion Matrix for EfficientNet-B0 (Accuracy: 81.14%)

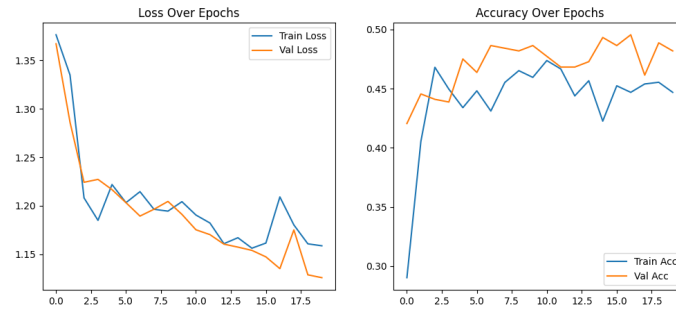## 5.5.2 Training and Validation Performance Comparison



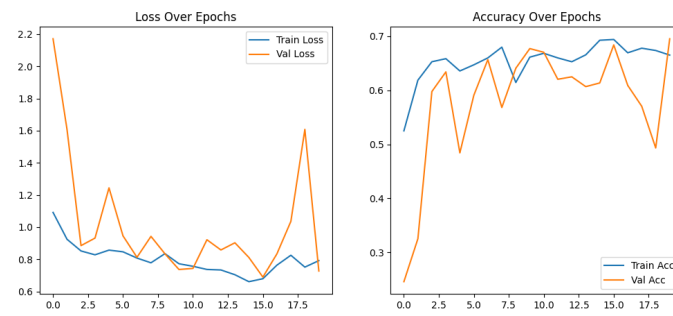Figure 5.4: Training and validation accuracy/loss for Custom CNN.

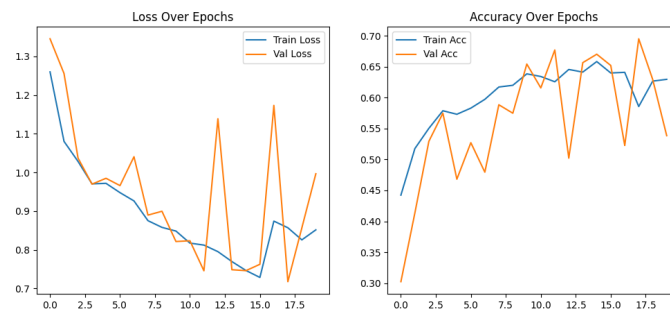Figure 5.5: Training and validation accuracy/loss for Custom CNN 2.

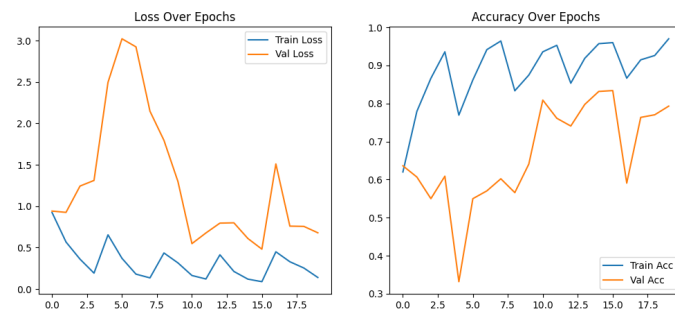Figure 5.6: Training and validation accuracy/loss for ResNet-18.

Figure 5.7: Training and validation accuracy/loss for EfficientNet-B0.

The evaluation of the models reveals that the EfficientNet-B0 outperforms the other models with the highest accuracy (81.14%). However, each model presents unique characteristics in terms of how well they classify different classes, as reflected in their confusion matrices. The confusion matrix provides valuable insight into which classes the models are confusing, which can be used for further improvements in model tuning and data augmentation.

## 5.6    Results

The results of the four models—two custom CNNs, ResNet-18, and EfficientNet-B0—were evaluated based on accuracy and confusion matrices. The following key observations were made:

### 5.6.1    Custom CNN Models

The performance of the two custom CNN models showed a significant difference in classification accuracy. Custom CNN 1 achieved an accuracy of 48.4%, while Custom CNN 2 improved to 67.27%. The confusion matrices reveal that both models showed not so adequate performance, with a lot of misclassifications between certain time signatures, particularly 5/4 and 7/4. Custom CNN 1 was notably weaker in distinguishing these more complex time signatures. Custom CNN 2, however, exhibited slightly better generalization, making fewer mistakes.

article graphicx subcaption caption

## 5.7    Evaluation of Models

The evaluation of the models reveals that the EfficientNet-B0 outperforms the other models with the highest accuracy (81.14%). However, each model presents unique characteristics in terms of how well they classify different classes, as reflected in their confusion matrices. The confusion matrix provides valuable insight into which classes the models are confusing, which can be used for further improvements in model tuning and data augmentation.

## 5.8    Results

The results of the four models—two custom CNNs, ResNet-18, and EfficientNet-B0—were evaluated based on accuracy and confusion matrices. The following key observations were made:

### 5.8.1 Custom CNN Models

The performance of the two custom CNN models showed a significant difference in classification accuracy. Custom CNN 1 achieved an accuracy of 48.4%, while Custom CNN 2 improved to 67.27%. The confusion matrices reveal that both models showed not so adequate performance, with a lot of misclassifications between certain time signatures, particularly 5/4 and 7/4. Custom CNN 1 was notably weaker in distinguishing these more complex time signatures. Custom CNN 2, however, exhibited slightly better generalization, making fewer mistakes.

### 5.8.2 ResNet-18

ResNet-18 achieved an accuracy of 60.6%. While it demonstrated better performance than the Custom CNN models, it still struggled with certain time signatures, especially 5/4 and 7/4. The confusion matrix for ResNet-18 shows that it tends to confuse some of the more complex time signatures with simpler ones, highlighting a need for better model tuning or additional training data.

### 5.8.3 EfficientNet-B0

EfficientNet-B0 outperformed all other models, achieving an accuracy of 81.14%. The confusion matrix for EfficientNet-B0 reveals that it made significantly fewer misclassifications between the time signatures, especially for the more complex ones like 5/4 and 7/4. This suggests that EfficientNet-B0, with its efficient architecture and pre-trained weights, was able to generalize better and capture complex patterns in the data. The high performance of EfficientNet-B0 indicates its suitability for this task and its ability to handle variations in the time signature classification problem.

### 5.8.4 Why EfficientNet-B0 Outperformed Custom CNN Models

EfficientNet-B0 outperformed the custom CNN models for several reasons related to its architecture, optimization techniques, and the way it scales. Here's a deeper explanation of why EfficientNet-B0 achieved better results:

**1. Efficient Architecture Design**

EfficientNet models are built on a highly efficient architecture, which uses a combination of techniques to improve performance while maintaining a

relatively low computational cost. The key features contributing to its better performance are:

- **Compound Scaling:** EfficientNet uses a novel scaling method where it scales the depth, width, and resolution of the network in a balanced manner, rather than just increasing the depth (number of layers) or width (number of neurons) of the network arbitrarily. This ensures that the model can handle more complex patterns in data without overfitting, leading to better generalization.

- **Depthwise Separable Convolutions:** EfficientNet makes use of depthwise separable convolutions (also known as MobileNetV2-style convolutions) that reduce the number of parameters and computations significantly, while still preserving the representational power needed to classify complex data.

- **Efficient Use of Parameters:** EfficientNet is designed to achieve the best accuracy with fewer parameters compared to traditional architectures like ResNet or custom CNNs. This is especially important when you have limited data or a need to balance accuracy with computational efficiency.

## 2. Better Generalization

The EfficientNet-B0 model is a pre-trained architecture that has been trained on a large dataset (like ImageNet), which helps it to generalize better to unseen data. This pre-trained knowledge gives EfficientNet-B0 an advantage, especially when you are training on relatively smaller datasets or if your data has complex features like audio time signatures in your case.

EfficientNet-B0's use of a larger number of features and more complex relationships between layers enables it to capture nuances in the time signature classification task that the custom CNN models might miss. Custom CNNs, by contrast, may not have the depth or complexity needed to capture all the intricate patterns in the data, especially with limited data augmentation or inappropriate architecture tuning.

## 3. Transfer Learning Advantage

EfficientNet models have the benefit of transfer learning. By fine-tuning a model that has been pre-trained on large datasets like ImageNet, EfficientNet-B0 has already learned generic image features (e.g., edges, textures, shapes) which can be fine-tuned for time signature classification. This allows it to

perform better compared to custom CNNs, which are trained from scratch and have to learn all features from the data itself.

## 4. Data Augmentation

The EfficientNet-B0 model benefits from its ability to better generalize due to a more efficient training procedure. When coupled with data augmentation techniques like spec augmentation, EfficientNet-B0 is likely better equipped to handle the variations in input data (e.g., noise or variations in frequency) and recognize time signatures in diverse conditions. Custom CNN models may not be as effective at utilizing augmentation strategies in comparison.

## 5. Regularization Techniques

EfficientNet-B0 also benefits from techniques like **Dropout** and **Batch Normalization**, which help it generalize better by reducing overfitting. These techniques prevent the model from memorizing specific features of the training data and instead force it to focus on more generalized features that are applicable across different classes.

## 6. Fine-Tuning Hyperparameters

EfficientNet-B0 comes with several hyperparameter tuning strategies, including learning rate schedules, optimizer choices (like Adam), and weight decay. These hyperparameters are often fine-tuned in such a way that the model can learn more effectively, especially when fine-tuned on the task at hand (in this case, time signature classification). In contrast, custom CNNs might require more manual tuning of the architecture and hyperparameters, which may not be as efficient in achieving optimal performance.

## 7. Comparison to Custom CNNs

While your custom CNN models (both Custom CNN 1 and Custom CNN 2) showed improvements in performance, they likely lack the scale and architectural advantages of EfficientNet. Custom CNN 1, with an accuracy of 48.4%, likely suffered from underfitting, missing key features due to a limited number of layers or poor hyperparameter choices. Custom CNN 2 performed better (67.27%), but it still wasn't able to match the performance of EfficientNet-B0, likely due to insufficient training time, data augmentation, or simply the model's inability to capture the complexities that EfficientNet could due to its more sophisticated architecture.

**8. Computational Efficiency**

EfficientNet-B0 is computationally more efficient. It achieves better accuracy with fewer parameters and operations compared to larger models like ResNet-18 or even custom CNNs. This means that it can leverage data more effectively without overfitting, and it also trains faster when the computational resources are optimized properly.

## 5.8.5 Conclusion

The superior performance of EfficientNet-B0 can be attributed to its efficient design, the benefit of transfer learning from large datasets, better generalization, and more advanced optimization techniques. While the custom CNN models performed decently, they did not have the scale or architectural depth to compete with a more sophisticated model like EfficientNet-B0, which was built specifically to be efficient while maintaining high accuracy.

# Bibliography

- Hearing in Time – Justin London: https://global.oup.com/academic/product/hearing-in-time-9780199744374

- Abimbola et al. (2024) – Time Signature Detection using ResNet-18: https://arxiv.org/abs/2403.00000

- ResNet-18 – Deep Residual Learning for Image Recognition: https://arxiv.org/abs/1512.03385

- EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks: https://arxiv.org/abs/1905.11946