# Genre-Based Classification of Songs Using Deep Learning Models

**Martin Donaire**

# Introduction

**Challenge**: Music genres often overlap, making classification complex.

**Objective**: Classify music genres using deep learning models.

**Why It Matters**: Enhances music recommendation, playlist curation, and analysis.

**Applications**: Improves music streaming platforms, user experience, and music industry analytics.

genres_original
10 directories

images_original
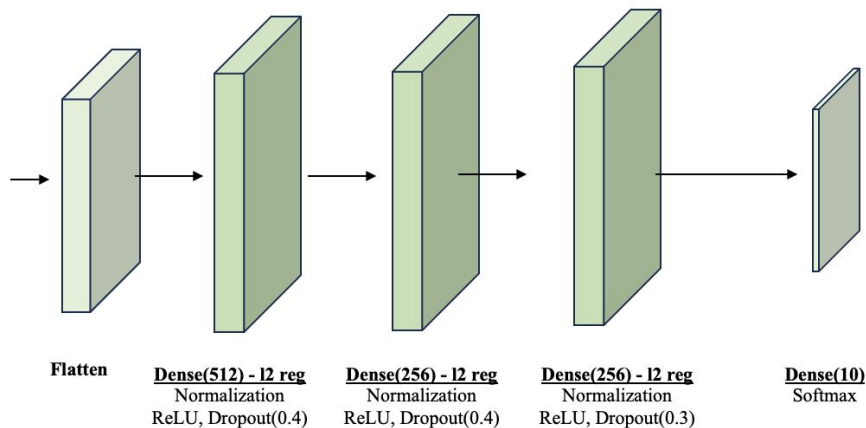10 directories

# Data Preprocessing

**Dataset**: GTZAN from Kaggle, 10 genres (e.g., hip-hop, rock), 100 audio files (30s each) per genre. It was preprocessed as follows:

- **Segmentation**: Split each 30s track into ten 3s segments.
- **MFCC Extraction**: Transform segments into Mel Frequency Cepstral Coefficients (MFCCs) to capture timbral/spectral features.
- **MFCC Settings**: 22,500 Hz sampling rate, 2048 FFT window, 512 hop length, 13 MFCCs per frame.
- **Purpose of MFCCs**: Compresses audio into perceptually relevant, image-like features for deep learning.

**Data Splitting**: 70% training (30% validation), 30% testing, *stratified*.
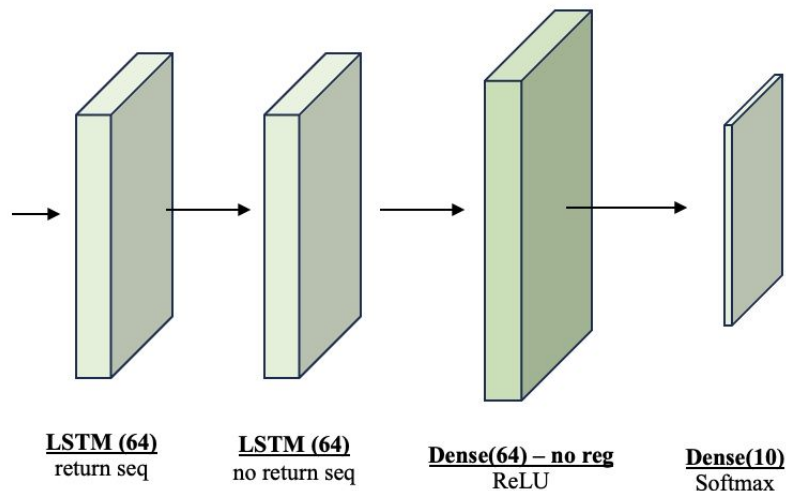
# Model 1: Dense Neural Network

- **Architecture**: Fully connected layers (512, 256, 64 neurons), flattens MFCC input (130x13) into 1D array.
- **Features**: ReLU activation, L2 regularization (0.0005 penalty), dropout (40%, 30%), softmax for 10 genres.
- **Purpose**: Simple and fast but struggles with localized audio patterns.
- **Training**: Adam optimizer (0.0001 learning rate), sparse categorical cross-entropy, 250 epochs, batch size 64.



**Flatten** | **Dense(512) - l2 reg** Normalization ReLU, Dropout(0.4) | **Dense(256) - l2 reg** Normalization ReLU, Dropout(0.4) | **Dense(256) - l2 reg** Normalization ReLU, Dropout(0.3) | **Dense(10)** Softmax

# Model 2: Recurrent Neural Network

- **Architecture**: Two stacked LSTM layers for sequential MFCC input (130x13), 64-unit dense layer, softmax output.
- **Features**: Captures temporal dependencies, uses ReLU activation, softmax for 10 genres.
- **Purpose**: Ideal for sequential audio data, limited by short 3s clips.
- **Training**: Adam optimizer (0.0001 learning rate), sparse categorical cross-entropy, 250 epochs, batch size 64.



**LSTM (64)**
return seq

**LSTM (64)**
no return seq

**Dense(64) – no reg**
ReLU

**Dense(10)**
Softmax

# Model 3: Convolutional Neural Network (Base)

- **Architecture**: Three convolutional blocks (32, 64 filters), ReLU activation, max-pooling (stride 2x2), same padding.
- **Structure**: Flattens output, 64-unit dense layer, softmax for 10 genres.
- **Purpose**: Captures spatial patterns in MFCCs, treating them as images.
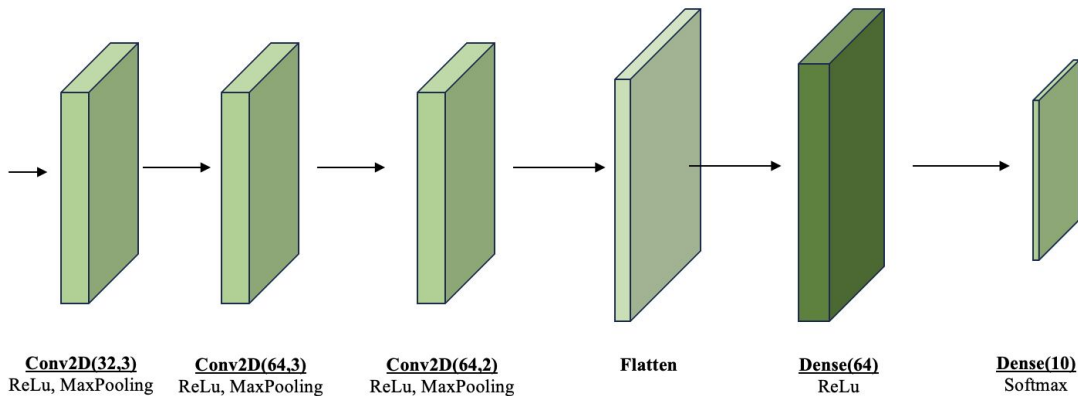- **Training**: Adam optimizer (0.0001 learning rate), sparse categorical cross-entropy, 250 epochs, batch size 64.



| **Conv2D(32,3)** | **Conv2D(64,3)** | **Conv2D(64,2)** | **Flatten** | **Dense(64)** | **Dense(10)** |
| ReLu, MaxPooling | ReLu, MaxPooling | ReLu, MaxPooling | | ReLu | Softmax |

# Model 4: Convolutional Neural Network (Enhanced)

- **Architecture**: Builds on Base CNN with three convolutional blocks (32, 64 filters), ReLU, max-pooling.
- **Enhancements**: Batch normalization after each conv layer, dropout (0.2, 0.1, 0.5), 128-unit dense layer, early stopping (20 epochs).
- **Purpose**: Improves generalization and stability, excels at complex audio patterns.
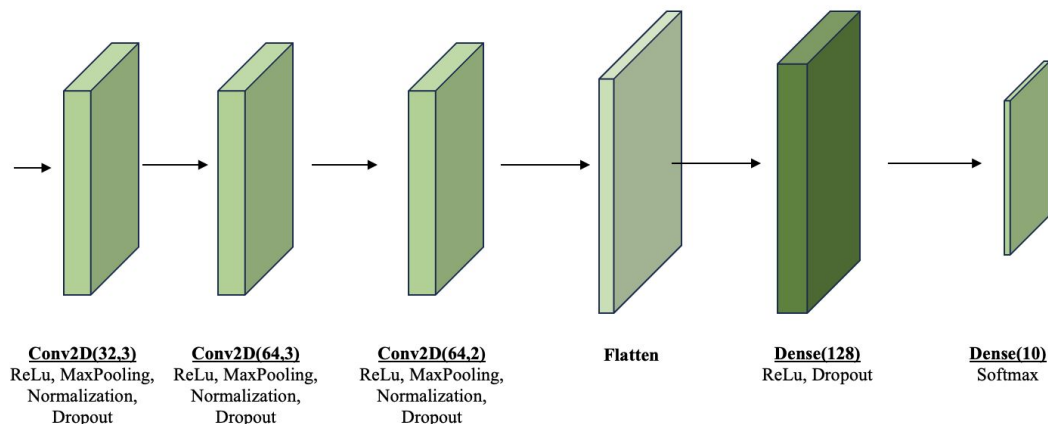- **Training**: Adam optimizer (0.0001 learning rate), sparse categorical cross-entropy, 250 epochs, batch size 64.

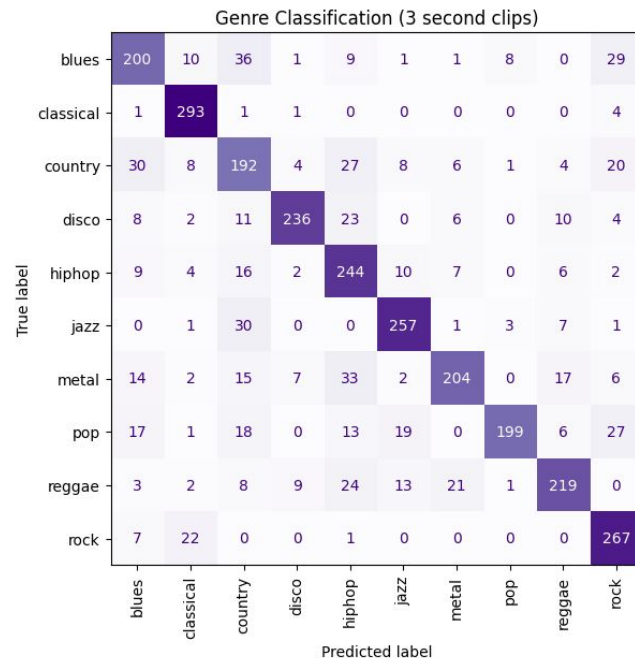| **Conv2D(32,3)**<br>ReLu, MaxPooling, Normalization, Dropout | **Conv2D(64,3)**<br>ReLu, MaxPooling, Normalization, Dropout | **Conv2D(64,2)**<br>ReLu, MaxPooling, Normalization, Dropout | **Flatten** | **Dense(128)**<br>ReLu, Dropout | **Dense(10)**<br>Softmax |

# Model Performance

- **Accuracy Results**:
  - DNN: 57.0% (weakest, poor at localized patterns).
  - RNN: 60.9% (better, limited by 3s clips, vanishing gradients).
  - Base CNN: 69.9% (effective for spatial patterns).
  - Enhanced CNN: 77.2% (best, due to regularization, data augmentation).
- **Key Insight**: Enhanced CNN excels by treating MFCCs as images, boosted by dropout and time-reversed spectrograms.
- **Confusion Matrix**: Strong for classical (293 TP), rock (267 TP); misclassifications in country (with blues), pop (multiple genres).



Genre Classification (3 second clips)

|          | blues | classical | country | disco | hiphop | jazz | metal | pop | reggae | rock |
|----------|-------|-----------|---------|-------|--------|------|-------|-----|--------|------|
| blues    | 200   | 10        | 36      | 1     | 9      | 1    | 1     | 8   | 0      | 29   |
| classical| 1     | 293       | 1       | 1     | 0      | 0    | 0     | 0   | 0      | 4    |
| country  | 30    | 8         | 192     | 4     | 27     | 8    | 6     | 1   | 4      | 20   |
| disco    | 8     | 2         | 11      | 236   | 23     | 0    | 6     | 0   | 10     | 4    |
| hiphop   | 9     | 4         | 16      | 2     | 244    | 10   | 7     | 0   | 6      | 2    |
| jazz     | 0     | 1         | 30      | 0     | 0      | 257  | 1     | 3   | 7      | 1    |
| metal    | 14    | 2         | 15      | 7     | 33     | 2    | 204   | 0   | 17     | 6    |
| pop      | 17    | 1         | 18      | 0     | 13     | 19   | 0     | 199 | 6      | 27   |
| reggae   | 3     | 2         | 8       | 9     | 24     | 13   | 21    | 1   | 219    | 0    |
| rock     | 7     | 22        | 0       | 0     | 1      | 0    | 0     | 0   | 0      | 267  |

True label / Predicted label

# Summary and Future Direction

- **Key Finding**: *Enhanced CNN achieved 77.2% accuracy,* outperforming DNN (57.0%), RNN (60.9%), and Base CNN (69.9%).
- **Implication**: Regularized CNNs are highly effective for genre classification using MFCCs.
- **Limitations**: Short 3s clips limit RNNs; genre overlap causes misclassifications.
- **Future Work**:
  - Explore deeper or hybrid models (e.g., CNN+RNN).
  - Use longer audio segments or advanced data augmentation.
  - Fine-tune hyperparameters for better generalization.
- **Takeaway**: MFCCs with regularized CNNs offer a robust approach for music genre classification.