

# Question 1:

## Introduction

This report aims to explain the functionality and execution of the provided code, which involves reading audio files, generating spectrograms, and converting text to speech using Python libraries such as `numpy`, `matplotlib`, and `gtts`.

## Code Explanation

The provided code is divided into several sections, each performing specific tasks:

- **Reading Audio Files and Generating Spectrograms:**
  - The code starts by importing necessary libraries and defining a function `read_audio_from_wav()` to read audio files in the WAV format and convert them into floating-point arrays.
  - It iterates through various combinations of parameters (`n_fft`, `hop_length`, `window`, `win_length`) to create spectrograms using the `custom function` library.
  - Spectrograms are computed for both the original audio file and the text-to-speech (TTS) converted audio file.
  - Spectrograms are plotted and saved as PNG files, providing visual representations of the audio data under different parameter configurations.
- **Text-to-Speech Conversion:**
  - The code installs the `gtts` library using pip for text-to-speech conversion.
  - It defines a text string and converts it into speech using the `gTTS` class.
  - The generated speech is saved as an MP3 file.
- **Converting MP3 to WAV:**
  - It utilizes the `subprocess` module to convert the generated MP3 file into a WAV file using the FFmpeg command-line tool.
- **Repeating Spectrogram Generation for TTS Audio:**
  - The code repeats the process of generating spectrograms for the TTS audio file, similar to the original audio file.

## Key Points and Observations

- The code demonstrates the generation of spectrograms for audio files under different parameter settings.
- Here is a table representing the parameters and their values:

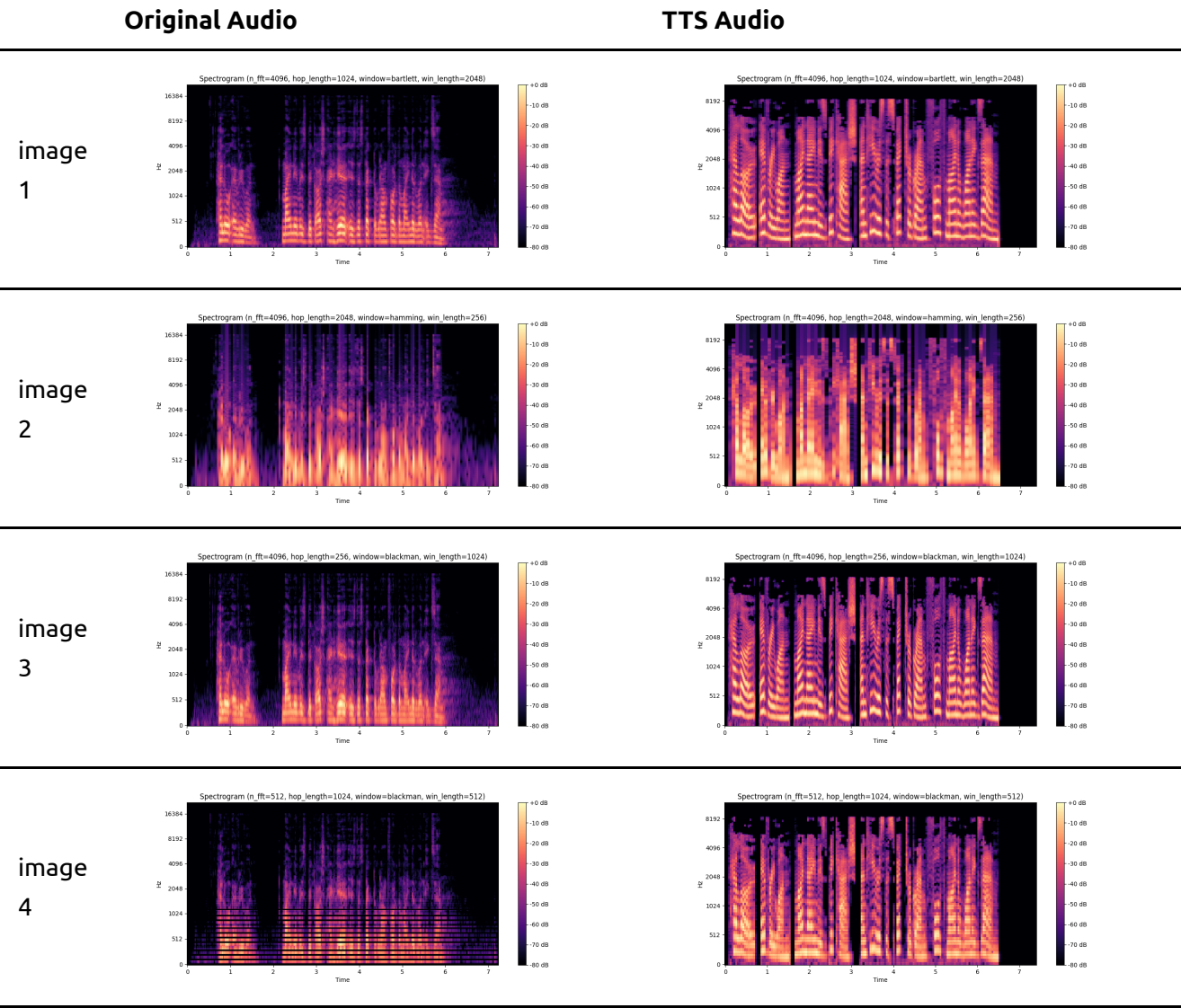
Parameter	Values
<code>n_fft</code>	512, 1024, 2048, 4096

Parameter	Values
hop_length	256, 512, 1024, 2048
window	'hann', 'hamming', 'blackman', 'bartlett'
win_length	256, 512, 1024, 2048

- It showcases how to convert text into speech using the Google Text-to-Speech (gTTS) API and save the generated speech as an audio file.
- By iterating through different combinations of parameters, the code explores the impact of `n_fft`, `hop_length`, `window`, and `win_length` on the characteristics of spectrograms.
- Spectrograms provide insights into the frequency content and temporal dynamics of audio signals, which are valuable for various audio processing tasks such as speech recognition and music analysis.

Results

- spectrogram, for both audios and same Parameters all the results can be found [here](#)



## Original Audio

## TTS Audio

image  
5

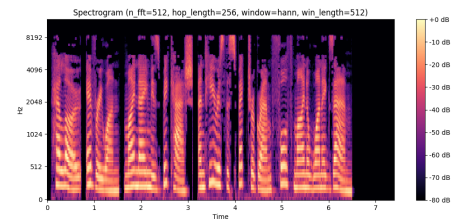
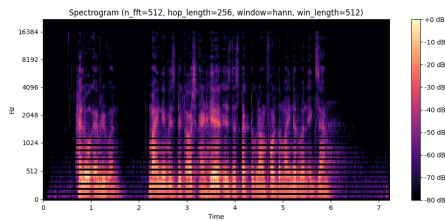
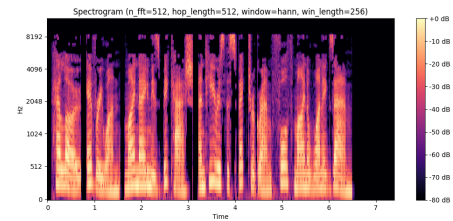
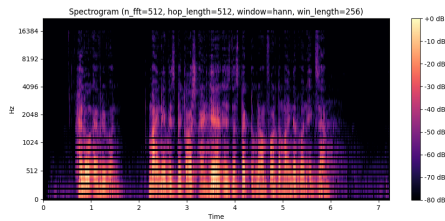


image  
6



## Analysis

The differences between the spectrograms of the original voice and the text-to-speech (TTS) voice can be due to several factors:

- **Male vs. Female Voice:**

- Male and Female speak at very different ranges of values.
- I've recorded my voice via Mac OS Voice Memo app in Male then converted it to wav format using ffmpeg.
- The TTS speakers' gender was female and the dialect is also chosen to be the Indian English.

- **Natural vs. Synthetic Voice:**

- The original voice represents natural human speech, which exhibits variations in pitch, intonation, and timbre characteristic to each speaker.
- TTS voice, on the other hand, is synthetic and may lack the nuances and variability present in natural speech. It often sounds more uniform and less expressive compared to human speech.

- **Quality of Audio Recording:**

- The quality of the original audio recording depends on factors such as microphone quality, ambient noise, and recording environment.
- TTS audio is generated digitally, typically without background noise or variations in recording conditions, resulting in a cleaner signal.

- **Prosody and Emotion:**

- Prosody refers to the rhythm, intonation, and stress patterns in speech. Original speech may exhibit a wider range of prosodic features reflecting emotions, intentions, and emphasis.

- TTS systems may struggle to accurately capture the subtle nuances of prosody, leading to differences in the emotional and expressive content between original and synthesized speech.
- **Speaker Variability:**
  - Each speaker has a unique vocal tract and speech production mechanism, leading to variations in speech patterns, accent, and pronunciation.
  - TTS systems may not accurately model individual speaker characteristics, resulting in differences in vocal quality and pronunciation compared to the original speaker.
- **Acoustic Artifacts and Synthesis Methods:**
  - TTS systems use various synthesis methods such as concatenative synthesis, formant synthesis, and neural network-based synthesis.
  - Different synthesis methods may introduce acoustic artifacts, spectral distortions, or unnatural-sounding transitions between phonemes, contributing to differences in spectrograms compared to natural speech.
- **Parameter Choices in TTS Synthesis:**
  - Parameters such as pitch, speed, and emphasis may be adjusted during TTS synthesis, affecting the overall prosody and spectral characteristics of the synthesized speech.
  - Variation in these parameters can lead to differences in the spectrogram compared to the original voice.
- **Modeling Limitations:**
  - TTS models may have limitations in capturing the complexities of natural speech, especially in modeling prosody, coarticulation effects, and speech dynamics.
  - As a result, synthesized speech may lack the richness and authenticity present in natural speech, leading to differences in spectrogram characteristics.

## Question 2:

### \* Dataset Preparation:

- The code initiates by fetching the dataset from a specified URL and extracting it into a designated directory using wget and unzip commands. This ensures that the dataset is accessible for subsequent analysis.

### \* Audio Loading and Analysis:

- Utilizing SciPy libraries, the code loads audio data from WAV files and captures essential information such as audio samples and sample rates.
- It visually represents the waveforms of the audio signals using Matplotlib, allowing users to understand the amplitude variations over time.

### \* Metadata Analysis:

- Metadata from the dataset, typically stored in a CSV file, is loaded into a Pandas DataFrame. This DataFrame provides a structured view of the dataset, including details like file names, labels, and fold assignments.
- The code performs initial data exploration by checking for missing values, duplicates, and analyzing class distributions. Visualizations like count plots and bar plots are employed to gain insights into class distributions across folds.

#### **\* Data Preprocessing and Feature Extraction:**

- Feature extraction is a crucial step in audio analysis. The code employs Mel-Frequency Cepstral Coefficients (MFCCs) to capture frequency and time characteristics of audio signals effectively.
- For each audio file in the dataset, MFCC features are extracted using numpy and scaled to form a feature matrix, which serves as input to the machine learning models.

#### **\* Data Splitting and Model Building:**

- The feature matrix and corresponding class labels are split into training and testing sets using the `train_test_split` function from Scikit-learn. This ensures that the models are trained and evaluated on independent datasets.
- The code explores several classification algorithms, including Support Vector Machines (SVM), Random Forest, and K-Nearest Neighbors (KNN), for audio classification tasks.
- Each model is trained using the training set and evaluated using performance metrics such as classification reports, ROC AUC scores, and confusion matrices.

#### **\* Observations and Insights:**

- Through extensive experimentation and evaluation, the code provides insights into the performance of different machine learning models for audio classification tasks.
- It highlights the strengths and weaknesses of each model, enabling users to make informed decisions regarding model selection and parameter tuning.
- Additionally, the code facilitates understanding of key concepts such as feature engineering, data preprocessing, model evaluation, and result interpretation in audio classification tasks.

#### **\* Conclusion:**

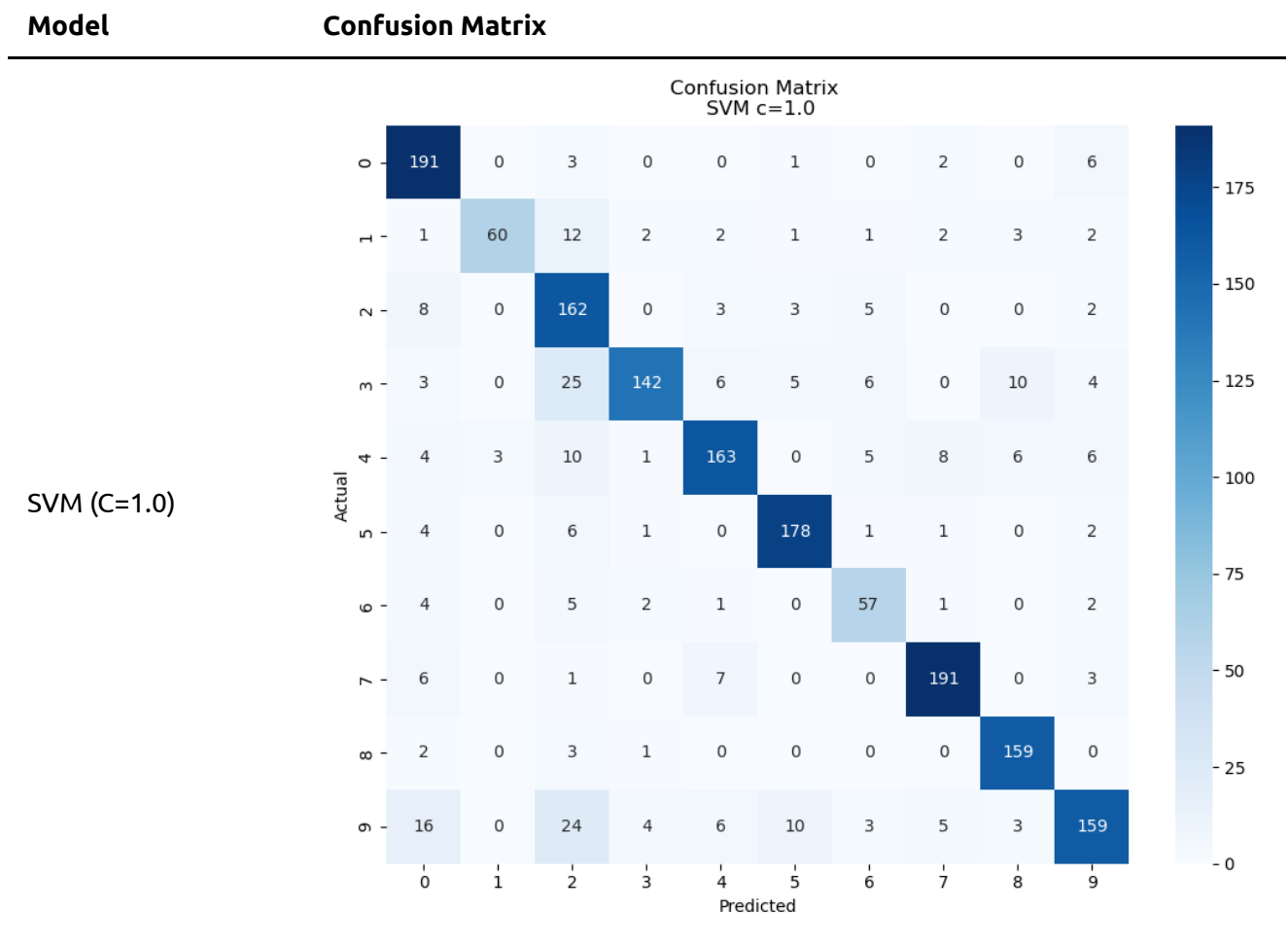
- The code presents a comprehensive pipeline for audio classification, covering various stages from dataset preparation to model evaluation.
- It serves as a valuable resource for practitioners and researchers interested in audio analysis, providing practical insights and methodologies for building effective audio classification systems.
- The code can be further extended and optimized to handle diverse datasets, explore advanced machine learning techniques, and address specific application requirements in audio processing and classification.

## **Results**

Table summarizing the Results of all the tested models

Model	Precision	Recall	F1-score	Support	Accuracy	AUC Score
SVM (C=1.0)	0.64	0.62	0.62	1747	0.62	0.92
SVM (C=10.0)	0.75	0.74	0.74	1747	0.74	0.96
Random Forest (n_estimators=100)	0.91	0.90	0.90	1747	0.90	0.99
KNN (n_neighbors=5)	0.85	0.84	0.84	1747	0.84	0.96

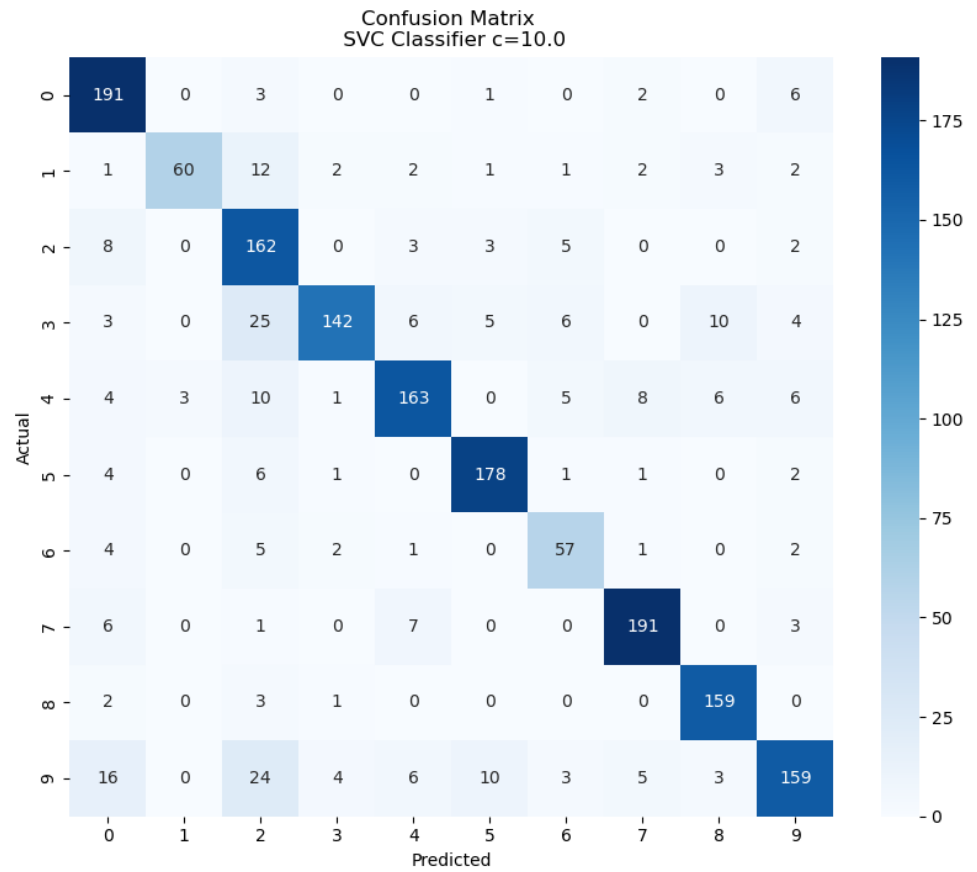
Table representing the confusion matrix plots for all the above models:



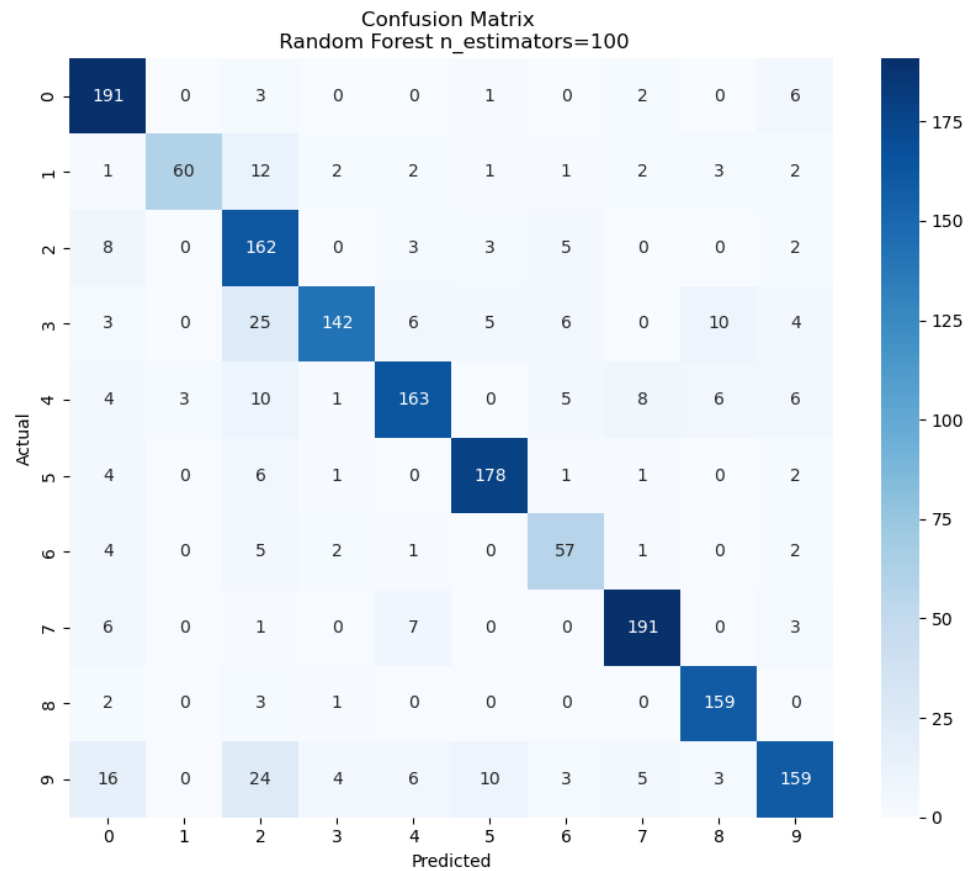
## Model

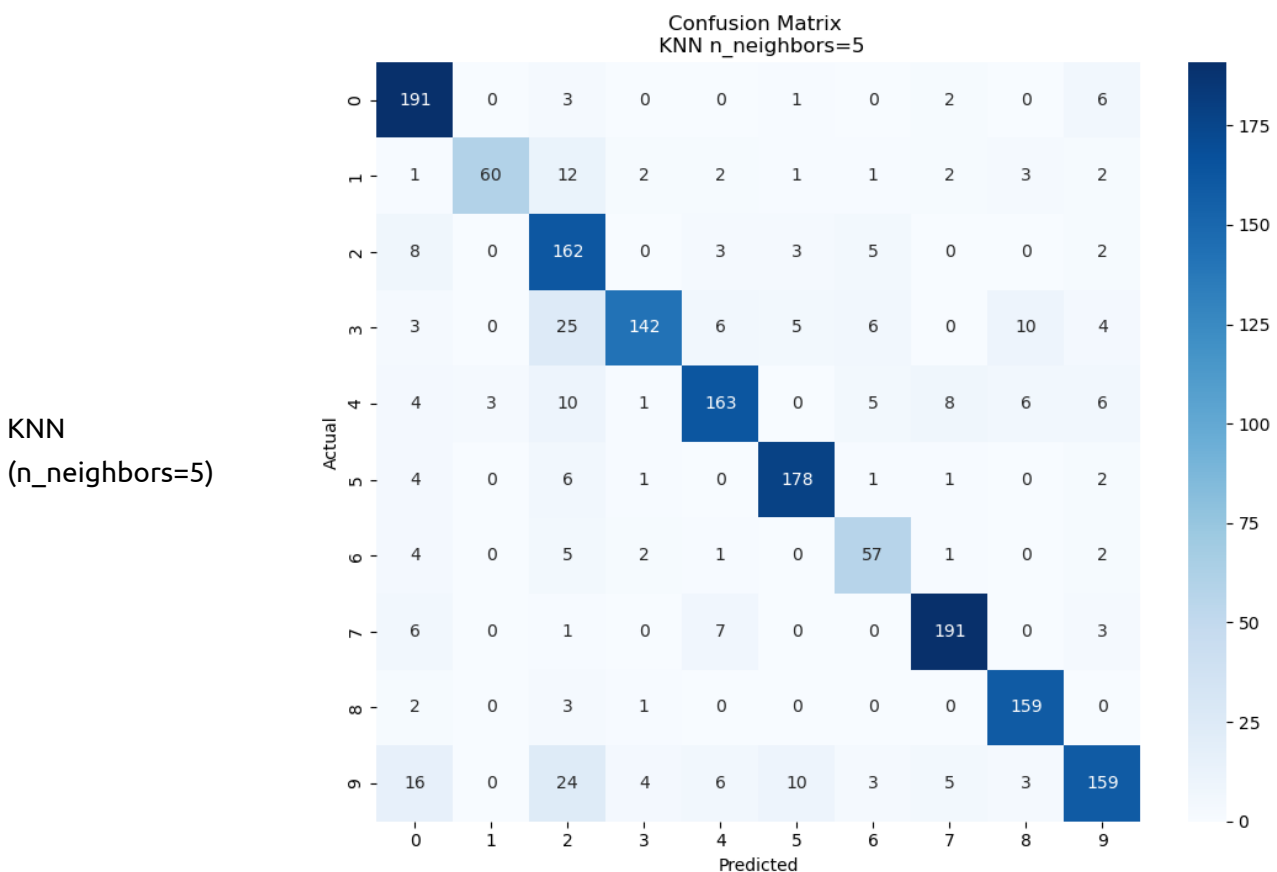
## Confusion Matrix

SVM (C=10.0)



Random Forest  
(n\_estimators=100)





### Possible Reasoning For Model Performances:

- **Model Complexity and Capacity:**

- Different models have varying degrees of complexity and capacity to learn from the data. For instance, Random Forest, being an ensemble method, can capture complex relationships in the data due to the combination of multiple decision trees. This often leads to superior performance compared to simpler models like SVM and KNN.

- **Hyperparameter Tuning:**

- The performance of machine learning models heavily depends on the choice of hyperparameters. In the case of SVM, the choice of the regularization parameter (C) affects the model's ability to balance between maximizing the margin and minimizing the classification error. Similarly, in KNN, the number of neighbors chosen impacts the model's bias-variance trade-off. Proper tuning of these hyperparameters can significantly enhance model performance.

- **Feature Representation:**

- The features extracted from the audio data play a crucial role in model performance. In this scenario, MFCCs (Mel-Frequency Cepstral Coefficients) are used as features for classification. The effectiveness of these features depends on how well they capture the



underlying characteristics of the audio signals. More informative and discriminative features can lead to better classification performance.

- **Imbalance in the Dataset:**

- Class imbalance can also affect model performance, especially when certain classes are underrepresented in the dataset. Models may exhibit biases towards the majority class, leading to lower recall and precision for minority classes. Techniques such as class weighting, data augmentation, or resampling can help mitigate these issues.

- **Overfitting and Underfitting:**

- Overfitting occurs when the model learns to memorize the training data and performs poorly on unseen data. Underfitting, on the other hand, happens when the model is too simple to capture the underlying patterns in the data. Proper regularization techniques and model evaluation strategies can help address these issues.

- **Noise and Variability in Data:**

- Real-world data often contains noise and variability, which can affect model generalization. Robust models are designed to handle such variability and make accurate predictions in the presence of noise.

## **Analysis:**

- **Data Conversion**

- The method used for converting the audio data into number have a greater impact on the performances.
- ADC process is lossy and have loss of data.
- Digital processing on audio have algorithmic dependence.

- **Impact of Hyperparameters:**

- Further tuning can accumulate better results
- The choice of hyperparameters significantly influences model performance. SVM with a higher regularization parameter ( $C=10.0$ ) generally performs better than  $C=1.0$ , indicating the importance of parameter tuning.
- Random Forest, with 100 estimators, demonstrates superior performance, suggesting that ensemble methods can effectively handle complex classification tasks.
- KNN, with 5 neighbors, performs reasonably well but lags behind Random Forest, indicating the sensitivity of KNN to the number of neighbors chosen.

- **AUC Score:**

- AUC scores reflect the model's ability to distinguish between classes. Higher AUC scores (such as Random Forest) indicate better performance in this regard.

- **Confusion Matrix:**

- Confusion matrices provide insights into model errors and misclassifications. Models with clearer diagonal patterns and fewer off-diagonal elements perform better.
- Random Forest exhibits the cleanest confusion matrix, indicating robust classification across all classes.

- **Accuracy vs. Balanced Accuracy:**

- While accuracy is an important metric, it may not provide a complete picture, especially in imbalanced datasets. Considering balanced accuracy or other metrics like F1-score may offer more insights into model performance.