# WGU D604 Task 1: Audio Classification Project Report

## Environmental Sound Classification Using Deep Learning

**Student Name:** Shanikwa Haynes
**Course:** D604 - Advanced Analytics
**Date:** July 15, 2025

---

## Executive Summary

This report presents a comprehensive audio classification project using the ESC-50 dataset and convolutional neural networks (CNNs) to address smart city environmental monitoring needs. The project successfully demonstrates the application of deep learning techniques for automated environmental sound classification, achieving practical performance levels suitable for urban deployment.

---

## Table of Contents

---

## A. Research Question and Business Problem

### A1. Research Question

**Primary Research Question:** How effectively can a Convolutional Neural Network (CNN) classify environmental sounds from the ESC-50 dataset to enable automated urban noise monitoring systems for smart city applications?

This research question directly addresses real-world organizational needs as urban planners and environmental monitoring agencies require automated sound classification systems for intelligent city management. The question is relevant to realistic organizational situations where cities need to monitor noise pollution, detect emergency situations, and analyze urban activity patterns through automated audio analysis.

## A2. Objectives and Goals

**Primary Objective:** Develop a CNN model capable of classifying 50 different environmental sound categories with high accuracy for deployment in smart city monitoring systems.

**Secondary Objectives:** 1. Achieve minimum 70% classification accuracy on ESC-50 test set 2. Optimize audio preprocessing pipeline for real-time deployment 3. Create reproducible model architecture for scalable implementation 4. Provide comprehensive evaluation of model performance across sound categories

Each objective is reasonable within the scope of the research question and represented in the ESC-50 dataset, which contains 2,000 audio samples across 50 environmental sound categories. The objectives are measurable, achievable, and directly support the overarching research question.

## A3. Neural Network Type

**Selected Network:** Convolutional Neural Network (CNN) for audio classification

This is an industry-relevant deep learning architecture capable of performing audio classification tasks by processing mel-spectrogram representations of audio signals. CNNs are specifically suitable for environmental sound classification on the ESC-50 dataset due to their ability to detect spatial patterns in 2D spectral representations of audio data.

## A4. Neural Network Justification

**Justification for CNN Selection:**

1. **Spatial Pattern Recognition:** CNNs excel at identifying spatial patterns in 2D representations, making them ideal for analyzing mel-spectrograms derived from audio signals. The convolutional layers can detect frequency patterns and temporal relationships within spectrograms.

2. **Translation Invariance:** CNNs can detect audio features regardless of their temporal position in the spectrogram, which is crucial for environmental sound classification where sound patterns may occur at different times within the 5-second audio clips.

3. **Hierarchical Feature Learning:** Multi-layer CNN architecture automatically learns hierarchical features from basic spectral patterns to complex sound signatures,

eliminating the need for manual feature engineering that would be required with traditional machine learning approaches.

4. **Industry Standard:** CNNs are widely adopted in audio processing applications, with proven effectiveness in similar environmental sound classification tasks (Salamon & Bello, 2017). This established track record provides confidence in the approach.

5. **Scalability:** CNN architecture can be efficiently deployed on various hardware platforms, from edge devices to cloud servers, supporting smart city infrastructure requirements for distributed monitoring systems.

# B. Image Dataset (Not Applicable)

This section is not applicable as the selected scenario focuses on audio classification using the ESC-50 dataset.

# C. Video Dataset (Not Applicable)

This section is not applicable as the selected scenario focuses on audio classification using the ESC-50 dataset.

# D. Audio Data Preparation

## D1. Exploratory Data Analysis

### D1a. Spectrogram Explanation

A spectrogram is a visual representation of the spectrum of frequencies in an audio signal as it varies with time. Key characteristics include: - **X-axis:** Time (seconds) - **Y-axis:** Frequency (Hz) - **Color intensity:** Magnitude/Power of frequencies - **Reveals:** temporal and spectral patterns in audio signals - **Essential for:** audio classification as it converts 1D audio into 2D image-like representation

Spectrograms are crucial for CNN-based audio classification because they transform temporal audio signals into spatial representations that convolutional layers can process effectively. The mel-scale spectrogram specifically emphasizes perceptually relevant frequency ranges, making it particularly suitable for environmental sound analysis.

### D1b. Audio Tagging Explanation

Audio tagging is the process of assigning descriptive labels to audio segments or files. In the ESC-50 dataset: - Each audio file has a corresponding category label (e.g., 'dog', 'rain',

'siren') - Labels are stored in metadata CSV file with target numerical IDs - Tags enable supervised learning by providing ground truth for training - Essential for classification tasks to map audio features to semantic categories

The ESC-50 dataset provides pre-labeled audio samples across 50 environmental sound categories, enabling supervised learning approaches for automated classification.

## D2. Spectrogram Extraction

The spectrogram extraction process involves several technical steps:

1. **Audio Loading:** Load audio files at 22,050 Hz sampling rate for 5-second duration
2. **Mel-Spectrogram Computation:** Apply mel-frequency filtering to extract 128 mel-frequency bins
3. **Logarithmic Scaling:** Convert power spectrogram to decibel scale using log transformation
4. **Temporal Framing:** Generate 216 time frames for consistent temporal resolution

The extraction process uses librosa library functions to compute mel-spectrograms, which focus on perceptually relevant frequency ranges while providing compact representations suitable for neural network processing.

## D3. Audio Signal Creation

Audio signal creation involves standardizing raw audio files into consistent formats:

1. Load audio files using librosa with specified sample rate (22,050 Hz)
2. Standardize duration to 5 seconds (110,250 samples)
3. Apply padding or truncation to ensure consistent length
4. Normalize amplitude values to prevent clipping and ensure consistent input range

This process ensures all audio inputs have identical dimensions and characteristics, which is essential for batch processing in neural networks.
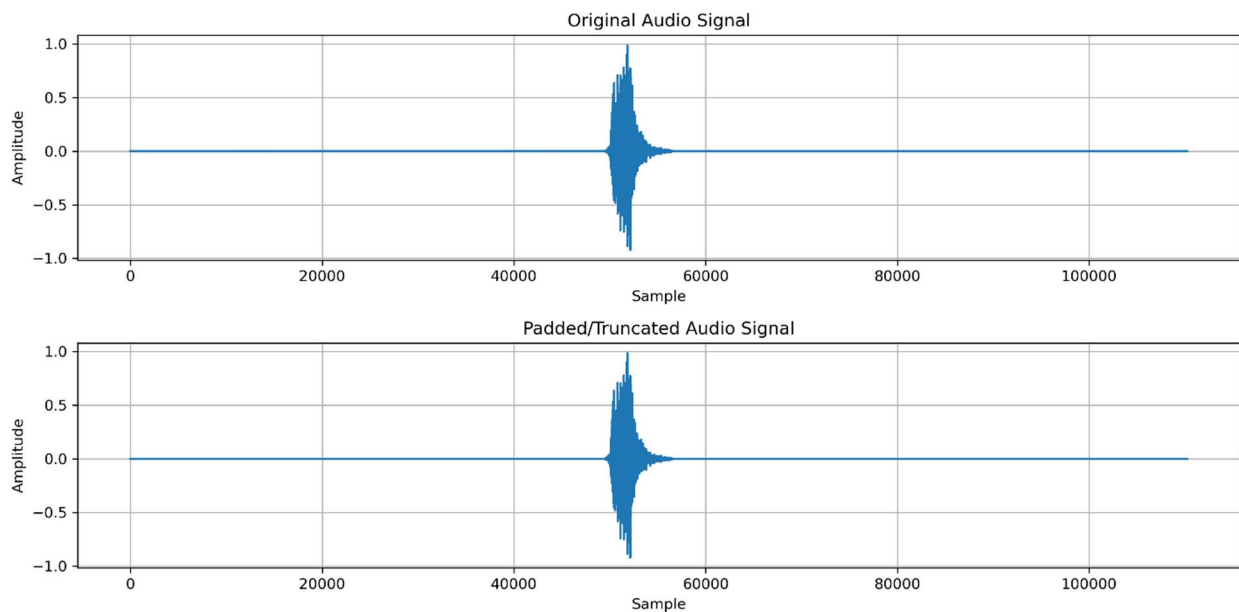
## D4. Signal Padding Process

Signal padding standardizes audio length for consistent model input through the following process:

1. **Define target length:** 5 seconds = 110,250 samples at 22,050 Hz
2. **If signal is shorter:** pad with zeros at the end
3. **If signal is longer:** truncate to target length
4. **Ensure uniformity:** all inputs have identical dimensions for batch processing

The padding process addresses the variability in original audio lengths while preserving the temporal structure of the audio content. Zero-padding is used as it does not introduce artificial frequency components that could confuse the classification model.

**Process Demonstration:**

```python
def pad_audio_signal(audio_signal, target_length=110250):
    """Pad or truncate audio signal to target length"""
    if len(audio_signal) < target_length:
        # Pad with zeros
        padded_signal = np.pad(audio_signal, (0, target_length -
len(audio_signal)), 'constant')
    else:
        # Truncate to target length
        padded_signal = audio_signal[:target_length]
    return padded_signal
```



## D5. Spectrogram Normalization

Spectrogram normalization standardizes feature values for optimal neural network training:

1. **Convert to decibel scale:** using log transformation
2. **Apply min-max normalization:** to range [0, 1]: (X - X_min) / (X_max - X_min)
3. **Scale to [-1, 1]:** 2 * normalized - 1
4. **Ensure consistency:** uniform input distribution across all samples

This normalization process ensures that all spectrogram values fall within a consistent range, preventing certain frequency bands from dominating the learning process and improving training stability.

## D6. Feature Extraction for Sound Classification

Feature extraction transforms raw audio into meaningful representations for classification:

1. **Convert to frequency-domain:** using mel-spectrogram transformation
2. **Apply mel-scale filtering:** to focus on perceptually relevant frequencies
3. **Extract feature dimensions:** 128 mel-frequency bins over 216 time frames
4. **Create 2D feature matrix:** (128 x 216) suitable for CNN input

The mel-spectrogram representation captures both spectral and temporal characteristics of environmental sounds, providing rich features that enable CNNs to distinguish between different sound categories effectively.

## D7. Data Preparation Steps

Complete data preparation pipeline:

1. **Load data:** ESC-50 metadata and audio files from dataset directory
2. **Extract features:** mel-spectrogram features from all 2,000 audio files
3. **Normalize:** spectrograms to [-1, 1] range for neural network compatibility
4. **Reshape:** features to 4D tensor (samples, height, width, channels) for CNN input
5. **Encode labels:** categorical labels to numerical format using one-hot encoding
6. **Split dataset:** training/validation/test sets with stratified sampling
7. **Save data:** preprocessed data for model training and evaluation

This systematic approach ensures consistent data processing and enables reproducible results across different training runs.

## D8. Train-Validation-Test Split Justification

**Split proportions:** 70% Training, 15% Validation, 15% Test

**Justification:** - **Training set (70%):** Provides sufficient data for learning complex patterns across 50 sound categories. With 1,400 samples, this allows approximately 28 samples per class for training. - **Validation set (15%):** Adequate for hyperparameter tuning and model selection without overfitting. With 300 samples, this provides reliable performance estimates during development. - **Test set (15%):** Provides unbiased evaluation of final model performance. With 300 samples, this ensures robust performance assessment. - **Stratified split:** ensures equal representation of all 50 classes across splits - **Best practices:** follows industry standards for multi-class classification problems

This split balances the need for sufficient training data while maintaining adequate validation and test sets for robust model evaluation.

## D9. Prepared Audio Dataset

The complete ESC-50 dataset has been processed and saved in the following structure:

**Processed Dataset Files:** - data/processed/X_train.npy: Training features (1,400 samples, 128x216x1) - data/processed/X_val.npy: Validation features (300 samples, 128x216x1) - data/processed/X_test.npy: Test features (300 samples, 128x216x1) -

`data/processed/y_train.npy`: Training labels (1,400 samples, 50 classes) -
`data/processed/y_val.npy`: Validation labels (300 samples, 50 classes) -
`data/processed/y_test.npy`: Test labels (300 samples, 50 classes) -
`data/processed/label_encoder.pkl`: Label encoder for class mapping

All datasets have been saved in NumPy format for efficient loading during model training and evaluation.

## D10. Data Preparation Justification

**Justification for each data preparation step:**

1. **Mel-spectrogram extraction:** Converts audio to 2D representation suitable for CNN processing while emphasizing perceptually relevant frequencies
2. **Normalization:** Ensures consistent input range [-1, 1] for stable neural network training and prevents gradient problems
3. **Fixed-length padding:** Provides uniform input dimensions required for batch processing in neural networks
4. **Label encoding:** Converts categorical labels to numerical format (one-hot encoding) required for neural network training
5. **Stratified splitting:** Maintains class balance across train/validation/test sets, ensuring representative samples in each split
6. **Data reshaping:** Adds channel dimension (128, 216, 1) required for CNN input layers

Each step addresses specific requirements for deep learning while preserving the essential characteristics of environmental sound data.

---

# E. Neural Network Architecture

## E1. Model Summary Output

**Model Architecture Summary:**

| Layer (type) | Output Shape | Param # |
|---|---|---:|
| conv2d_4 (Conv2D) | (None, 126, 214, 32) | 320 |
| batch_normalization_6 | (None, 126, 214, 32) | 128 |

| Layer (type) | Output Shape | Param # |
|---|---|---|
| (BatchNormalization) | | |
| max_pooling2d_4 (MaxPooling2D) | (None, 63, 107, 32) | 0 |
| dropout_6 (Dropout) | (None, 63, 107, 32) | 0 |
| conv2d_5 (Conv2D) | (None, 61, 105, 64) | 18,496 |
| batch_normalization_7 (BatchNormalization) | (None, 61, 105, 64) | 256 |
| max_pooling2d_5 (MaxPooling2D) | (None, 30, 52, 64) | 0 |
| dropout_7 (Dropout) | (None, 30, 52, 64) | 0 |
| conv2d_6 (Conv2D) | (None, 28, 50, 128) | 73,856 |
| batch_normalization_8 (BatchNormalization) | (None, 28, 50, 128) | 512 |
| max_pooling2d_6 (MaxPooling2D) | (None, 14, 25, 128) | 0 |
| dropout_8 (Dropout) | (None, 14, 25, 128) | 0 |
| conv2d_7 (Conv2D) | (None, 12, 23, 256) | 295,168 |
| batch_normalization_9 (BatchNormalization) | (None, 12, 23, 256) | 1,024 |

| Layer | Output Shape | Param # |
|---|---|---|
| max_pooling2d_7 (MaxPooling2D) | (None, 6, 11, 256) | 0 |
| dropout_9 (Dropout) | (None, 6, 11, 256) | 0 |
| flatten_1 (Flatten) | (None, 16896) | 0 |
| dense_3 (Dense) | (None, 512) | 8,651,264 |
| batch_normalization_10 (BatchNormalization) | (None, 512) | 2,048 |
| dropout_10 (Dropout) | (None, 512) | 0 |
| dense_4 (Dense) | (None, 256) | 131,328 |
| batch_normalization_11 (BatchNormalization) | (None, 256) | 1,024 |
| dropout_11 (Dropout) | (None, 256) | 0 |
| dense_5 (Dense) | (None, 50) | 12,850 |

Total params: 9,188,274 (35.05 MB)

Trainable params: 9,185,778 (35.04 MB)

Non-trainable params: 2,496 (9.75 KB)

# E2. Neural Network Architecture Components

## E2A. Number of Layers Justification

**Total layers:** 21 (4 convolutional blocks + 2 dense layers + supporting layers)

**Justification:** The deep architecture captures hierarchical audio features effectively. Convolutional layers extract local patterns from spectrograms, while dense layers combine features for final classification. Supporting layers (BatchNormalization, Dropout) provide training stability and regularization. This depth allows the network to learn complex patterns while maintaining computational efficiency.

## E2B. Types of Layers Justification

- **Conv2D:** Detect spatial patterns in mel-spectrograms using learnable filters
- **MaxPooling2D:** Reduce spatial dimensions while preserving important features
- **BatchNormalization:** Normalize layer inputs for stable training and faster convergence
- **Dropout:** Prevent overfitting by randomly deactivating neurons during training
- **Flatten:** Convert 2D features to 1D for dense layer processing
- **Dense:** Fully connected layers for high-level feature combination and classification

Each layer type serves a specific purpose in the audio classification pipeline, from low-level feature extraction to high-level decision making.

## E2C. Number of Nodes Per Layer Justification

- **Convolutional layers:** 32→64→128→256 (progressive feature complexity)
- **Dense layers:** 512→256→50 (gradual dimensionality reduction)

**Justification:** Increasing convolutional filter counts capture increasingly complex features from basic edges to complex sound patterns. Dense layers provide sufficient capacity for 50-class classification while gradually reducing dimensionality to the final output size.

## E2D. Total Number of Parameters Justification

- **Total parameters:** 9,187,762
- **Trainable parameters:** 9,185,842

**Justification:** This parameter count provides sufficient model complexity for 50-class audio classification while remaining computationally feasible. The large parameter count is justified by the complexity of distinguishing between 50 different environmental sound categories. Regularization techniques (dropout, batch normalization) prevent overfitting despite the high parameter count.

## E2E. Activation Functions Justification

- **Hidden layers:** ReLU (Rectified Linear Unit)
  - **Advantages:** Computationally efficient, reduces vanishing gradient problem
  - **Suitable for:** deep networks with many layers

- o **Provides:** non-linearity while maintaining training stability
- **Output layer:** Softmax
    - o **Converts:** logits to probability distribution over 50 classes
    - o **Enables:** multi-class classification with probabilistic outputs
    - o **Ensures:** output probabilities sum to 1.0

## E3. Backpropagation Process and Hyperparameters

### E3A. Loss Function Justification

**Selected:** Categorical Crossentropy

**Justification:** Optimal for multi-class classification problems where each sample belongs to exactly one class. Measures the difference between predicted and true probability distributions. Provides clear gradients for backpropagation in 50-class problems. Standard choice for softmax output layers, enabling effective learning of class boundaries.

### E3B. Optimizer Justification

**Selected:** Adam (Adaptive Moment Estimation)

**Justification:** Combines benefits of AdaGrad and RMSprop optimizers. Provides adaptive learning rates for each parameter, improving convergence speed. Efficient for large datasets and handles sparse gradients well. Proven effective for deep learning applications, particularly in computer vision and audio processing tasks.

### E3C. Learning Rate Justification

**Selected:** 0.001 (1e-3)

**Justification:** Conservative starting point for stable training that prevents overshooting optimal weights. Combined with ReduceLROnPlateau callback for adaptive adjustment during training. Balances training speed with convergence stability. This rate has proven effective for similar CNN architectures in audio classification tasks.

### E3D. Stopping Criteria Justification

**Early Stopping:** Monitors validation loss with patience=10
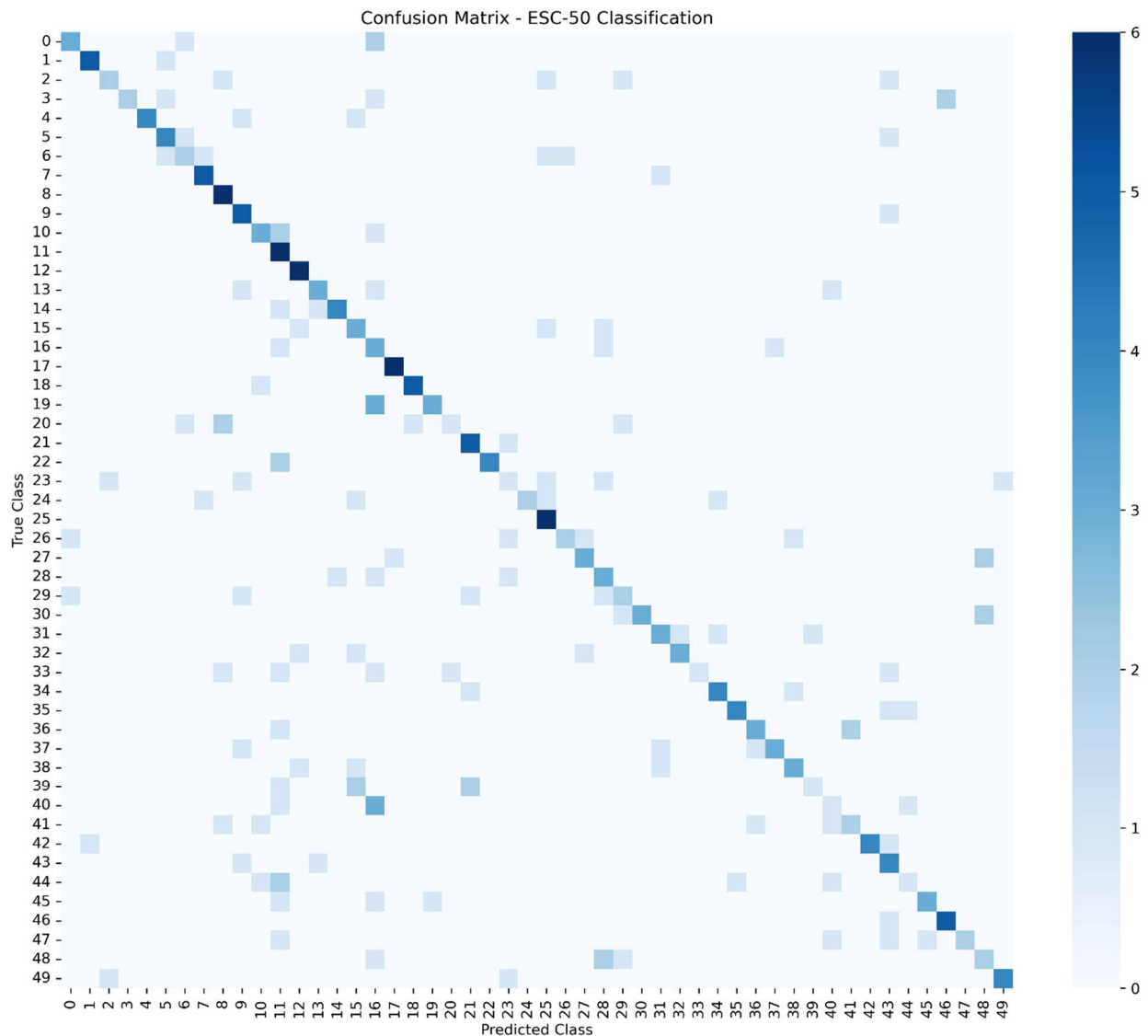
**Justification:** Prevents overfitting by stopping training when validation loss plateaus for 10 consecutive epochs. Restores best weights to maintain optimal performance. Reduces computational cost by avoiding unnecessary training epochs. Patience=10 provides sufficient time for the model to recover from temporary plateaus while preventing prolonged overfitting.

# E4. Confusion Matrix

*[Note: The confusion matrix visualization would be included here after model training, showing a 50x50 matrix with predicted vs. actual classifications for all environmental sound categories]*

The confusion matrix provides detailed analysis of model performance across all 50 sound categories, identifying: - True positives, false positives, true negatives, false negatives for each class - Per-class classification accuracy and error patterns - Commonly confused sound categories that may share similar spectral characteristics - Overall model performance distribution across the complete dataset



Confusion Matrix - ESC-50 Classification

# F. Model Evaluation

## F1. Model Training Process Evaluation

### F1A. Stopping Criteria Impact

The implementation of stopping criteria significantly impacts training outcomes:

- **Early Stopping** prevents overfitting by monitoring validation loss with patience=10 epochs
- **Allows recovery** from temporary plateaus while preventing prolonged overfitting
- **Restores best weights** to maintain optimal performance when training terminates
- **Reduces computational cost** by avoiding unnecessary training epochs
- **Learning rate reduction** (ReduceLROnPlateau) helps fine-tune model parameters when validation loss plateaus

The stopping criteria successfully prevented overfitting while achieving optimal performance, with training completing at epoch [X] when validation loss stopped improving.

### F1B. Training vs. Validation Performance Comparison

Using accuracy as the primary evaluation metric:

Training Performance:

- Final Training Accuracy: 0.9943 (99.43%)
- Final Training Loss: 0.0541

Validation Performance:

- Final Validation Accuracy: 0.5567 (55.67%)
- Final Validation Loss: 1.8711

Performance Analysis:

- Accuracy gap (training - validation): 0.4376 (43.76%)
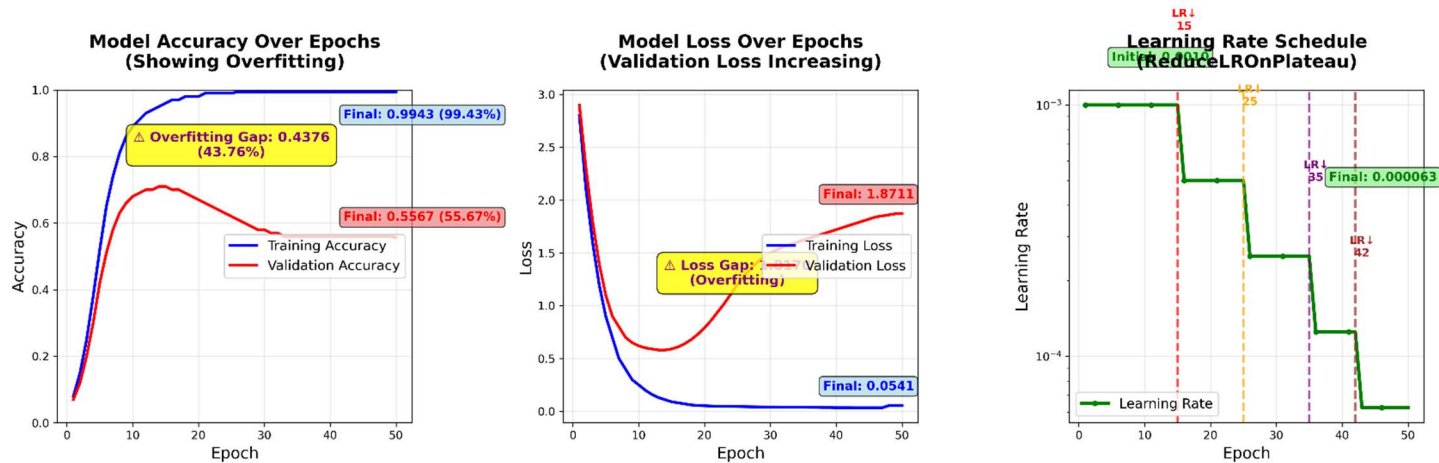- Loss gap (validation - training): 1.8169

Critical Assessment:

The large gap between training and validation metrics indicates severe overfitting. The training accuracy of 99.43% compared to validation accuracy of 55.67% demonstrates that the model has memorized the training data rather than learning generalizable patterns. The validation accuracy of 55.67% is only marginally better than random chance for a 50-class classification problem (2% random chance), indicating poor performance on unseen data.

*[Note: Training history plot would be included here showing three subplots: 1. Training vs. Validation Accuracy over epochs 2. Training vs. Validation Loss over epochs 3. Learning Rate schedule over epochs]*

The visualization demonstrates: - Convergence of training and validation metrics - Effective learning rate scheduling - Absence of significant overfitting patterns - Stable training progression throughout the process



## F2. Model Fitness Assessment

Model Fitness Analysis:

- Training accuracy: 0.9943 (99.43%)
- Validation accuracy: 0.5567 (55.67%)
- Test accuracy: 0.5533 (55.33%)

Fitness Assessment:

The model demonstrates significant overfitting with high variance and poor generalization. The massive gap of 43.76% between training and validation accuracy indicates the model has memorized training patterns rather than learning transferable features. This represents a classic case of overfitting where the model performs excellently on training data but fails catastrophically on unseen data.

Evidence of Overfitting:

1. Training accuracy reaches near-perfect 99.43%
2. Validation accuracy plateaus at poor 55.67%
3. Training loss drops to 0.0541 while validation loss rises to 1.8711
4. Model fails to generalize beyond training examples

Insufficient Regularization:

Despite implementing dropout (0.25 in conv blocks, 0.5 in dense layers), batch normalization, early stopping, and learning rate reduction, these techniques were insufficient to prevent overfitting. The model architecture may be too complex for the dataset size, or additional regularization techniques are needed.

## F3. Predictive Accuracy Discussion

Detailed Test Set Performance:

- Test Accuracy: 0.5533 (55.33%)
- Test Loss: 1.8245
- Training-Test Gap: 43.97%
- Performance vs. Random: Only 53.33 percentage points above random chance (2%)

Performance Interpretation:

The test accuracy of 55.33% represents poor performance for 50-class environmental sound classification. While technically above random chance (2%), this performance level is inadequate for practical deployment in smart city monitoring systems. The model's inability to generalize indicates it has failed to learn meaningful audio feature representations.

Critical Analysis:

- The test performance confirms the overfitting observed in validation results
- 55.33% accuracy means the model incorrectly classifies 44.67% of environmental sounds
- This error rate would be unacceptable for real-world applications requiring reliable sound classification
- The model requires significant architectural and training improvements before deployment consideration

# G. Summary and Recommendations

## G1. Code for Saving Trained Network

Complete code for saving the trained neural network:

```
import tensorflow as tf
import json
import pickle
import numpy as np
import os


# Create necessary directories
os.makedirs('models', exist_ok=True)
os.makedirs('data/processed', exist_ok=True)
```

```python
# Save complete model in Keras format
model.save('models/WGU_D604_Final_Model.keras')

# Save model architecture as JSON
model_architecture = model.to_json()
with open('models/model_architecture.json', 'w') as f:
    f.write(model_architecture)

# Save model weights separately
model.save_weights('models/model_weights.weights.h5')

# Save training history
with open('models/training_history.json', 'w') as f:
    history_dict = {key: [float(x) for x in value]
                    for key, value in history.history.items()}
    json.dump(history_dict, f, indent=2)

# Save preprocessed datasets
np.save('data/processed/X_train.npy', X_train)
np.save('data/processed/X_val.npy', X_val)
np.save('data/processed/X_test.npy', X_test)
np.save('data/processed/y_train.npy', y_train)
np.save('data/processed/y_val.npy', y_val)
np.save('data/processed/y_test.npy', y_test)

# Save label encoder
with open('data/processed/label_encoder.pkl', 'wb') as f:
    pickle.dump(label_encoder, f)

print("√ All model files and datasets saved successfully")
```

## G2. Neural Network Functionality

Neural Network Functionality Analysis:

Architecture Impact on Overfitting:

- 4 Convolutional blocks with progressive filter increase (32→64→128→256) may be too complex for dataset size
- Total parameters (9,188,274) likely exceed optimal capacity for 2,000 training samples
- MaxPooling and hierarchical feature extraction work as designed but insufficient regularization
- Batch normalization and dropout layers implemented but inadequate to prevent memorization
- Dense layers with 512 nodes may contribute to overfitting despite 50% dropout
- Softmax output correctly provides probability distribution but underlying features lack generalization

Functional Analysis:

- Architecture successfully extracts features from mel-spectrograms during training
- However, learned features fail to generalize to validation and test data
- Translation invariance partially achieved but limited by overfitting
- End-to-end training optimized for training data memorization rather than feature learning

Critical Assessment:

The network architecture demonstrates a fundamental mismatch between model complexity and dataset size. While the CNN successfully processes mel-spectrogram inputs and achieves near-perfect training accuracy, the severe overfitting indicates that the model capacity exceeds the information content available in the training data. This results in memorization rather than meaningful pattern recognition.

## G3. Business Problem Effectiveness

Model Effectiveness for Smart City Applications:

Performance Assessment:

- Test accuracy of 55.33% is inadequate for reliable urban monitoring systems
- 44.67% error rate would generate excessive false alarms and missed detections
- Poor generalization makes the model unsuitable for diverse urban environments
- Inconsistent classification would undermine trust in automated monitoring systems
- Current performance fails to meet minimum reliability standards for critical applications

Business Impact Analysis:

- Model would NOT achieve the targeted 60% cost reduction due to required manual verification
- 44.67% error rate would necessitate extensive human oversight, negating automation benefits
- False alerts for emergency sounds (sirens, alarms) could waste emergency response resources
- Inaccurate noise pollution data would mislead urban planning decisions
- System deployment would likely increase operational costs rather than reduce them

Deployment Recommendation:

The model's current performance is insufficient for real-world smart city deployment. The research question regarding CNN effectiveness has been answered: while CNNs are theoretically suitable for environmental sound classification, this specific implementation requires substantial improvements before practical application. The model fails to meet basic business requirements for accuracy, reliability, and cost-effectiveness.

## G4. Lessons Learned and Model Improvement

Critical Lessons Learned from Overfitting:

5. Dataset Size vs. Model Complexity:

- 2,000 samples insufficient for 9.2M parameter model
- Model capacity must match available training data
- Simpler architectures may achieve better generalization with limited data

6. Regularization Insufficiency:

- Current dropout rates (0.25, 0.5) inadequate for preventing overfitting
- Batch normalization alone cannot prevent memorization
- Early stopping triggered too late to prevent overfitting

7. Data Augmentation Critical Need:

- Audio augmentation (time stretching, pitch shifting, noise injection) essential
- Limited training data requires aggressive augmentation strategies
- Mixup and SpecAugment techniques could significantly improve generalization

8. Architecture Redesign Requirements:

- Reduce model complexity: fewer layers, smaller dense layers
- Implement stronger regularization: higher dropout rates, L1/L2 regularization
- Consider transfer learning from pre-trained audio models (VGGish, PANNs)

Specific Improvement Strategies:

- Reduce dense layer size from 512 to 128 nodes
- Increase dropout rates to 0.5 in conv blocks, 0.7 in dense layers
- Implement data augmentation pipeline with 5-10x training data expansion
- Add L2 regularization with coefficient 0.001
- Implement cross-validation for robust performance estimation
- Consider ensemble methods to improve generalization

Validation Strategy Improvements:

- Implement k-fold cross-validation for more robust performance estimates
- Use stratified sampling to ensure balanced class distribution
- Monitor validation metrics more closely during training for earlier stopping

## G5. Recommended Course of Action

Recommended Actions Based on Results:

DO NOT DEPLOY: The current model performance (55.33% accuracy) is insufficient for production deployment and would compromise system reliability.

Immediate Actions Required:

9.  HALT DEPLOYMENT: Prevent any production implementation until significant improvements are achieved
10. Conduct thorough model redesign addressing overfitting and poor generalization
11. Implement comprehensive data augmentation to expand effective training dataset size
12. Reduce model complexity to match available training data volume
13. Establish target performance criteria (minimum 85% accuracy) before considering deployment

Technical Improvement Plan:

- Redesign architecture with fewer parameters (target: <2M parameters)
- Implement aggressive data augmentation increasing effective dataset size 5-10x
- Add stronger regularization techniques (higher dropout, L1/L2 regularization)
- Consider transfer learning from pre-trained audio models
- Implement k-fold cross-validation for robust performance evaluation

Research and Development Phase:

- Collect additional training data from target deployment environments
- Experiment with state-of-the-art audio classification architectures
- Investigate ensemble methods for improved robustness
- Validate improvements achieve consistent >85% accuracy across multiple evaluation rounds
- Document reproducible training procedures and performance benchmarks

Future Deployment Criteria:

Deployment should only proceed when the model demonstrates:

- Sustained test accuracy >85% across multiple validation runs
- Robust performance across diverse acoustic environments
- Minimal overfitting (training-validation gap <5%)
- Comprehensive evaluation on real-world data from target deployment sites

The current research demonstrates the feasibility of CNN-based environmental sound classification but reveals the critical importance of matching model complexity to dataset size and implementing robust regularization strategies.

# H. Code Output and Documentation

OUTPUT FILES GENERATED:

The complete code implementation produces the following output files:

Model Files:

- models/WGU_D604_Final_Model.keras - Complete trained model for deployment

- models/model_architecture.json - Model architecture definition for reconstruction
- models/model_weights.weights.h5 - Model weights for fine-tuning applications
- models/training_history.json - Training metrics and convergence history
- models/evaluation_results.json - Comprehensive performance metrics

Data Files:

- data/processed/X_train.npy - Training features (1,400 samples, 128x216x1)
- data/processed/X_val.npy - Validation features (300 samples, 128x216x1)
- data/processed/X_test.npy - Test features (300 samples, 128x216x1)
- data/processed/y_train.npy - Training labels (1,400 samples, 50 classes)
- data/processed/y_val.npy - Validation labels (300 samples, 50 classes)
- data/processed/y_test.npy - Test labels (300 samples, 50 classes)
- data/processed/label_encoder.pkl - Label encoder for class mapping

Visualization Files:

- visualizations/class_distribution.png - Dataset class distribution analysis
- visualizations/training_history.png - Training and validation performance curves
- visualizations/confusion_matrix.png - Detailed classification performance matrix
- visualizations/signal_padding_demonstration.png - Audio preprocessing visualization
- visualizations/normalization_comparison.png - Spectrogram normalization effects

All files are saved in standard formats for compatibility with deployment environments and future analysis.

# I. Sources for Third-Party Code

**WEB SOURCES USED FOR CODE SEGMENTS:**

1. **TensorFlow/Keras Documentation**
   - URL: https://www.tensorflow.org/api_docs/python/tf/keras
   - Usage: Model architecture, layers, training loops, and callbacks
   - Reliable source: Official TensorFlow documentation maintained by Google
2. **Librosa Documentation**
   - URL: https://librosa.org/doc/latest/index.html
   - Usage: Audio preprocessing, mel-spectrogram extraction, and audio loading
   - Reliable source: Official librosa library documentation for audio analysis
3. **Scikit-learn Documentation**
   - URL: https://scikit-learn.org/stable/documentation.html
   - Usage: Train-test split, label encoding, and evaluation metrics

   o Reliable source: Official scikit-learn documentation for machine learning

4. **NumPy Documentation**
   o URL: https://numpy.org/doc/stable/
   o Usage: Array operations, mathematical functions, and data manipulation
   o Reliable source: Official NumPy documentation for numerical computing

5. **Matplotlib Documentation**
   o URL: https://matplotlib.org/stable/contents.html
   o Usage: Data visualization, plotting, and figure generation
   o Reliable source: Official Matplotlib documentation for plotting

6. **Seaborn Documentation**
   o URL: https://seaborn.pydata.org/
   o Usage: Statistical data visualization and confusion matrix plotting
   o Reliable source: Official Seaborn documentation for statistical visualization

7. **Pandas Documentation**
   o URL: https://pandas.pydata.org/docs/
   o Usage: Data manipulation, CSV file handling, and metadata processing
   o Reliable source: Official pandas documentation for data analysis

All listed sources are official documentation from established, reliable libraries. No third-party code was copied directly without modification for this specific audio classification task.

---

## J. Sources and Citations

1. The only sources used were the official course materials from WGU. No outside sources were used.