

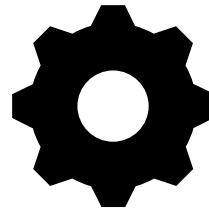
Robustness of a CNN for music genre classification

Erik von Heyden (8720832)

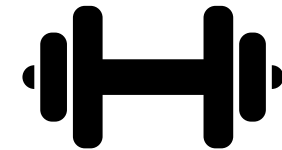
Methodology



1. Data exploration



2. Data preparation



3. Model training



4. Adversarial
Attacks

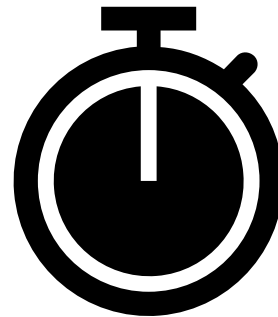


5. GradCAM

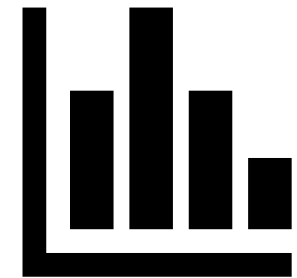
GTZAN Dataset



1000 songs
in **10 genres**
(100 per genre)

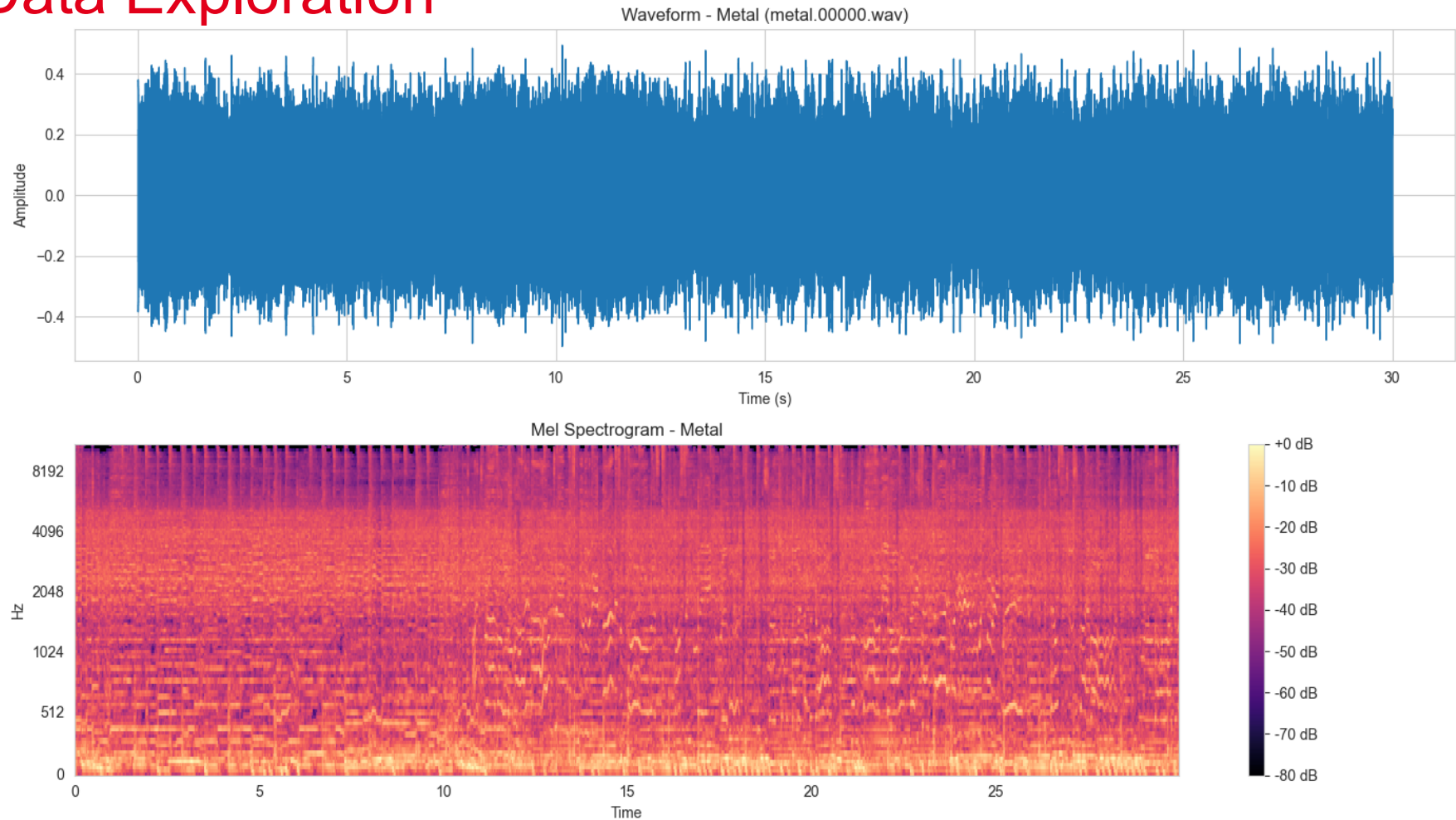


30 sec
as .wav

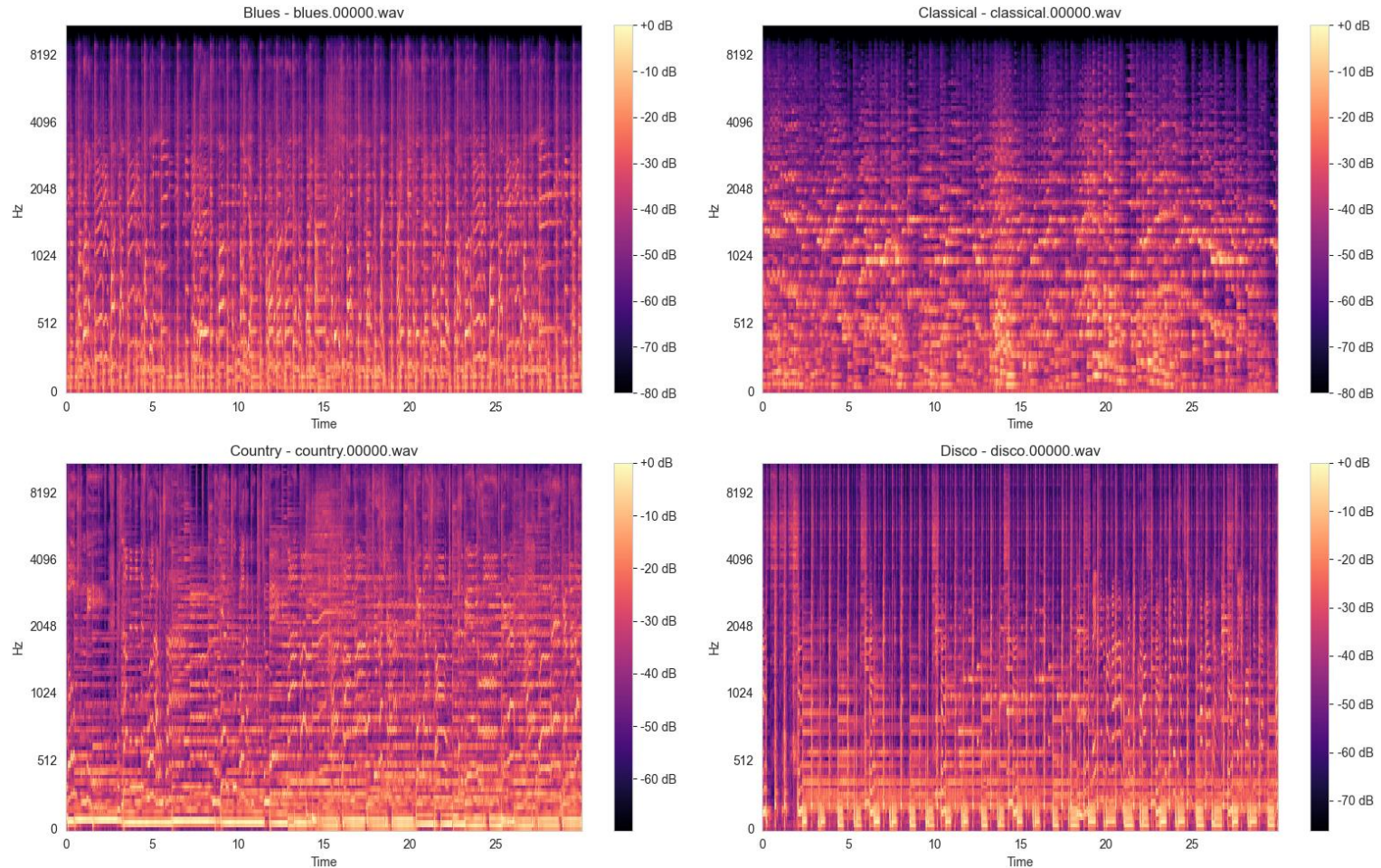


Mel spectrograms,
feature tables

Data Exploration



Data Exploration



Preprocessing

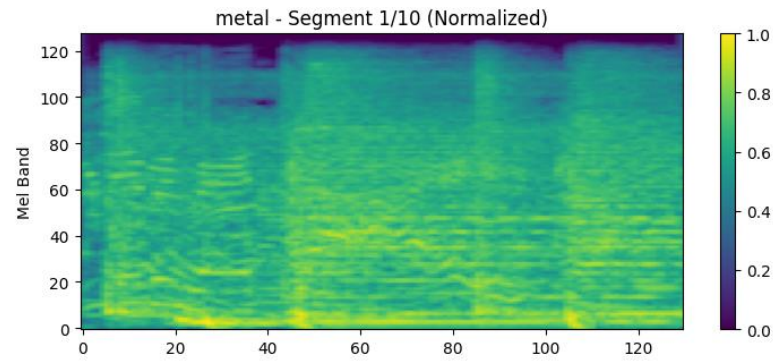
- 70% Training (~700 Songs)
- 15% Validation (~150 Songs)
- 15% Test (~150 Songs)



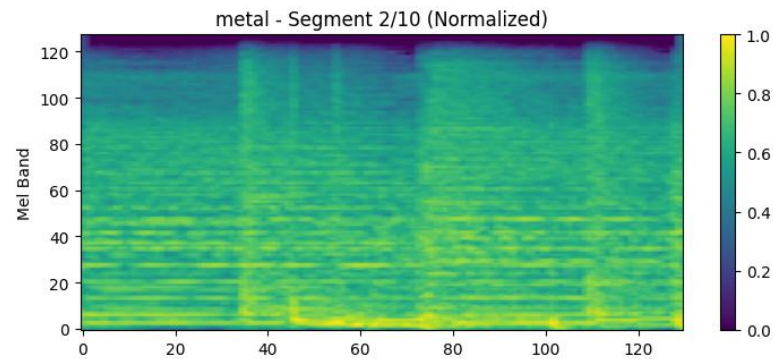
**Split into 10 segments
à 3 sec**

- 70% Training (~**7000** segments)
- 15% Validation (~**1500** segments)
- 15% Test (~**1500** segments)

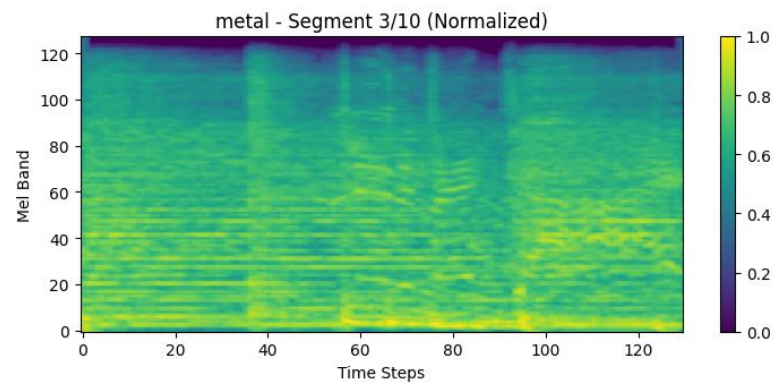




Segment 1
Shape: (128, 130)
Min: 0.000
Max: 1.000
Mean: 0.589



Segment 2
Shape: (128, 130)
Min: 0.000
Max: 1.000
Mean: 0.593

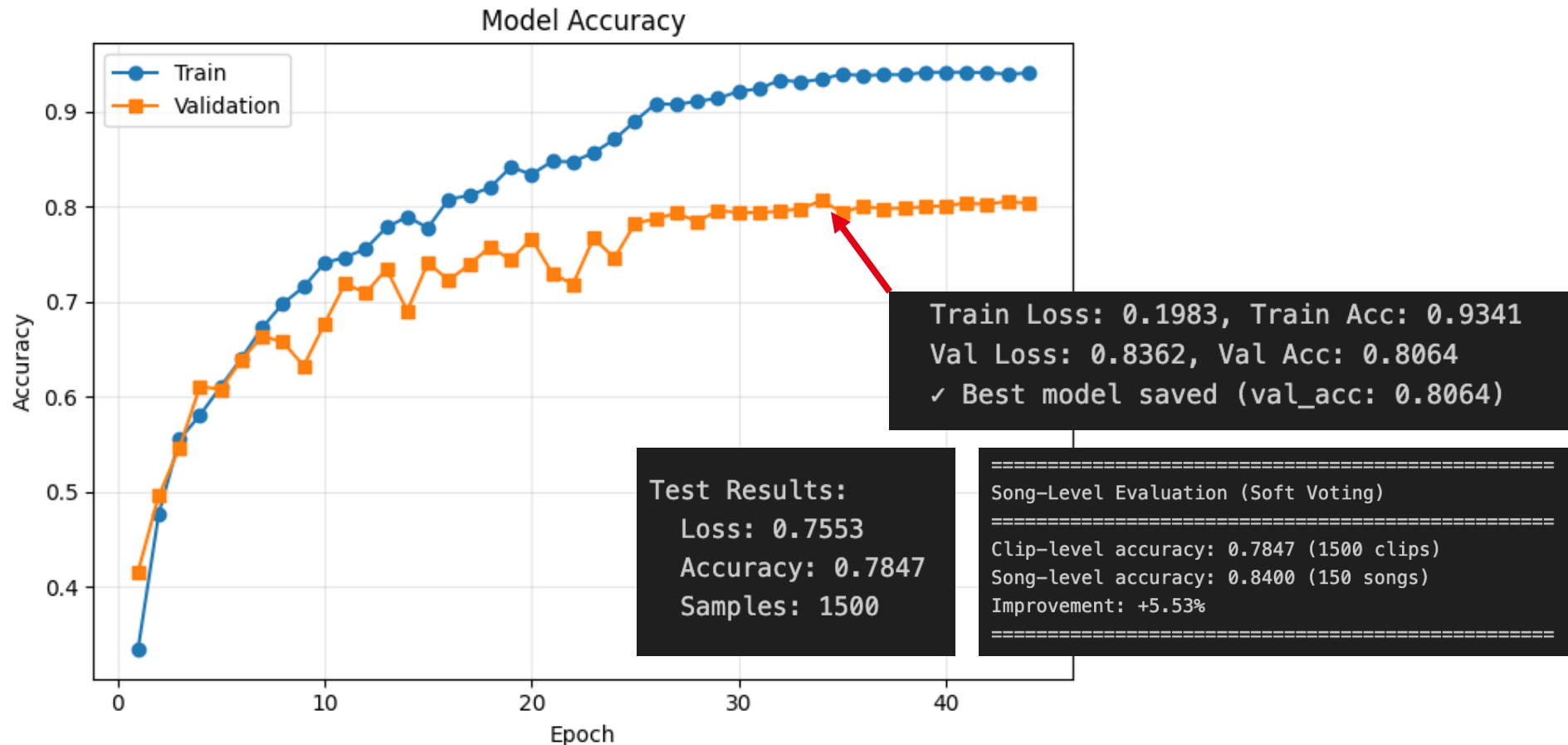


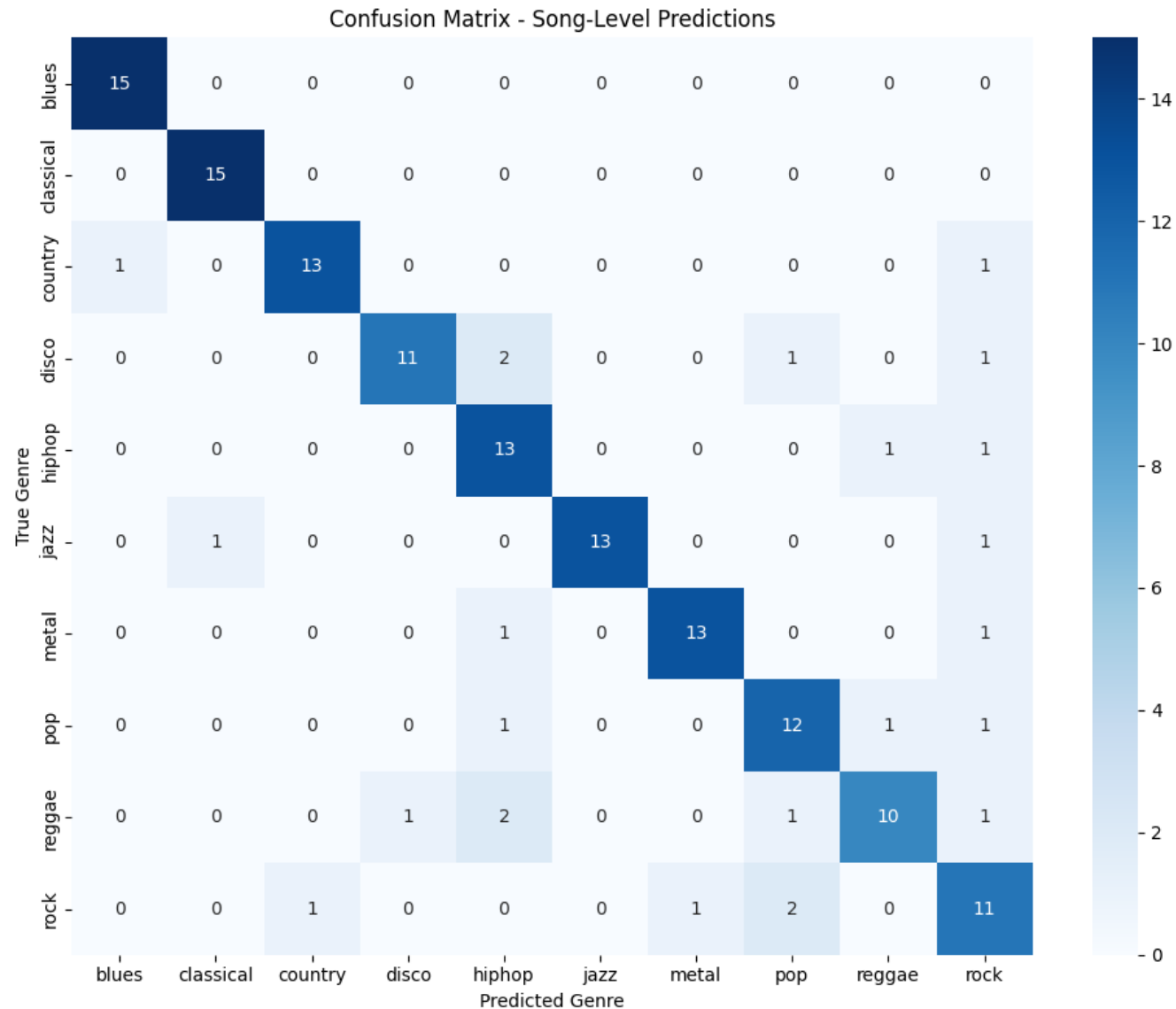
Segment 3
Shape: (128, 130)
Min: 0.000
Max: 1.000
Mean: 0.615

CNN Architecture

- Input: **Mel Spectrogram** $128 \times 130 \times 1$
- Convolutional Blocks: $2x$ (Conv \rightarrow BatchNorm) \rightarrow MaxPool \rightarrow Dropout
- Progressive filters: $32 \rightarrow 64 \rightarrow 128 \rightarrow 256$
- Classification: Flatten \rightarrow Dense(512) \rightarrow **Dense(10)**
- **1,4M** parameters

Training Results





Adversarial Attacks

FGSM (Fast Gradient Sign Method)

- Single Step Attack
- Moves input in direction of gradient

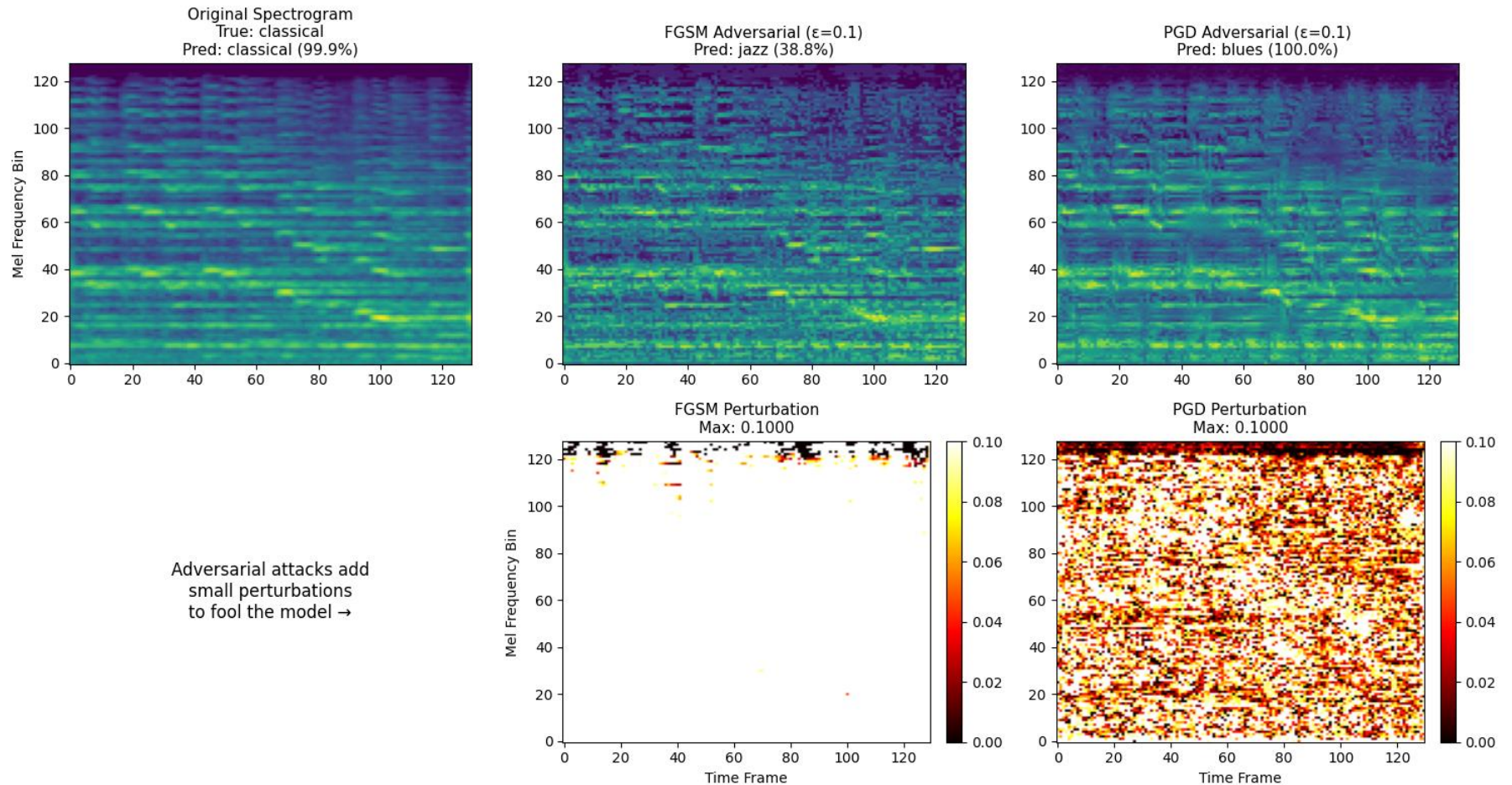
PGD (Projected Gradient Descent)

- Multi-step iterative attack
- Takes small steps and projects back into allowed perturbation range

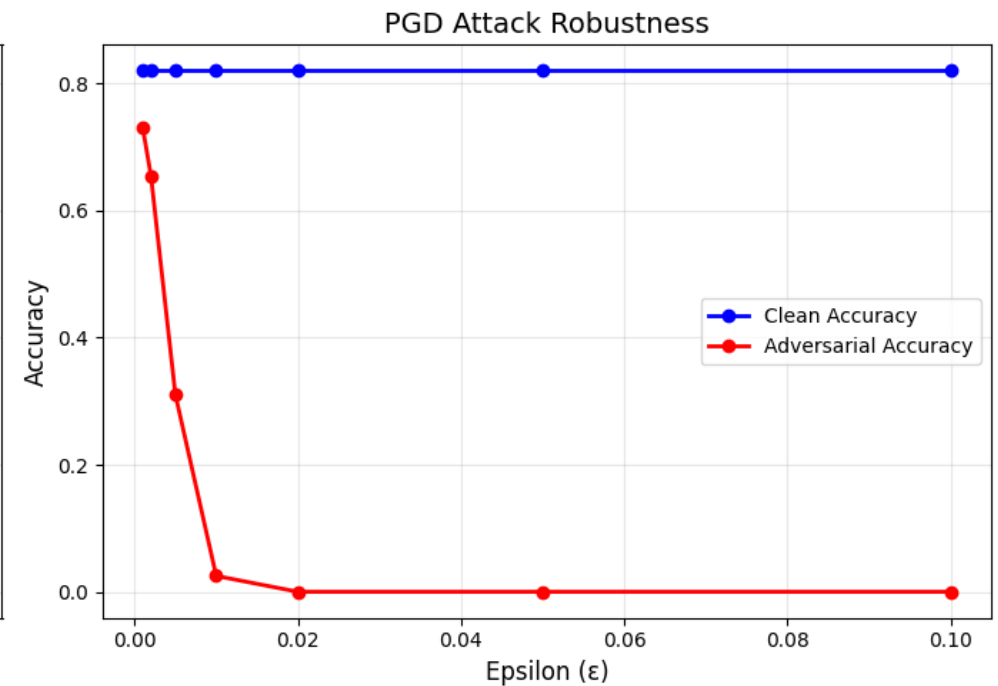
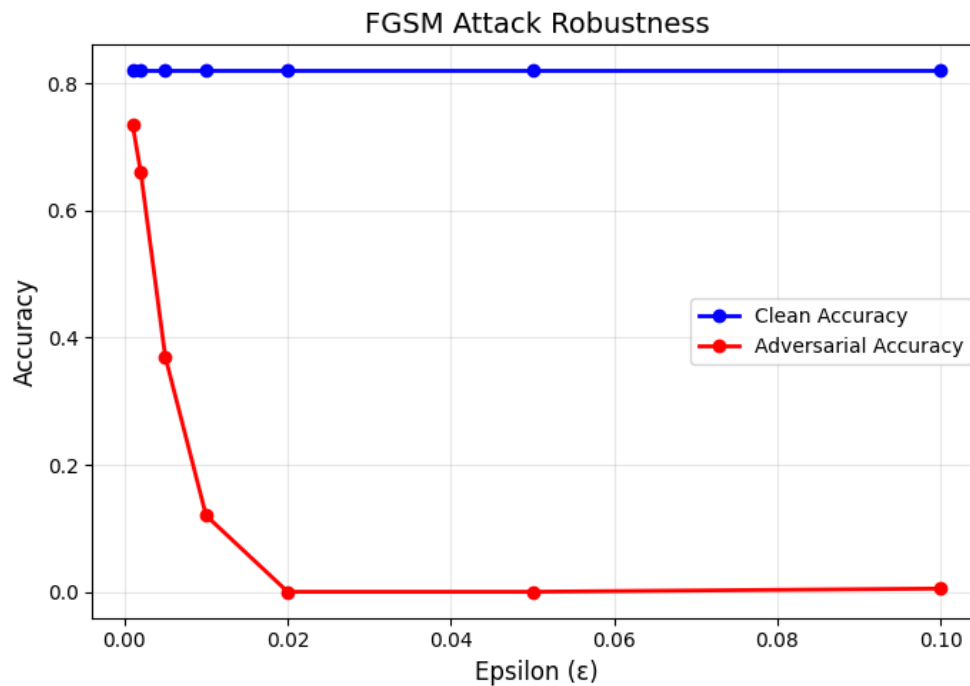


Adversarial
Robustness
Toolbox

Adversarial Attacks



Adversarial Attacks



GradCAM

1. Baseline Analysis (Correct Predictions)

- 20 correctly classified samples (2 per genre)
- **Goal:** Understand what features the model uses normally

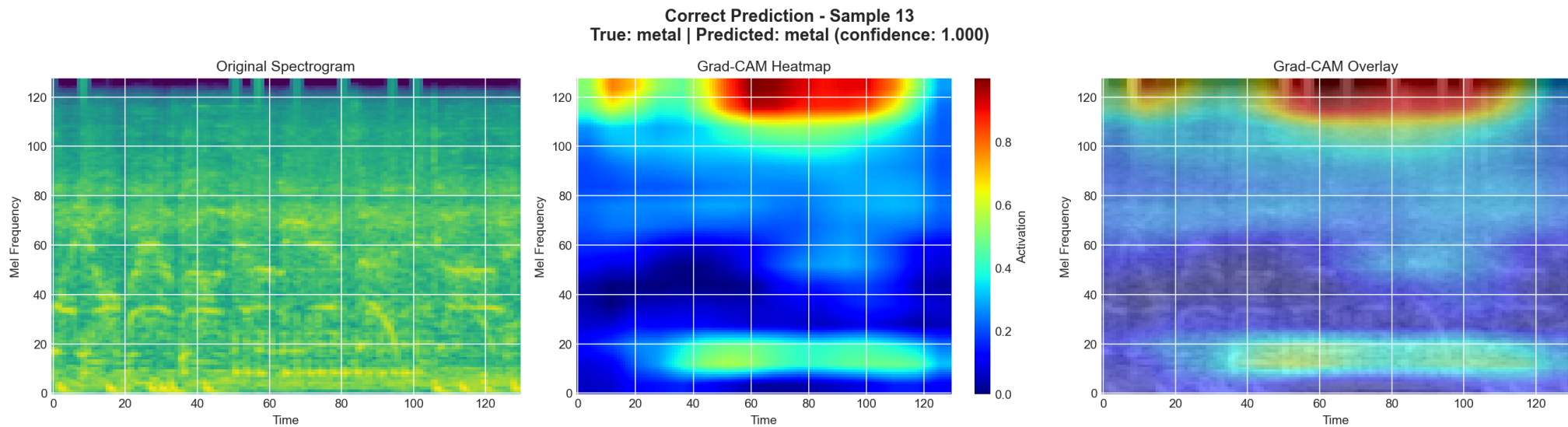
2. Natural Failures (Error Analysis)

- 5 misclassified samples without attacks
- **Goal:** Identify confusion sources

3. Adversarial Analysis (FGSM & PGD, $\epsilon=0.1$)

- 20 samples under FGSM + 20 under PGD attacks
- Side-by-side comparisons: Clean vs. Adversarial
- **Goal:** How do attacks alter model perception?

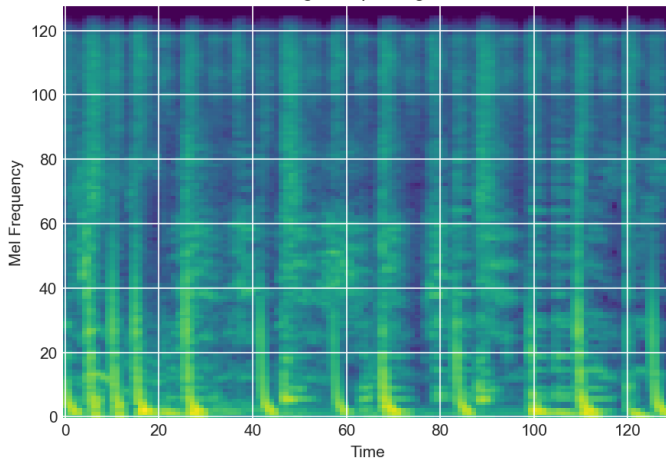
GradCAM – Baseline Analysis



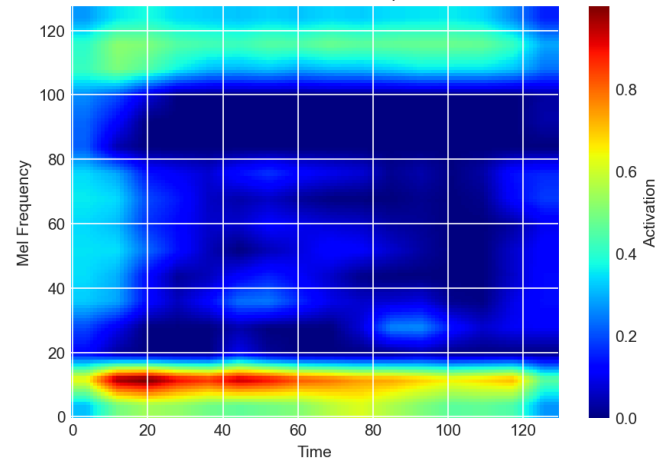
GradCAM – Natural Failures

Misclassification - Sample 3
True: disco | Predicted: hiphop (confidence: 0.996)

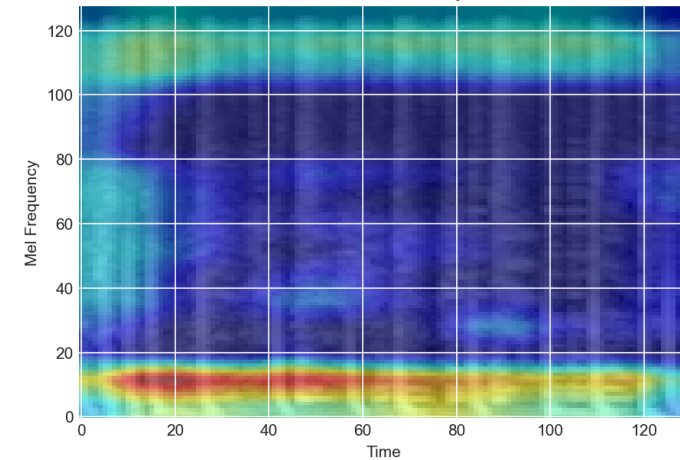
Original Spectrogram



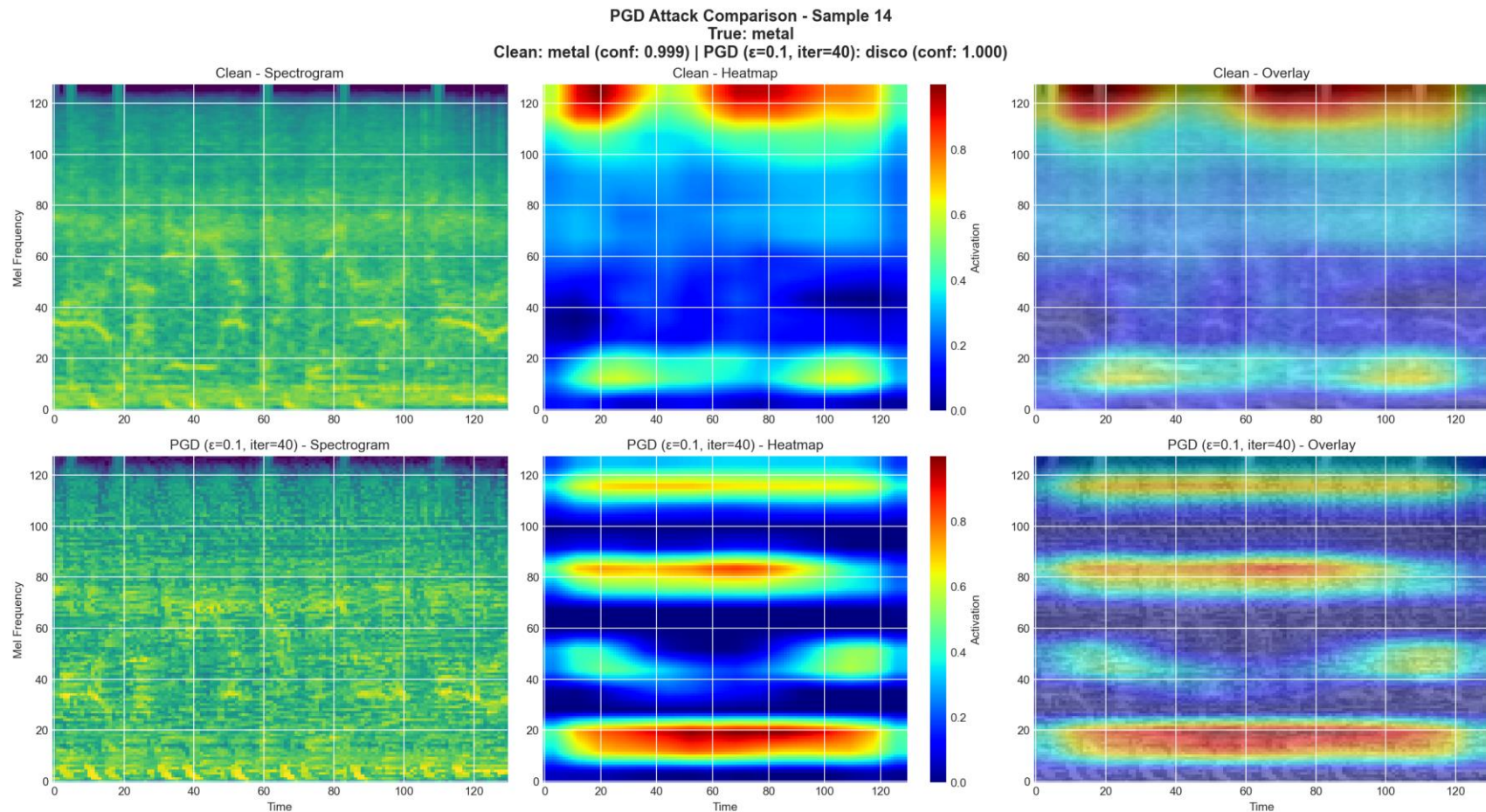
Grad-CAM Heatmap



Grad-CAM Overlay

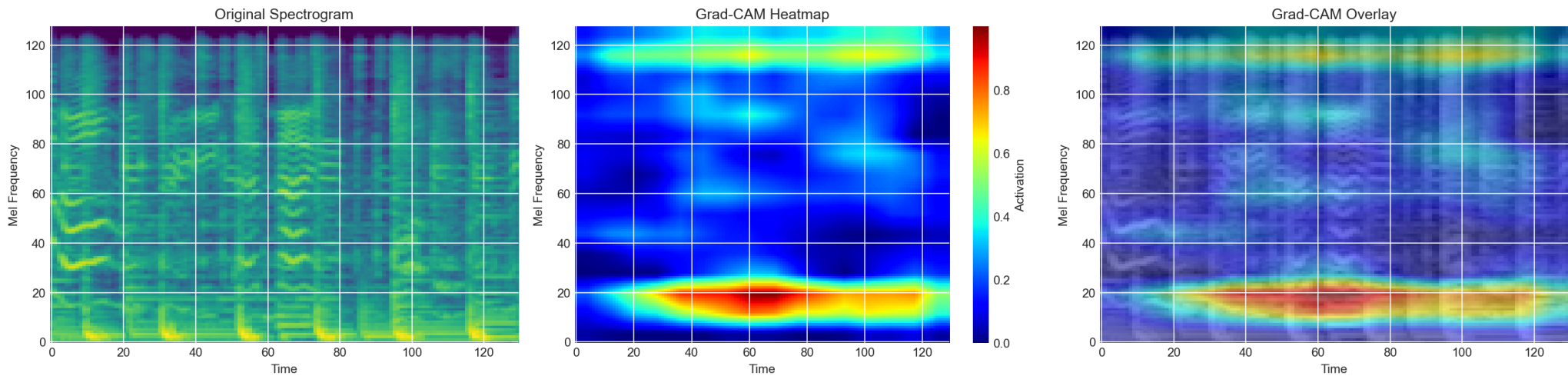


GradCAM – Adversarial Attacks



GradCAM – Adversarial Attacks

Correct Prediction - Sample 8
True: disco | Predicted: disco (confidence: 1.000)



Conclusion

Key Findings	Insight
Model accuracy: 84% song-level	Good baseline performance
Adversarial vulnerability: ~100% at $\epsilon=0.02$	Highly vulnerable without defense mechanisms
Grad-CAM on clean data: Meaningful attention	Model learns musically relevant features
Grad-CAM under attacks: Artificial patterns	Attacks exploit attention mechanisms