# Heart Sound Classification using CNN & Transfer Learning

SANMAY KUMAR JENA-8025340035

*THAPAR UNIVERSITY ,PATIALA, PUNJAB*

## 1. Problem Definition & Scope of the Problem

This project focuses on the automated classification of heart sound recordings into two diagnostic categories: (Normal) and (Abnormal). Heart sounds provide critical information for early detection of cardiac abnormalities such as murmurs, arrhythmias, and structural defects. Traditional auscultation is subjective and requires clinical expertise. Therefore, this project aims to build a reproducible and automated deep-learning-based diagnostic system.

Scope of the project:
- Preprocess raw PCG (Phonocardiogram) audio signals.
- Convert time-domain waveforms into frequency-domain spectrograms using STFT.
- Train a baseline CNN model and compare it with transfer learning models (MobileNetV2 and ResNet50).
- Evaluate using accuracy, confusion matrix, and precision–recall metrics.
- Present a robust methodology capable of generalizing across unseen heart sound samples.

## 2. Literature Review

Heart sound analysis has gained significant attention over the last decade. Classical approaches involve time-domain analysis, frequency-domain filtering, and handcrafted feature extraction methods such as MFCCs, Wavelets, Hilbert transforms, and envelope extraction.

Recent advancements in deep learning have enabled automatic feature extraction using CNNs directly on spectrogram representations. Several research papers demonstrate that converting PCG signals into 2D images leads to superior performance compared to raw waveforms.

Key insights from literature:
- Frequency-domain features outperform raw waveforms.
- CNNs trained on spectrograms provide high diagnostic precision.
- Transfer learning models such as MobileNetV2 and ResNet50 show strong feature extraction capability even with small datasets.
- STFT-based log spectrograms are widely used in biomedical audio classification.
- Stratified dataset splitting improves model reliability and reduces bias

## 3. Model Descriptions (Proposed Architecture / Methodology Adopted)

The methodology adopted in this project follows a complete audio-to-deep-learning pipeline:

1. Data Extraction:
- Dataset provided as a ZIP file containing audio files and multiple annotation CSVs.
- A robust CSV parsing mechanism is implemented to detect filename and label columns.
- Labels are normalized to two classes: Normal (0) and Abnormal (1).

2. Preprocessing:
- Audio files may contain varying sampling rates; therefore, all waveforms are resampled to 2000 Hz.
- Signals are trimmed or padded to a fixed duration of 10 seconds.
- Short-Time Fourier Transform (STFT) is generated using FFT size = 256 and hop length = 128.
- Resulting spectrogram magnitude values are normalized and converted to 128×128 images.
- CNN input shape: (128, 128, 1) for grayscale spectrograms.

3. Dataset Preparation:
- Data is split into Train (70%), Validation (15%), and Test (15%) using stratified sampling.
- For transfer learning models, grayscale images are auto-converted to RGB.
- TensorFlow Dataset API is used with batching, shuffling, and prefetching.

4. Models Used:
A. Custom CNN Model
- Conv2D → MaxPool2D (16 filters)
- Conv2D → MaxPool2D (32 filters)
- Conv2D → GAP (64 filters)
- Dense(128) + Dropout(0.3)
- Output: Dense(2, softmax)
- Optimizer: Adam, Loss: Sparse Categorical Crossentropy

B. MobileNetV2 (Transfer Learning)
- Pretrained on ImageNet, frozen feature extractor.
- Added GAP + Dense(128) + Dropout + Softmax layers.

C. ResNet50 (Transfer Learning)
- Pretrained backbone, frozen layers.
- Custom classification head similar to MobileNetV2.

All models are trained for 6 epochs (recommended 20–50 epochs for final performance).

# 4.*Flow Diagram and Technical Aspects*

Raw Audio → Resampling → STFT Computation → Spectrogram Normalization → Image Resize → Dataset

Pipeline → CNN/Transfer Learning Training → Evaluation & Visualization.

Technical features include:
- Automatic CSV detection using heuristics
- Robust handling of stereo/mono audio
- TF Dataset pipeline for high throughput
- Matplotlib-based metric visualization

# *5. Results & Comparative Analysis*

Dataset Details:
- Maximum samples used: 3000
- Classes: Normal (0), Abnormal (1)
- Input image size: 128×128 grayscale spectrograms

System Configurations:
- Python, TensorFlow, SciPy, NumPy, Matplotlib
- GPU/TPU acceleration recommended

Training Details:
- Batch size: 16
- Optimizer: Adam
- Epochs: 6
- Loss: Sparse Categorical Crossentropy

Model Performance: **CNN**
The project generates the following evaluation outputs:
- Confusion matrix images for each model
- Precision–Recall curves for the Abnormal class
- JSON-based classification reports with precision, recall, F1-score

Precision-Recall (CNN) - Abnormal (AP=0.685)

```
"Normal": {
    "precision": 0.8693467336683417,
    "recall": 0.969187675070028,
    "f1-score": 0.9165562913907285,
    "support": 357.0
},
"Abnormal": {
    "precision": 0.7884615384615384,
    "recall": 0.44086021505376344,
    "f1-score": 0.5655172413793104,
    "support": 93.0
},
"accuracy": 0.86,
"macro avg": {
    "precision": 0.8289041360649401,
    "recall": 0.7050239450618957,
    "f1-score": 0.7410367663850195,
    "support": 450.0
},
"weighted avg": {
    "precision": 0.852630459992269,
    "recall": 0.86,
    "f1-score": 0.8440082210550355,
    "support": 450.0
}
}
```

Model Performance: **RESNET50**





Typical observations:
- CNN shows strong ability to differentiate Normal vs Abnormal patterns.
- Transfer learning models may improve generalization but require RGB conversion overhead.
- Precision–Recall analysis highlights positive-class detection capability.

# 6. Conclusions and Future Scope

The project demonstrates that heart sound classification using spectrogram-based CNN models is effective and computationally efficient. The baseline CNN achieves competitive performance even with a comparatively small dataset.

Future Scope:
- Increase epoch count (20–50) for stronger convergence.

Model Performance: **MOBILENET**

- Apply audio augmentation: noise addition, pitch shift, time-stretch.
- Train end-to-end raw waveform models (1D-CNNs, WaveNet, CRNN).
- Use larger pretrained models such as EfficientNet.
- Deploy as a mobile or web-based diagnostic application.

## 7. References

[1] PhysioNet / CinC Challenge Dataset
[2] Keras API Documentation
[3] Springer & IEEE papers on PCG classification
[4] MobileNetV2 and ResNet50 original research publications