

Rainforest Species Identification

Ben Rapp, Phil Genovese, Jared Huzar, Craig Fox

Problem

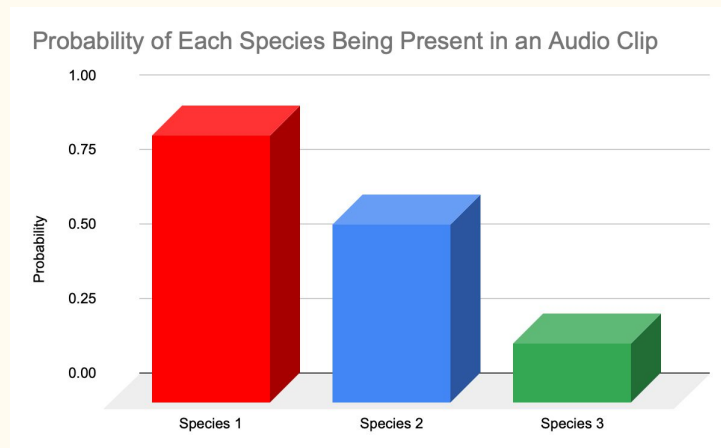
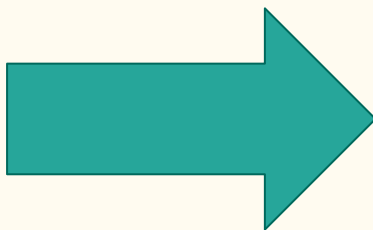
- Identifying species presence is a good proxy for rainforest health
- It is easier to collect meaningful audio data instead of video
- Expert identifiers are expensive and slow
- Basic machine learning models require huge amounts of data especially as the number of potential species increase



Puerto Rican Tody - eBird

Objective

Develop a machine learning model which can accurately predict relative probabilities of each species being present in an audio clip



Machine Learning Hurdles

Issues

Small Training Set

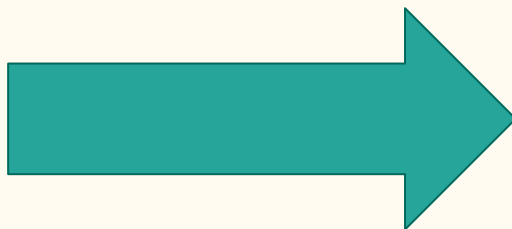
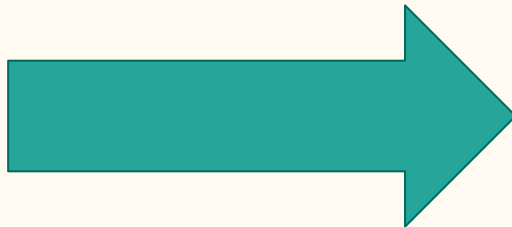
- 1,216 TPs
- 7,781 FPs

Noisy Data

- Insects
- Other Species

Audio Classification Understudied

- Few popular architectures



Solutions

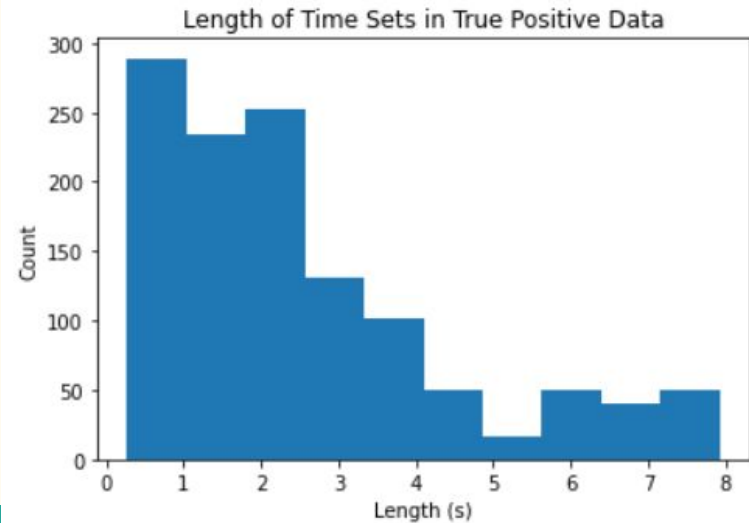
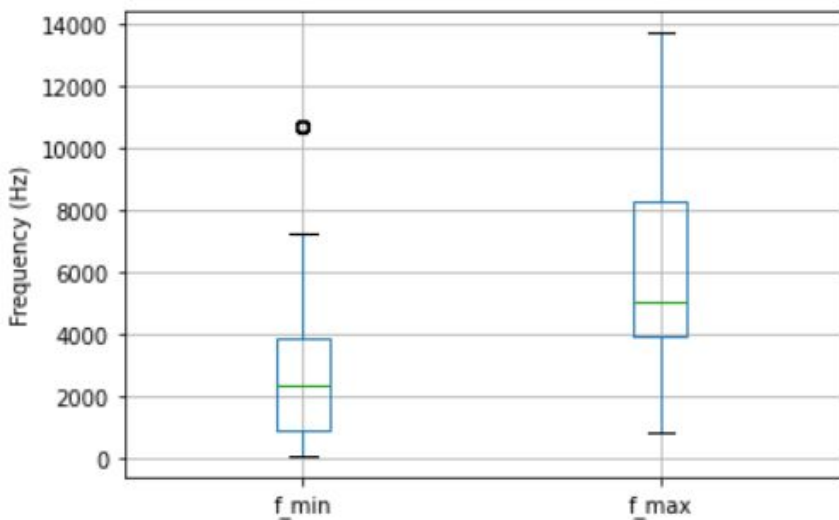
Data Augmentation

Smoothing and Other Preprocessing

Image Classification

Data Exploration

- Even distribution of species in training set
- Wide variance in background noise
- Species are a mix of frogs and birds



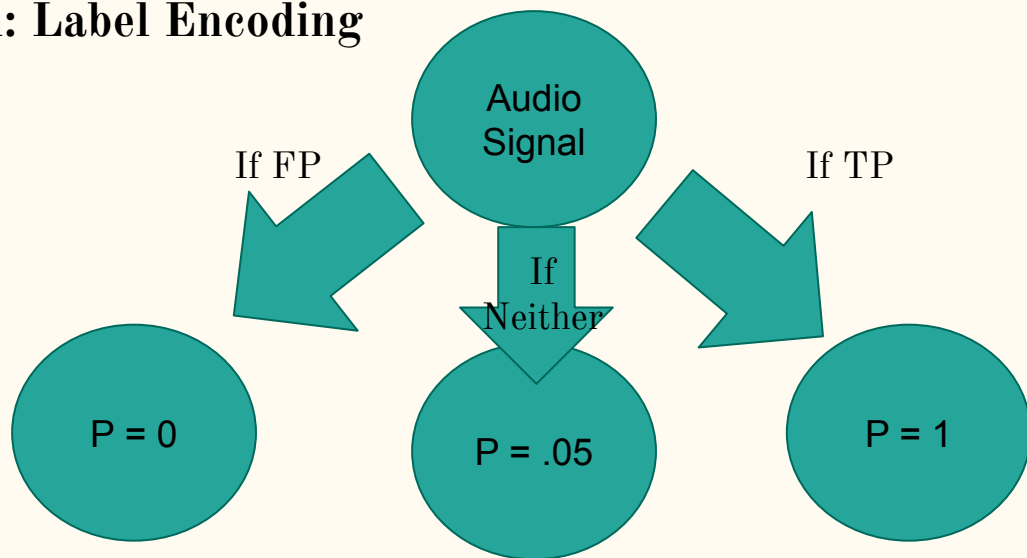
Preprocessing

- Temporal Splicing to 8 seconds surrounding a clip
- Frequency Splicing from 0 to 14000 Hz
- Nearest Neighbors Filtering
- Created “noisy versions” of the data set
 - Duplicates of the original with .2% and .5% maximum amplitude added
- Split into two 75% length clips with 50% overlap
- Converted audio into image

False Positive Data Labels

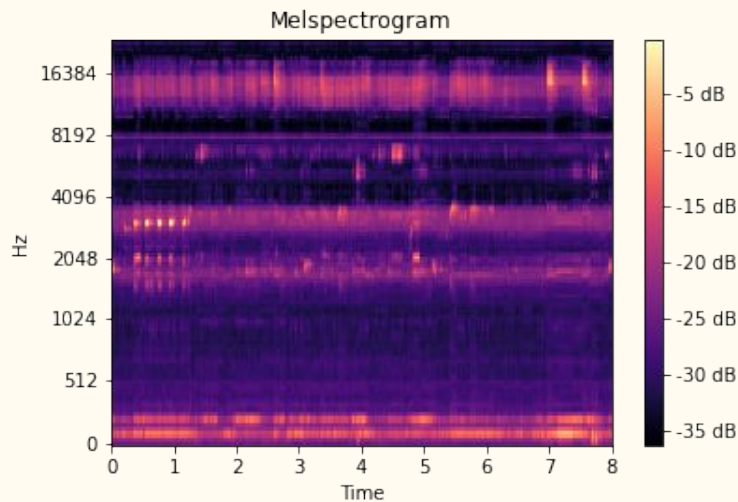
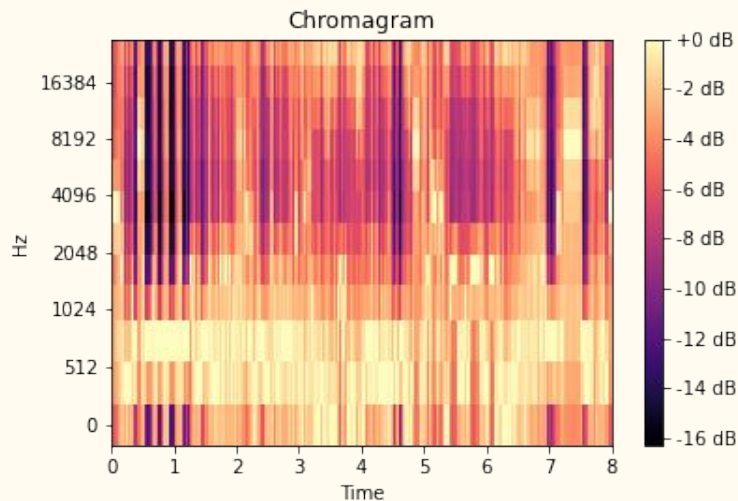
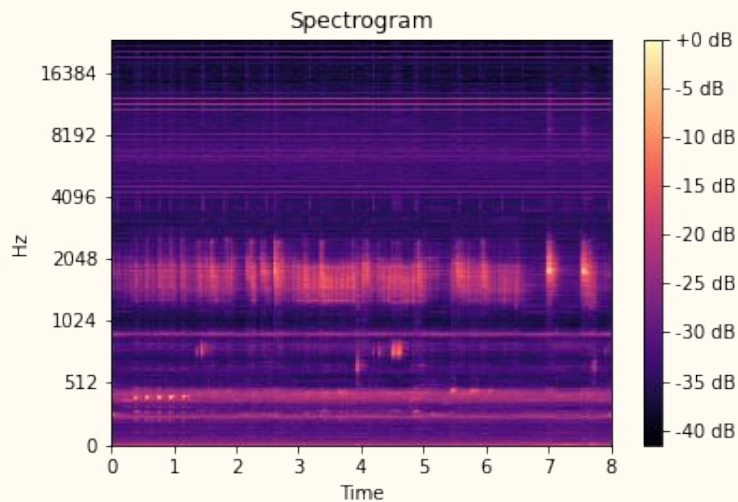
- Experts manually annotated false positives predicted by a baseline model
- Signals which are likely to be classified incorrectly by the model

Our Solution: Label Encoding



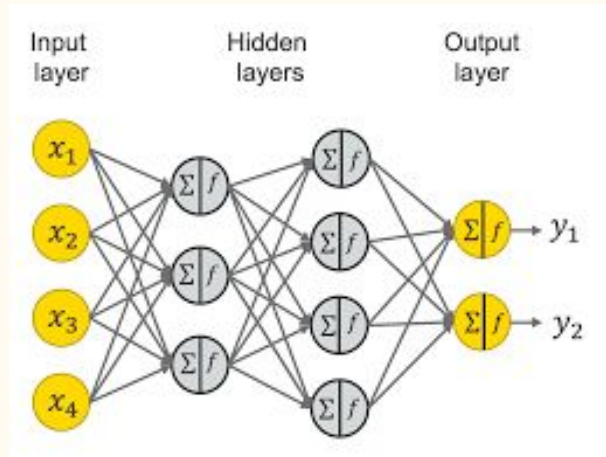
Audio to Image

- Spectrogram
- Melspectrogram
- Chromagram



Neural Networks and Transfer Learning

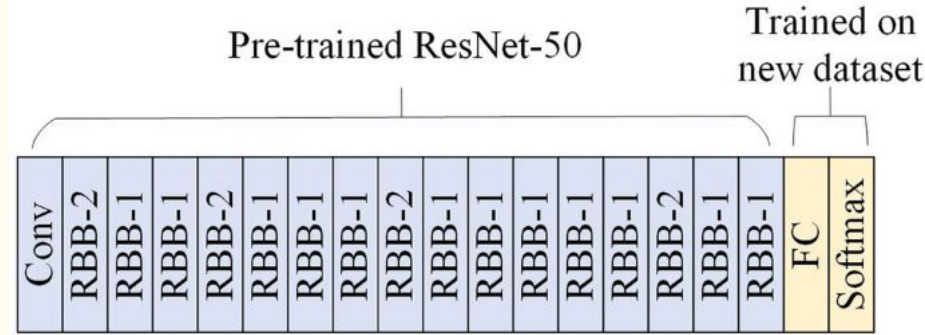
Neural networks use a set of algorithms modeled after the human brain to identify patterns in data.



Transfer learning uses weights from neural networks trained on other datasets



Model 1 - ResNet50



50 layer deep convolutional neural network which is pretrained on millions of images.

- ResNet is the most popular image classification model

```
model = models.resnet50(pretrained=True)
model.fc = nn.Sequential(
    nn.Linear(2048, 1024),
    nn.ReLU(),
    nn.Dropout(p=0.2),
    nn.Linear(1024, 1024),
    nn.ReLU(),
    nn.Dropout(p=0.2),
    nn.Linear(1024, num_species)
)
```

Optimal Hyperparameters and Model Components:

Learning rate = .001

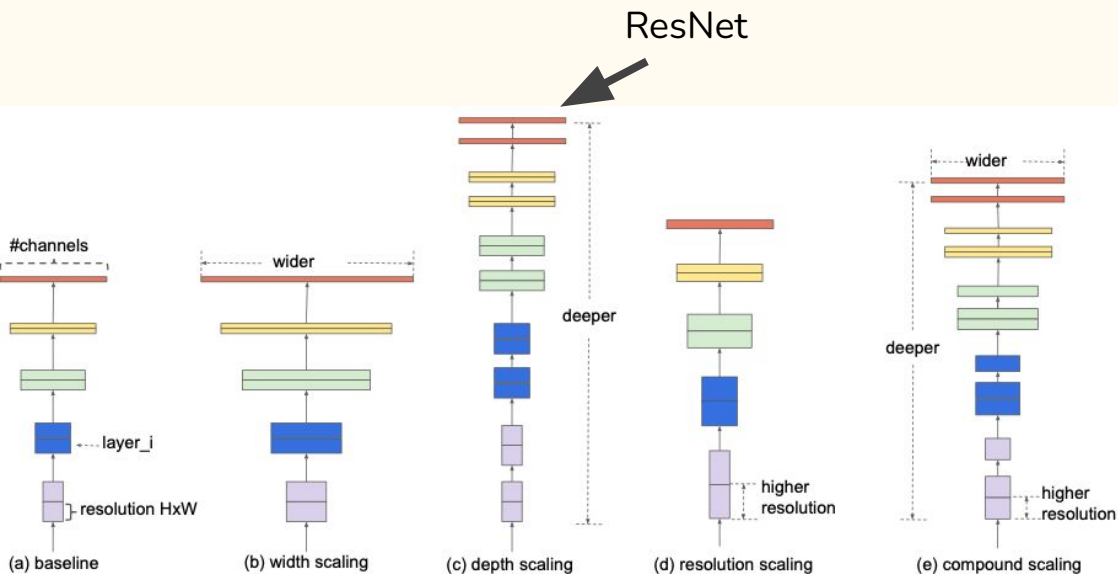
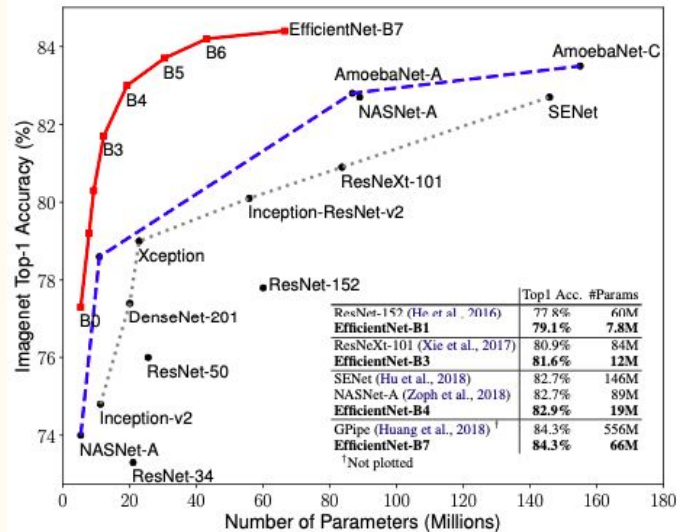
Optimizer = SGD

Momentum = .9

Loss = BCE w/Logits Loss

Model 2 - EfficientNet7

- Scales width, depth, and resolution using equal ratios



Optimal Hyperparameters and Model Components:

Learning rate = .001

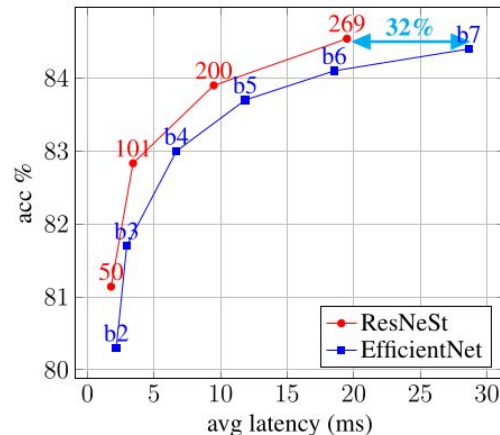
Optimizer = SGD

Momentum = .9

Loss = BCE w/Logits Loss

Model 3 - ResNeSt

- Split-Attention NN designed to improve on EfficientNet and ResNet
 - “Individual salient attributes for different visual features”
 - Correlated vs Independent* visual features
- Expects 3-channel RGB images, normalized using mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]



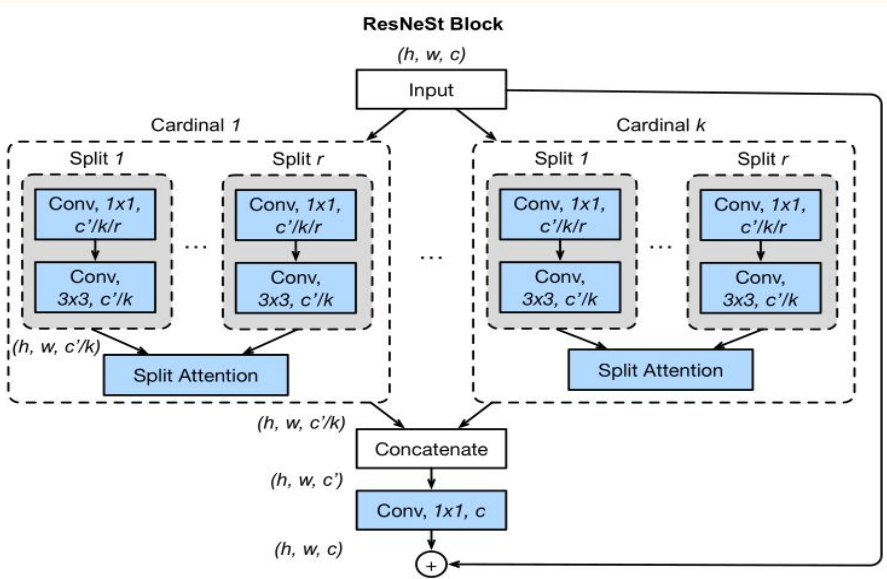
Optimal Hyperparameters and Model Components:

Learning rate = .001

Momentum = .7

Optimizer = SGD

Loss = Cross Entropy



Metric

Label-weighted label-ranking average precision (LWL RAP)

- Generalization of mean reciprocal rank measure
- Each label receives equal weight
- Each test observation is weighted differently based on number of true labels

Test Prediction 1

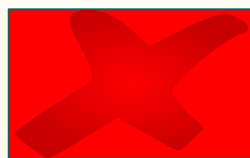
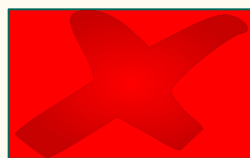


1/1



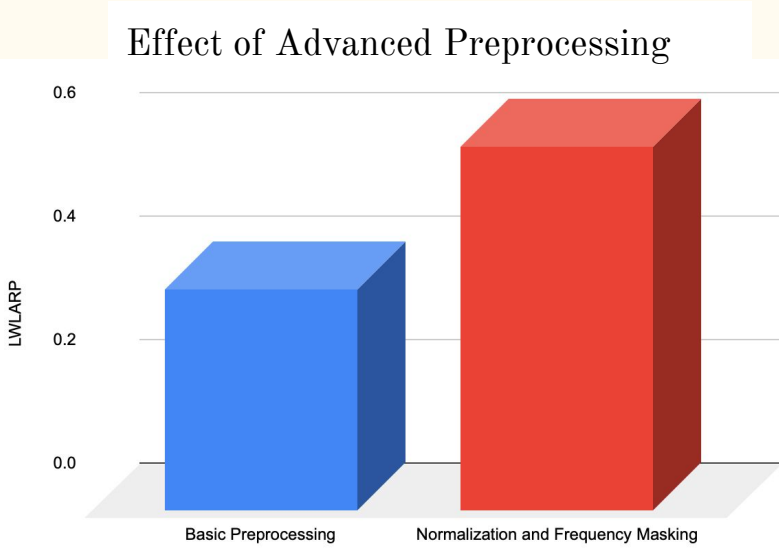
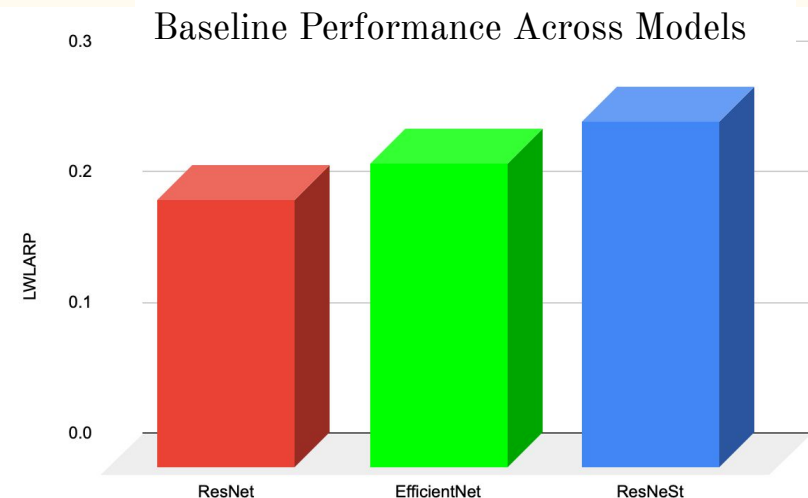
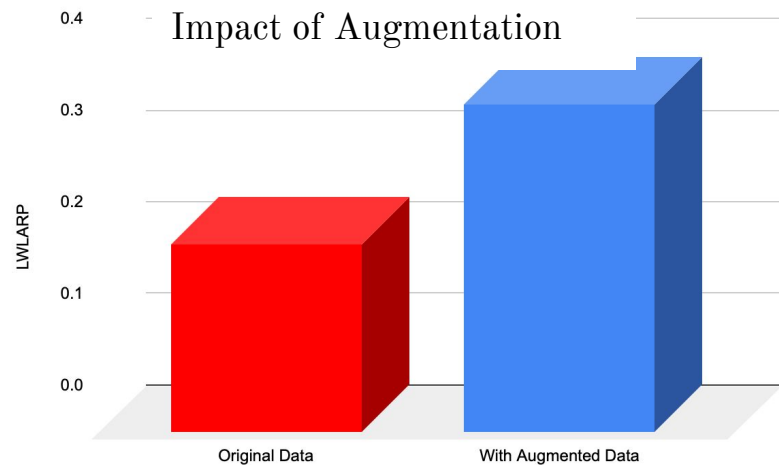
2/3

Test Prediction 2

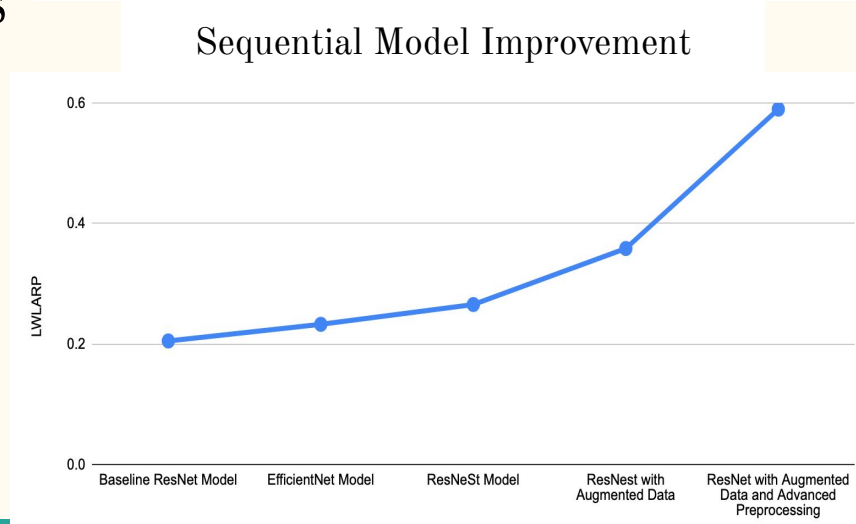


1/3

$$\begin{aligned}\text{LWL RAP} &= (1 + .67 + .33)/3 \\ &= 0.67\end{aligned}$$



Results



Conclusions

- Converting audio classification into a image classification problem is effective
- Data augmentation is useful to increase training data size
- For this problem, preprocessing is more important than model fine-tuning