# Intelligent Audio Systems

Prof. Cheng Siong CHIN

cheng.chin@ncl.ac.uk
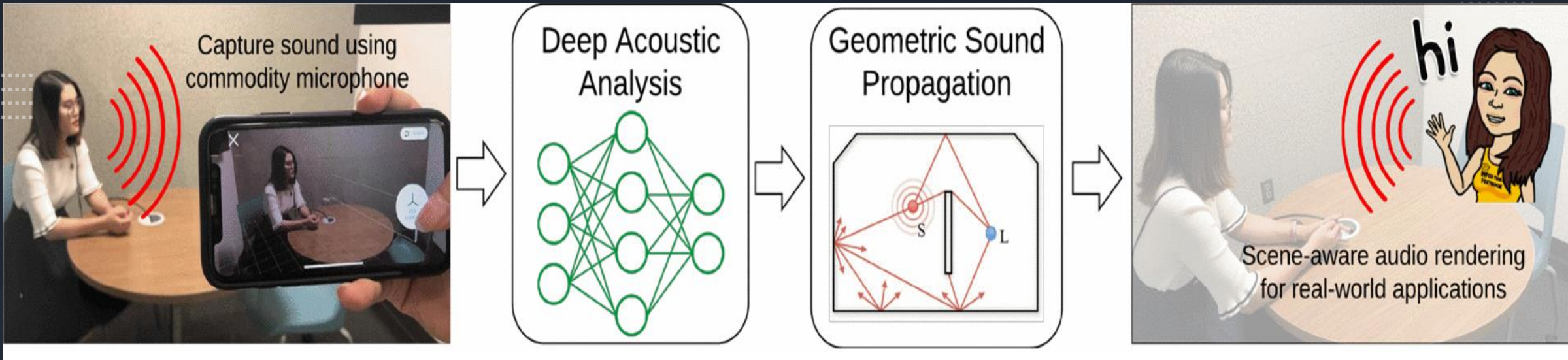
# Dr Cheng Siong Chin

## Research Interests:

Intelligent Systems Design of complex systems in uncertain environment involving Predictive Analytics (data mining, predictive modelling and machine learning).

## Research Projects:

❖ 2019-present: EMA grant on **AI** System for **Energy Storage i**n Hot and Humid Climate

❖ 2020-present: EDB-IPP grant on **Smart Manufacturing** for Yield Improvement

❖ 2018- present: EDB-IPP grant on **Smart** and High Precision Leakage Localization.

❖ 2018 to present: EDB-IPP grant on **Intelligent** Leakage Warning System.

❖ 2016 to 2019: SMI on **Intelligent** Software Tool for **Noise** Modelling and Prediction.

❖ 2013 to 2016: SMI on the **Battery** Power System for **ROV**.

❖ 2013 to 2015: Defence Innovative Research Programme Project on **AUV** Docking Hoop Control

❖ 2013 to 2017: EDB-IPP grant on **Noise** and Vibration Control of Offshore Structure.

❖ 2013 to 2018: EDB-IPP grant on Vibration and Psycho-**Acoustic** Parameters in Hard Disk Drive.



Email: cheng.chin@ncl.ac.uk

# Intelligent Systems→ Problem, Need and Industrial Relevance

- Intelligent Systems heavily rely of context awareness
- Context awareness is the ability of a system or system component to gather information about its environment at any given time and adapt behaviors accordingly.
- Being able to detect and classify sounds and events from the surroundings is an important aspect of context awareness.
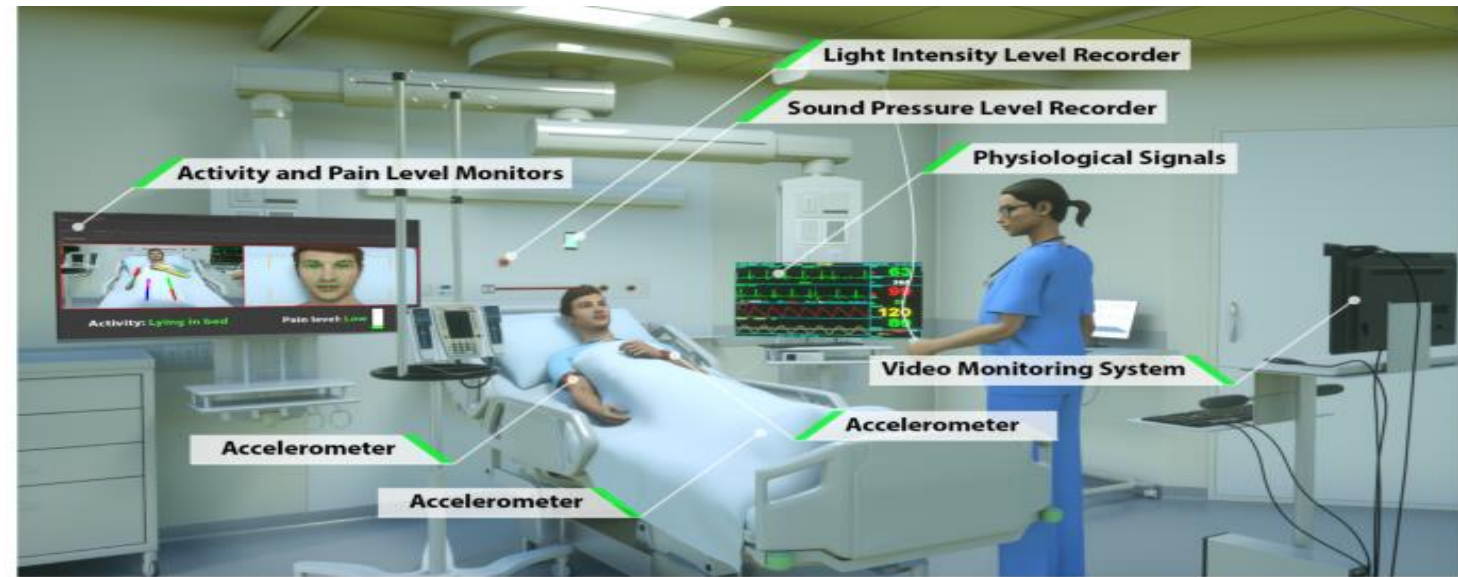
# Background

❑ Sound Classification is one of the most widely used applications in Audio Deep Learning.

❑ It involves learning to classify sounds and to predict the category of that sound.

❑ This type of problem can be applied to many practical scenarios e.g. classifying music clips to identify the genre of the music, acoustic scenes classification, and sound events detection.
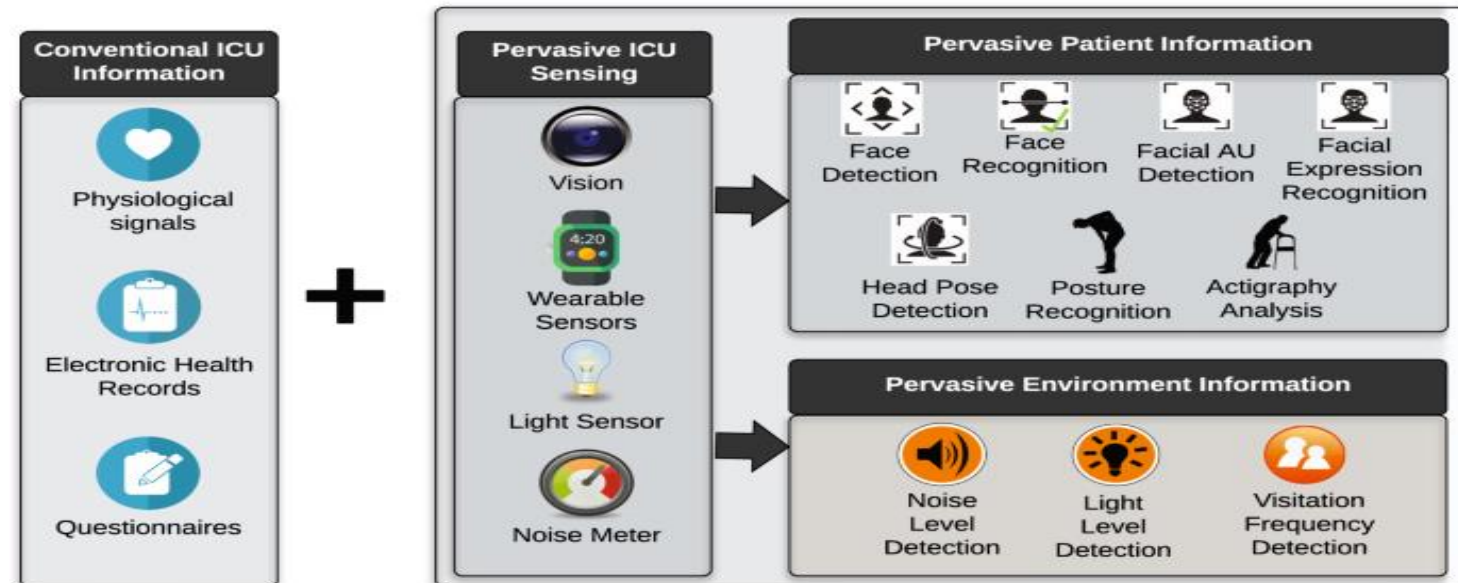
❑ It serves many important applications.

For example, to capture the acoustic characteristics of real-world rooms using commodity devices and use the captured characteristics to generate similar sounding sources with virtual models.

# Background

- AI in the critical care setting could reduce nurses' workload to allow them to spend time on more critical tasks, and could also augment human decision-making by offering low-cost and high capacity intelligent data processing.

- How pervasive sensing technology and AI can be used for monitoring patients and their environment in the ICU.

- Wearable accelerometer sensors, a light sensor, a sound sensor, and a high-resolution camera to capture data on patients and their environment in the ICU
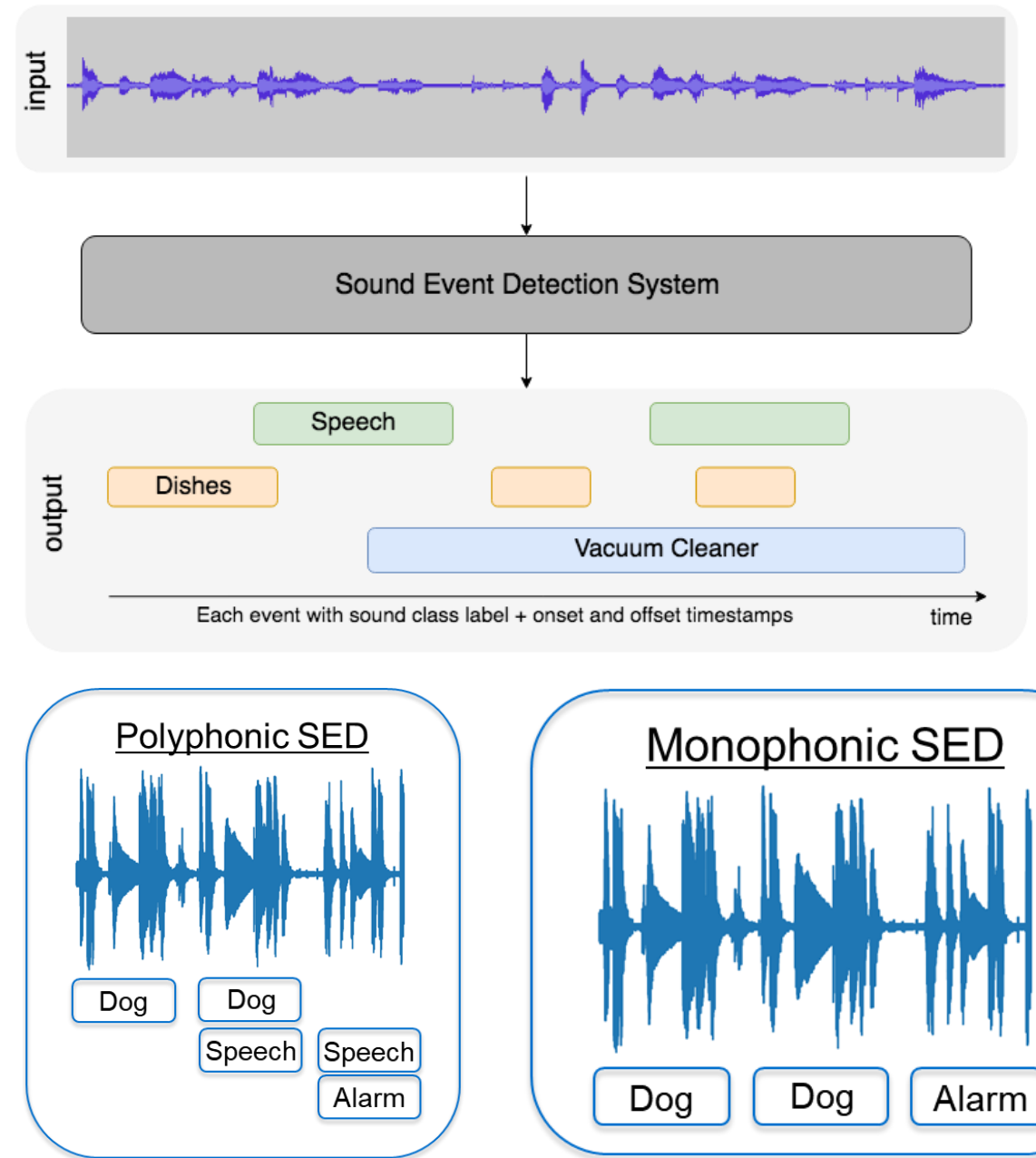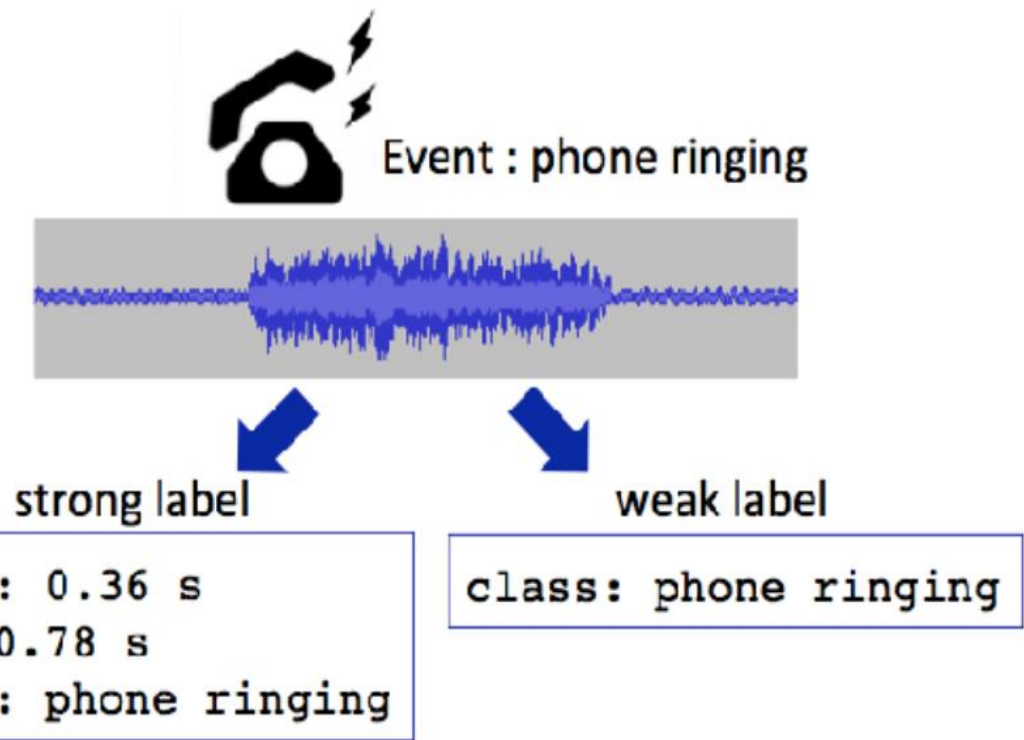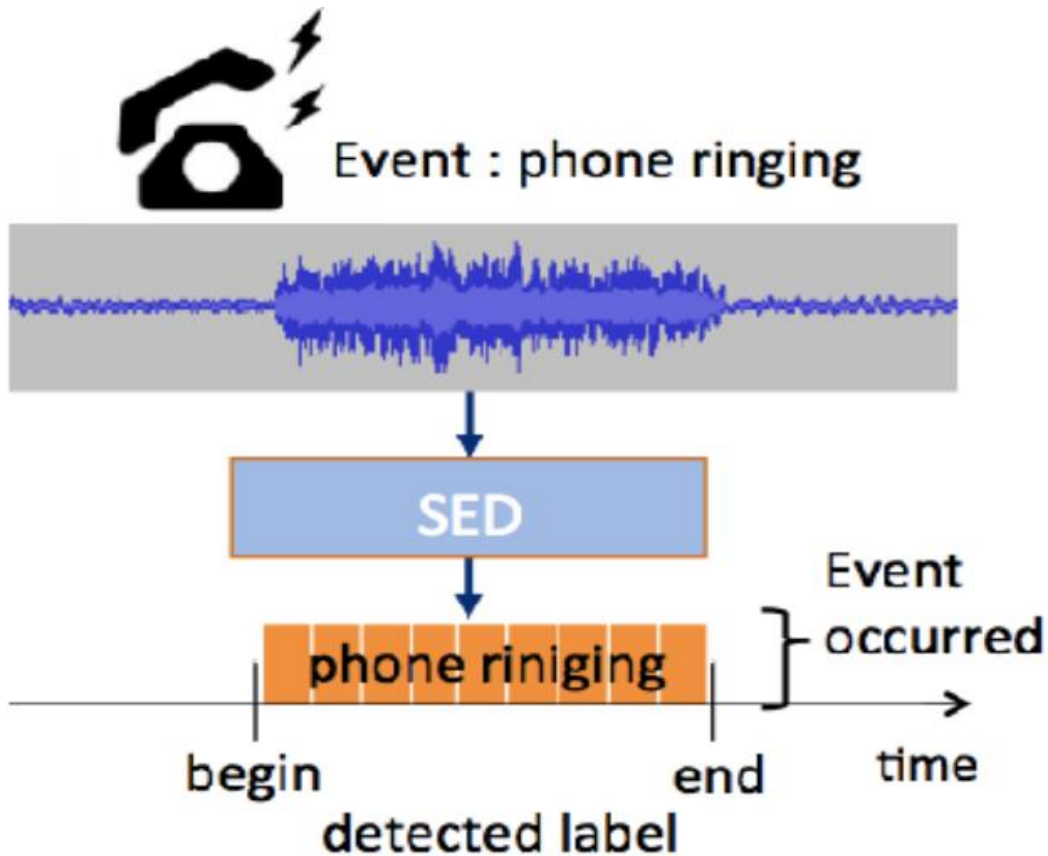


https://www.nature.com/articles/s41598-019-44004-w

# Sound Events Detection (SED)

# Sound Event Detection

❑ A Sound Event Detection (SED) refers to a system that is capable of
  ❑ Detecting the type of acoustic event present in an audio clip (i.e., audio tagging).
  ❑ Returning the onset and offset of the identified acoustic event (i.e., temporal localization).

❑ A SED system can be categorized as a monophonic or polyphonic system.

❑ A polyphonic SED system is more appropriate in real-life application as an audio clip is more likely to contain multiple sound events.

# Sound Event Detection



Event : phone ringing

SED

phone riniging — Event occurred

begin — end — time

detected label

https://www.semanticscholar.org/paper/Weakly-Labeled-Learning-Using-BLSTM-CTC-for-Sound-Matsuyoshi-Komatsu/401296f7d4058ef1702a77874c95a759f3868fc4

Event : phone ringing

strong label

begin: 0.36 s
end: 0.78 s
class: phone ringing

weak label

class: phone ringing

Weakly labeled data- where only the event tags are known with certainty → can compromise the solution.

# Sound Event Detection

- ❑ How is a SED system useful?
    - ❑ Not affected by the degree of illumination.
    - ❑ Some events can only be detected by sound.
    - ❑ Sound can capture the immediate attention of an individual.
    - ❑ An audio clip requires lesser computational resources than an image or video

# Challenges for Sound Event Detection

❑ System development may require a large amount of strongly labeled data, where event tags and their corresponding onset and offsets are known with certainty

    ❑ Difficult and time consuming to collect.

    ❑ Accuracy of onset and offset annotation is ambiguous due to the fade in and fade out effect.

    ❑ Such dataset is usually limited to minutes or a few hours.
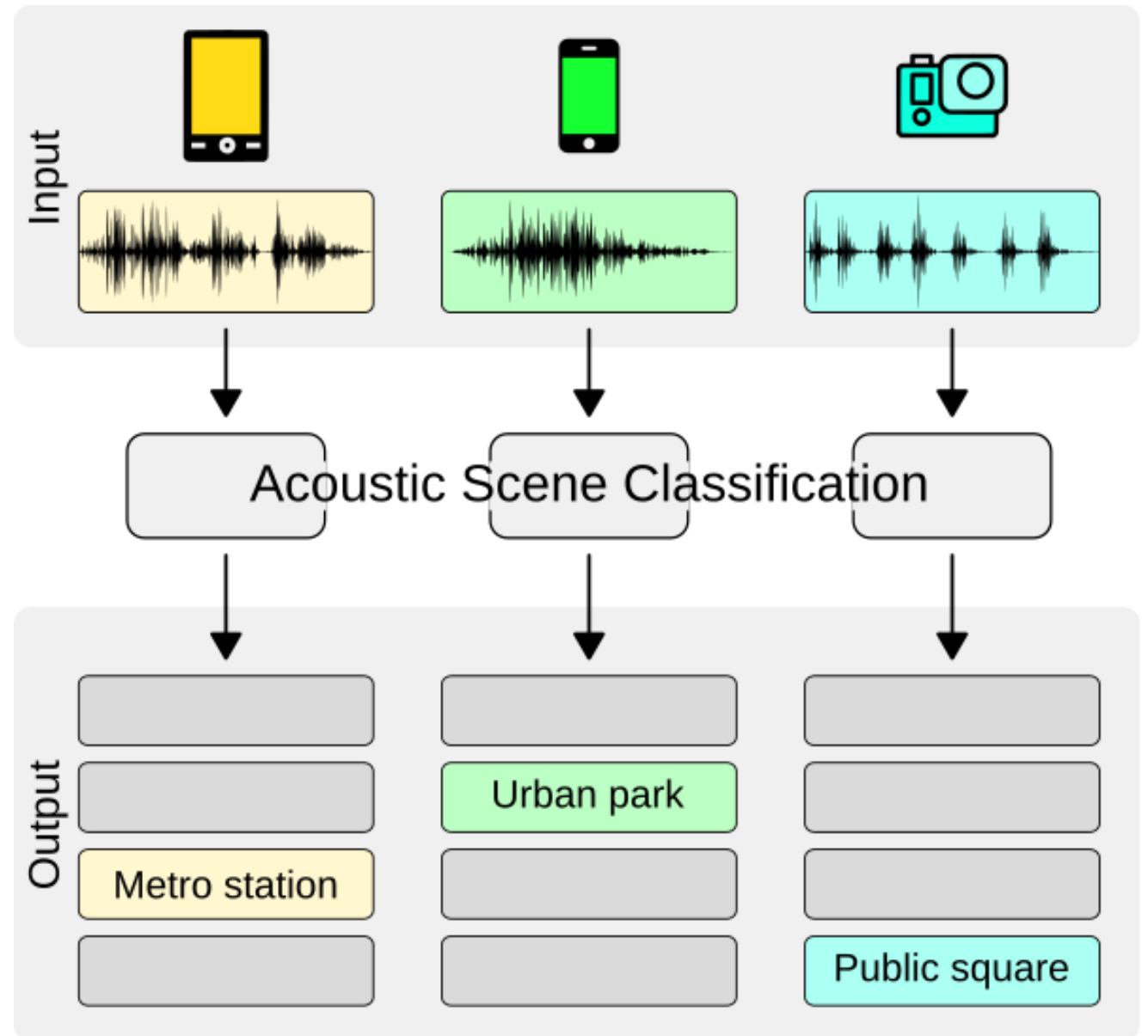
# DEMO

- Sound Classification System- apply the audio tagging system to build a sound event detection (SED) system.

- The SED prediction is obtained by applying the audio tagging system on consecutive 2-second segments.

# VIDEO

Qiuqiang Kong, Yin Cao, Turab Iqbal, Yuxuan Wang, Wenwu Wang, Mark D. Plumbley. "PANNs: Large-Scale Pretrained Audio Neural Networks for Audio Pattern Recognition." arXiv preprint arXiv:1912.10211 (2019).
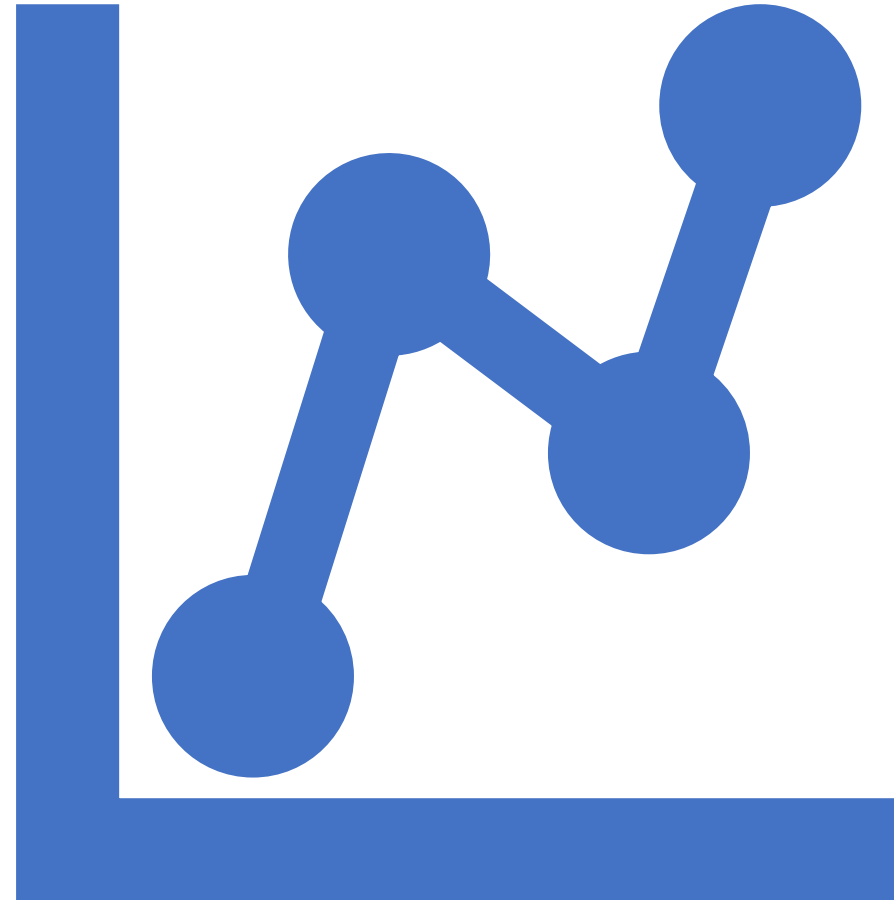
# Acoustic Scenes Classification (ASC)

# Acoustic Scene Classification

❑Acoustic scene classification (ASC) categorizes an audio file based on the environment in which it has been recorded.

❑ For example, →adjustment of a smartphone ring volume when its owner moves from a quiet acoustic environment into a noisier one.
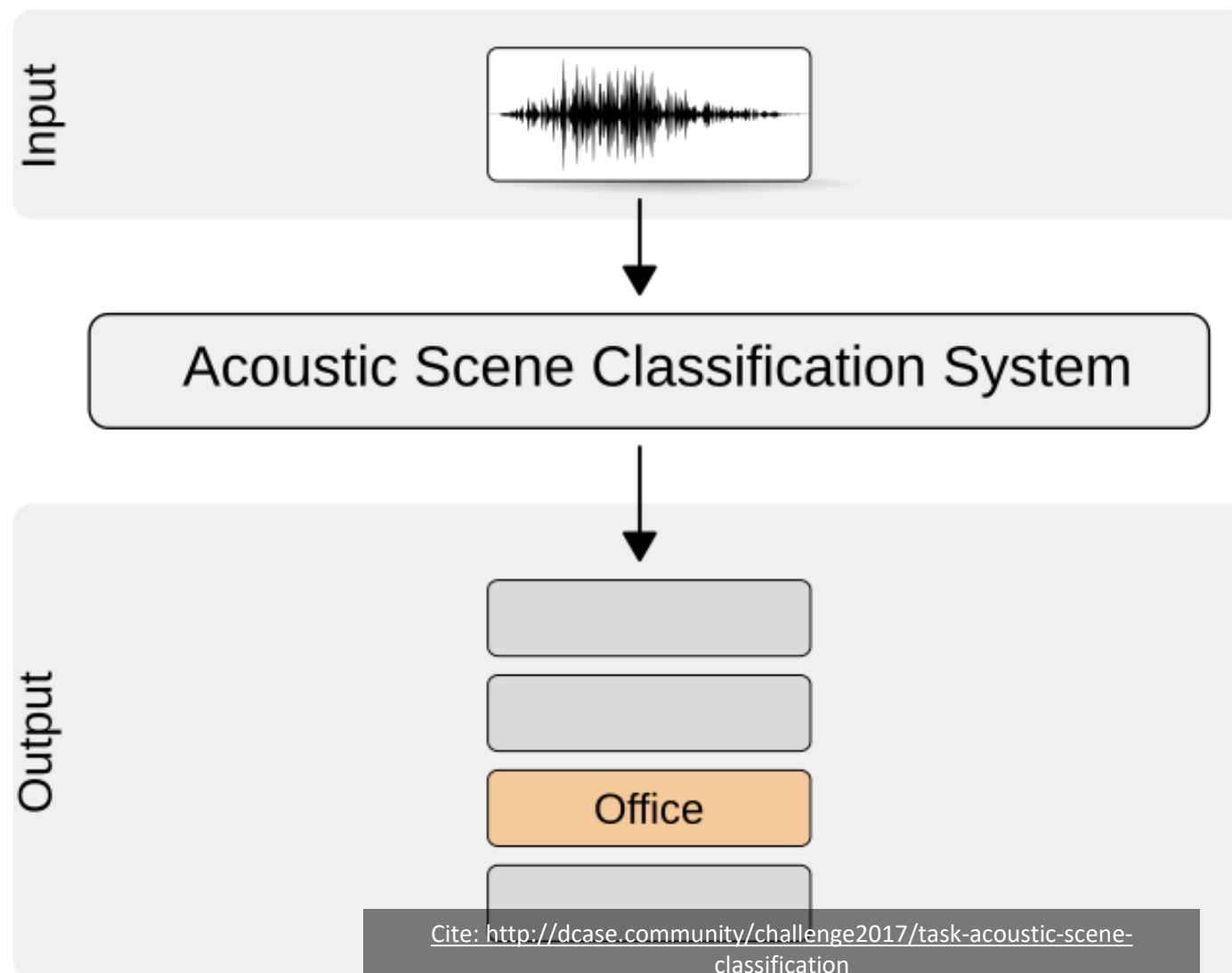
# Challenges for Acoustic Scene Classification

---

❑ Time consuming to collect for each scene.

❑ Dataset is usually limited to minutes or a few hours. Thus, short audio segments provide less information, making ASC difficult.

❑ Acoustic scene can share very similar acoustic profile making it difficult to be classified.

❑ Ambient sound scenes typically comprise multiple short events occurring on top of a somewhat stationary background.

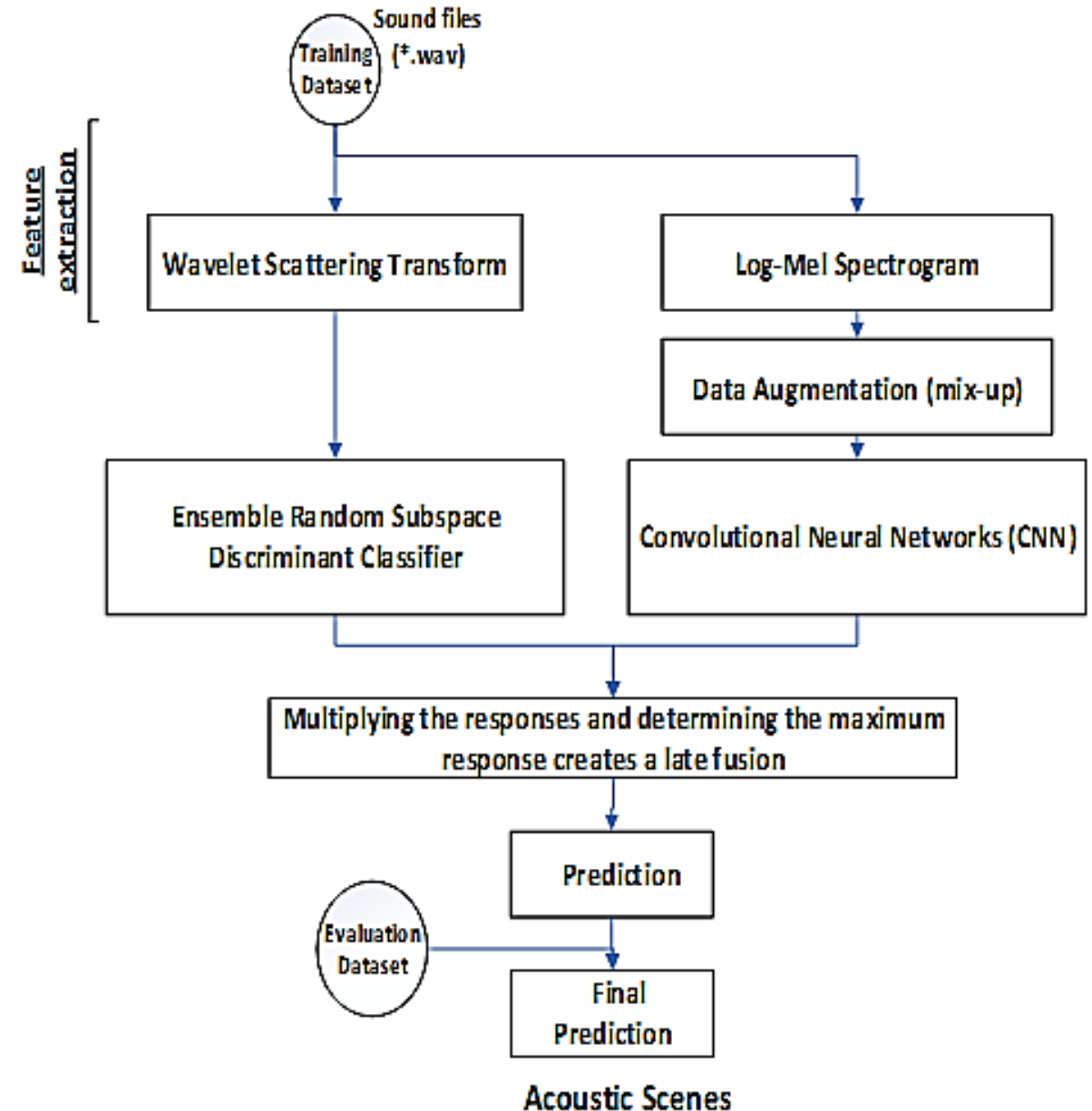❑ foreground-background ambient sound scene separation is challenging.

# Introduction

➢ The top-performing systems in Detection and Classification of Acoustic Scenes and Events (DCASE) challenge consist of multiple models using CNN.

➢ Instead of log-Mel spectrogram, for feature extraction, wavelet time scattering has been used.

➢ The extracted features are insensitive to translation, rotation, and small deformation.

➢ Wavelet scattering transform produce a more accurate classifier in dealing with variation in the ASC.

**Input**

## Acoustic Scene Classification System

**Output**

Office

Cite: http://dcase.community/challenge2017/task-acoustic-scene-classification

# Introduction

➤ **Overview:** We will start with sound files, convert them into spectrograms, input them into a CNN plus Linear Classifier model, and produce predictions about the class to which the sound belongs.

➤ The **wavelet features** are used for ensemble classifiers (different discriminant analysis learners-linear discriminant analysis (LDA), Quadratic discriminant analysis (QDA), Regularized linear discriminant analysis (RDA) ).

➤ To reduce the overfitting of the dataset in CNN, a Mixup algorithm is used. They are used to train a CNN.

➤ A multi-model late fusion is then used to fuse → log-Mel CNN & Wavelet ensemble classifier.

➤ It will provide a time-frequency representation of audio to capture Spectro-temporal modulation patterns for identifying various acoustic scenes.



Feature extraction

Training Dataset → Sound files (*.wav)

Wavelet Scattering Transform

Log-Mel Spectrogram

Data Augmentation (mix-up)

Ensemble Random Subspace Discriminant Classifier

Convolutional Neural Networks (CNN)

Multiplying the responses and determining the maximum response creates a late fusion

Prediction

Evaluation Dataset

Final Prediction

**Acoustic Scenes**

# Datasets

1) bus
2) forest path
3) home
4) city center
5) cafe
6) lakeside beach (outdoor)
7) library (indoor)
8) car
9) grocery store
10) urban park (outdoor)
11) office: multiple persons (indoor)
12) metro station (indoor)
13) train (traveling, vehicle)
14) residential area (outdoor)
15) tram (traveling, vehicle)

The TUT Acoustic Scenes 2017 dataset was used for the experiments.

It consists of 10-seconds audio segments from 15 acoustic scenes

Each acoustic scene has 312 segments giving a total of 52 (=312*10/60) minutes of audio.

Sound recordings were performed via different devices at 24-bit resolution and 44100Hz sampling rate. The microphones are worn during recording.
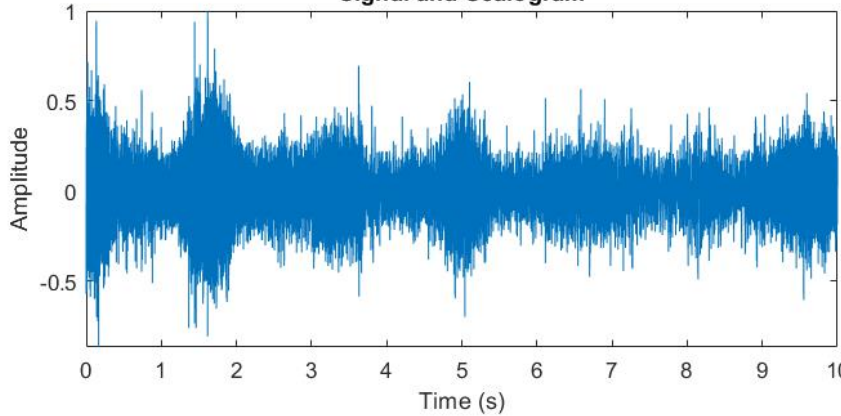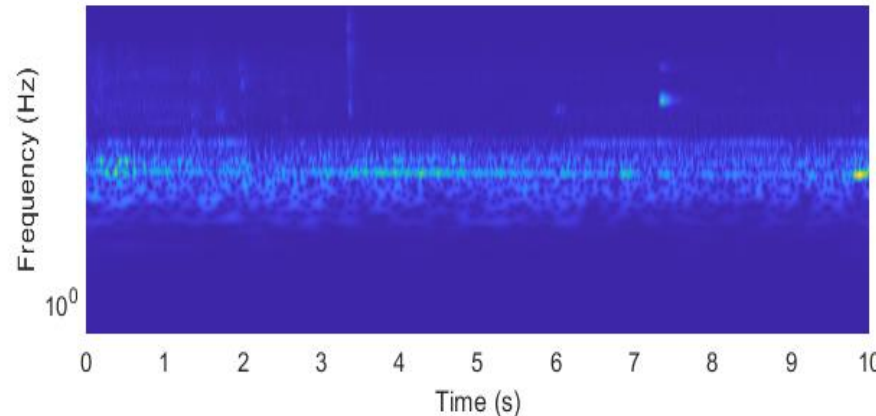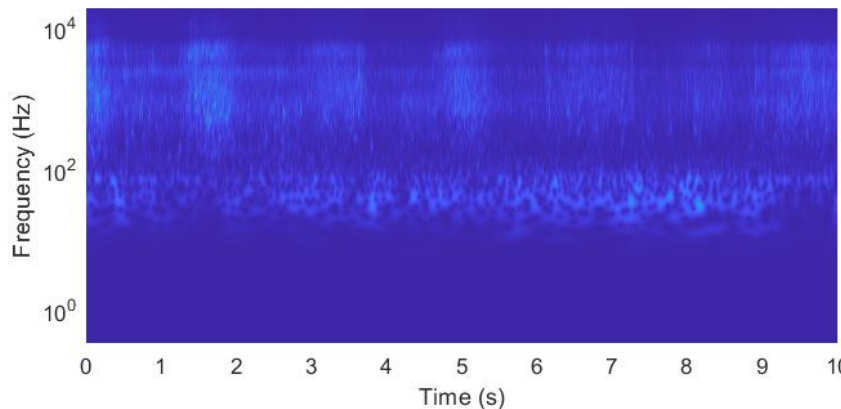
https://doi.org/10.5281/zenodo.400515

# Sample Data

**Beach**

**Train**

> The scalogram is the absolute value of the continuous wavelet transform (CWT) of a signal, plotted as a function of time and frequency.

> The scalogram analyzes real-world signals with features occurring at different scales — for example, signals with slowly varying events punctuated by abrupt transients.

> Give better time localization for short-duration, high-frequency events, and better frequency localization for low-frequency, longer-duration events.



Signal and Scalogram



Signal and Scalogram

Sample Data

Beach

# Sample Data

| train_3043 | train_3044 | train_3045 | train_3046 | train_3047 | train_3048 | train_3049 | train_3050 | train_3051 | train_3052 | train_3053 | train_3054 |
| train_3055 | train_3056 | train_3057 | train_3058 | train_3059 | train_3060 | train_3061 | train_3062 | train_3063 | train_3064 | train_3065 | train_3066 |
| train_3067 | train_3068 | train_3069 | train_3070 | train_3071 | train_3072 | train_3073 | train_3074 | train_3075 | train_3076 | train_3077 | train_3078 |
| train_3079 | train_3080 | train_3081 | train_3082 | train_3083 | train_3084 | train_3085 | train_3086 | train_3087 | train_3088 | train_3089 | train_3090 |
| train_3091 | train_3092 | train_3093 | train_3094 | train_3095 | train_3096 | train_3097 | train_3098 | train_3099 | train_3100 | train_3101 | train_3102 |

# Data Augmentation-Mixup

The technique literally mixing up the features and their corresponding labels.

$$\tilde{x} = \lambda x_i + (1 - \lambda)x_j, \qquad \text{where } x_i, x_j \text{ are raw input vectors}$$

$$\tilde{y} = \lambda y_i + (1 - \lambda)y_j, \qquad \text{where } y_i, y_j \text{ are label encodings}$$

Note that the lambda values are values with the [0, 1] range

# 1) Convolutional Neural Network (CNN)

- The Batch Normalization (BN) and rectified linear unit (ReLU) are used.

- The ReLU increases the non-linearity in the images.

- The batch normalization learning is used as a regularization to prevent overfitting.

- The activation function and BN are located before the convolution layer to improve the acoustic classification accuracy.

- The max-pooling layers come after the convolution process.

| Description of each layer |
|---|
| imageInputLayer- 128×42×2 |
| batchNormalizationLayer |
| convolution2dLayer- 32 filters (3×3) and zero padding |
| batchNormalizationLayer |
| reluLayer |
| convolution2dLayer- 32 filters (3×3) and zero padding |
| batchNormalizationLayer |
| reluLayer |
| maxPooling2dLayer- pool size 3×3, stride 2×2 and zero padding |
| convolution2dLayer- 32 filters (3×3) and zero padding |
| batchNormalizationLayer |
| reluLayer |
| convolution2dLayer- 32 filters (3×3) and zero padding |
| batchNormalizationLayer |
| reluLayer |
| maxPooling2dLayer- pool size 3×3, stride 2×2 and zero padding |
| convolution2dLayer- 128 filters (3×3) and zero padding |
| batchNormalizationLayer |
| reluLayer |
| convolution2dLayer- 128 filters (3×3) and zero padding |
| batchNormalizationLayer |
| reluLayer |
| maxPooling2dLayer- pool size 3×3, stride 2×2 and zero padding |
| convolution2dLayer- 256 filters (3×3) and zero padding |
| batchNormalizationLayer |
| reluLayer |
| convolution2dLayer- 256 filters (3×3) and zero padding |
| batchNormalizationLayer |
| reluLayer |
| averagePooling2dLayer-pool size 16 ×6 |
| dropoutLayer(0.5) |
| fullyConnectedLayer(15) |
| softmaxLayer |
| classificationLayer |

# 1)Convolutional Neural Network (CNN)

➢ The feature map that includes a prominent feature is obtained from the output of the max-pooling layer.

➢ The average pooling reduces the activation by combining the non-maximal activations.

➢ The last few layers consist of a dropout layer that removes 50% of the visible and hidden units to reduce overfitting.

➢ The fully connected layer is compiled the data to form the output for the last second layer that uses the softmax activation function to obtain probabilities of the input from the 15 classes.

➢ Lastly, the last classification layer produces the final classification

| Description of each layer |
|---|
| imageInputLayer- $128\times42\times2$ |
| batchNormalizationLayer |
| convolution2dLayer- 32 filters ($3\times3$) and zero padding |
| batchNormalizationLayer |
| reluLayer |
| convolution2dLayer- 32 filters ($3\times3$) and zero padding |
| batchNormalizationLayer |
| reluLayer |
| maxPooling2dLayer- pool size $3\times3$, stride $2\times2$ and zero padding |
| convolution2dLayer- 32 filters ($3\times3$) and zero padding |
| batchNormalizationLayer |
| reluLayer |
| convolution2dLayer- 32 filters ($3\times3$) and zero padding |
| batchNormalizationLayer |
| reluLayer |
| maxPooling2dLayer- pool size $3\times3$, stride $2\times2$ and zero padding |
| convolution2dLayer- 128 filters ($3\times3$) and zero padding |
| batchNormalizationLayer |
| reluLayer |
| convolution2dLayer- 128 filters ($3\times3$) and zero padding |
| batchNormalizationLayer |
| reluLayer |
| maxPooling2dLayer- pool size $3\times3$, stride $2\times2$ and zero padding |
| convolution2dLayer- 256 filters ($3\times3$) and zero padding |
| batchNormalizationLayer |
| reluLayer |
| convolution2dLayer- 256 filters ($3\times3$) and zero padding |
| batchNormalizationLayer |
| reluLayer |
| averagePooling2dLayer-pool size $16\times6$ |
| dropoutLayer(0.5) |
| fullyConnectedLayer(15) |
| softmaxLayer |
| classificationLayer |

# 2) Wavelet Scattering

The next step involves feature extraction using wavelet scattering for subsequent ensemble classifiers.

The parameters of the transform are the filter-bank (using 1D Morlet wavelets) resolutions $Q1=1$ and $Q2=4$.

The duration 0.75s of the averaging filter (or invariance scale) is used for the modulation structure duration.

The sampling frequency is 44100 Hz.

# 3) Ensemble Classifiers

➢ The proposed ensemble classifiers include different discriminant analysis learners, such as linear discriminant analysis (LDA), Quadratic discriminant analysis (QDA), and Regularized linear discriminant analysis (RDA) with other predictors covariance treatments.

➢ The random subspace learning method is used to increase the acoustic classification accuracy.

➢ In the random subspace, the feature subspaces are chosen randomly from the original feature space.

➢ The final prediction of these individual classifiers is then obtained using majority voting.

# 4) Fusion of CNN and Classifiers

➢The fusion of the CNN and classifier prediction results indicates the relative confidence of their prediction.

➢Multiplying the responses and determining the maximum response creates a late fusion system that inherent in the merits of each method.
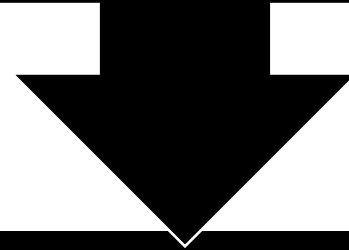
$$class\_pred^i = \mathrm{argmax}\left(prob^i_{CNN}, prob^i_{ensem\_class}\right)$$

where $prob^i_{CNN}$ and $prob^i_{ensem\_class}$ are the probabilities of sound recording $i$ from CNN and ensemble classifiers, respectively.

# Results and Discussion

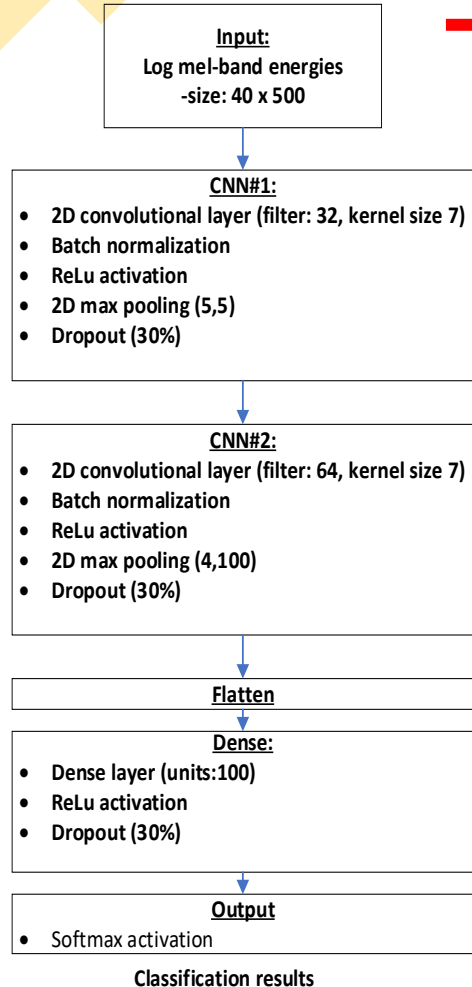**Run experiment using: Intel® Core i7-9750H CPU, 2.6GHz, 6 Cores, and Geforce RTX 2060**

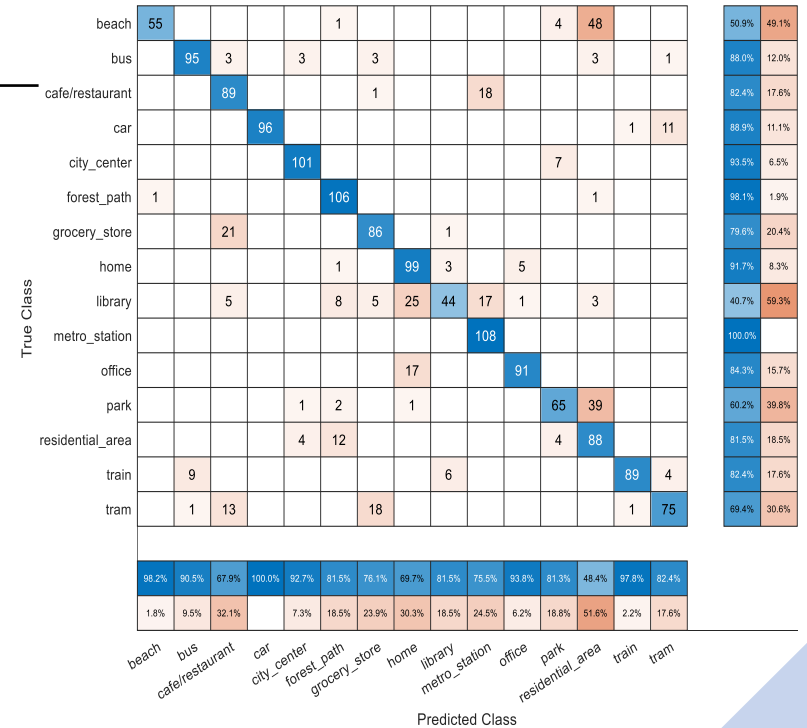**The configurations of the proposed model**

| Stochastic gradient descent with momentum optimizer with a learning rate: 0.05s | Size of the mini-batch for each training iteration: 128 | Momentum: 0.9 | Maximum number of epochs: 8 | Factor for L2 regularization: 0.005 | Number of epochs for dropping the learning rate: 2 | Multiplicative factor applied to the learning rate for each epoch: 0.2 |

# Results and Discussion
# -Comparison with Baseline Model

**Input:**
Log mel-band energies
-size: 40 x 500

**CNN#1:**
- 2D convolutional layer (filter: 32, kernel size 7)
- Batch normalization
- ReLu activation
- 2D max pooling (5,5)
- Dropout (30%)

**CNN#2:**
- 2D convolutional layer (filter: 64, kernel size 7)
- Batch normalization
- ReLu activation
- 2D max pooling (4,100)
- Dropout (30%)

**Flatten**

**Dense:**
- Dense layer (units:100)
- ReLu activation
- Dropout (30%)

**Output**
- Softmax activation

**Classification results**

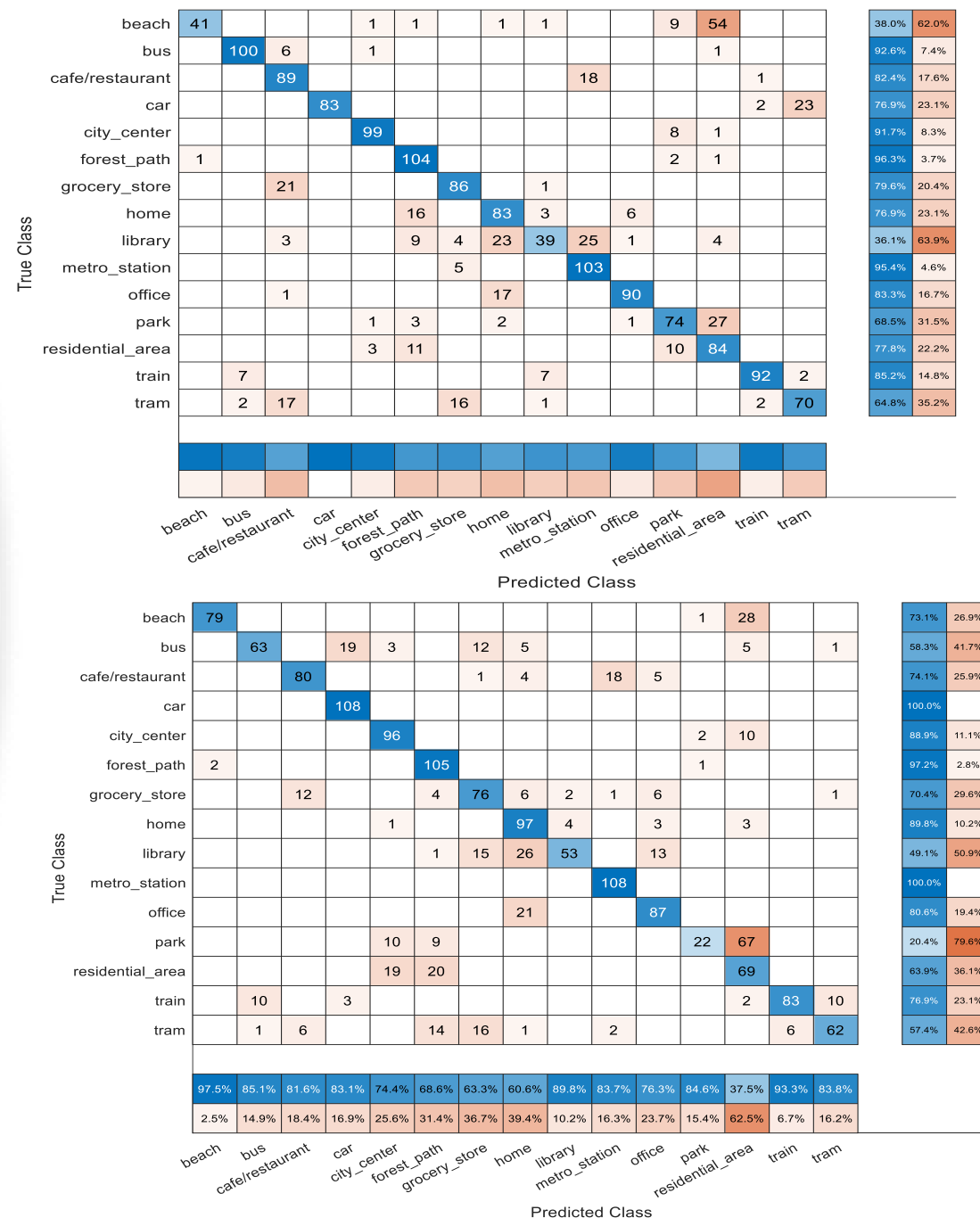| Scenes | Acoustic Classification Accuracy (%) | | | |
|--------|------------------------|-----------|----------------------|-------------|
| | Baseline Model [32] | CNN Model | Ensemble Classifiers Model | Fused Model |
| beach | 40.7 | 73.1 | 37.9 | 50.9 |
| bus | 38.9 | 58.3 | 92.5 | 87.9 |
| cafe/restaurant | 43.5 | 74.0 | 82.4 | 82.4 |
| car | 64.8 | 100 | 76.8 | 88.8 |
| city-center | 79.6 | 88.8 | 91.6 | 93.5 |
| forest path | 85.2 | 97.2 | 96.2 | 98.1 |
| grocery store | 49.1 | 70.3 | 79.6 | 79.6 |
| home | 76.9 | 89.8 | 76.8 | 91.6 |
| library | 30.6 | 49.0 | 36.1 | 40.7 |
| metro station | 93.5 | 100 | 95.3 | 100 |
| office | 73.1 | 80.5 | 83.3 | 84.2 |
| park | 32.4 | 20.3 | 68.5 | 60.1 |
| residential area | 77.8 | 63.8 | 77.7 | 81.4 |
| train | 72.2 | 76.8 | 85.1 | 82.4 |
| tram | 57.4 | 57.4 | 64.8 | 69.4 |
| Average | 61.0 | 73.3 | 76.3 | 79.4 |



➢ The ensemble classifiers have a higher acoustic classification accuracy than CNN.

➢ Compared to the baseline model (consists of 2 layers × 50 hidden units, 20% dropout), the fused model exhibits 18.4% higher accuracy.

# Results and Discussion -Comparison with CNN only



➤ Although the result of the scene (i.e., beach) using ensemble classifiers (37.96%) is quite poor as compared to CNN (73.14%), the fused model managed to increase the acoustic classification accuracy to 50.92%.

➤ Conversely, the scene (i.e. park) using CNN model is relatively low compared to the ensemble classifiers. The fused model increases it to 60.18%.

➤ The confusion matrix of CNN, ensemble classifiers, and the fused model are shown.

➤ The average acoustic classification accuracy of the fused model is computed as 79.43%.

# Conclusion

The multi-model late fusion system model consisting of the log-Mel spectrogram for convolutional neural network and wavelet time scattering for ensemble of subspace discriminant classifiers was proposed.

Based on the dataset from the TUT Acoustic Scenes, it demonstrated that the fused model gives good acoustic classification accuracy of 79.43%.

The proposed multi-model late fusion system exhibits 18.4% higher acoustic classification accuracy than the baseline model despite relatively low performance in a few scenes such as the beach and library.

Nevertheless, the multi-model late fusion system shows good acoustic classification accuracy for most of the scenes.

# Related Publications (2019-Present)

1) T.K. Chan, C.S. Chin, Multi Branch Convolutional Macaron Net for Sound Event Detection, <u>IEEE Transactions on Audio, Speech and Language Processing</u>, vol. 29, pp. 2972-2985, 2021.

2) T.K. Chan, C.S. Chin, Y. Li, Semi-Supervised NMF-CNN For Sound Event Detection, <u>IEEE Access</u>, vol. 9, pp. 130529-130542, 2021.

3) T.K. Chan, C.S. Chin, Lightweight Convolutional-iConformer For Sound Event Detection, <u>IEEE Transactions on Audio, Speech and Language Processing</u>, 2021, submitted.

4) T. K. Chan and C. S. Chin, A Comprehensive Review of Polyphonic Sound Event Detection, <u>IEEE Access</u>, 8, 103339-103373, 2020.

5) TK Chan, CS Chin, Detecting Sound Events Using Convolutional Macaron Net With Pseudo Strong Labels, <u>IEEE 23rd International Workshop on Multimedia Signal Processing</u>, Tampere, Finland, 6-8 October 2021.

6) CS. Chin, JF. Xiao, Max-Fusion of Random Ensemble Subspace Discriminant with Aggregation of MFCCs and High Scalogram Coefficients for Acoustics Classification, <u>20th IEEE/ACIS International Summer Semi-Virtual Conference on Computer and Information Science, (ICIS 2021-Summer),</u> Shanghai, China, 23-25 Jun 2021.

7) CS Chin, JH Zhang, Wavelet Scattering Transform for Multiclass Support Vector Machines in Audio Devices Classification System, <u>20th IEEE/ASME International Conference on Advanced Intelligent Mechatronics</u>, Aula Conference Centre TU Delft, Delft, The Netherlands, 12-16 Jul 2021.

# Related Publications (2019-Present)

8) XY Kek, CS Chin, Y. Li, An Investigation on Multiscale Normalised Deep Scattering Spectrum with Deep Residual Network for Acoustic Scene Classification, <u>22nd IEEE/ACIS International Fall Virtual Conference on Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing</u>,Taichung, Taiwan, 24-26 November 2021.

9) XY Kek, CS Chin, Y.Li, Deep Scattering Spectrum with Mobile Network for Low Complexity Acoustic Scene Classification, Detection and Classification of Acoustic Scenes and Events 2021(<u>DCASE2021</u>), Technical Report, 2021.

10) CS Chin, JH Zhang, Late Fusion of Convolutional Neural Network with Wavelet-based Ensemble Classifier for Acoustic Scene Classification, <u>Intelligent Systems Conference (IntelliSys) 2021</u>, Springer, Amsterdam, The Netherlands, 2-3 Sept 2021.

11) CS Chin, XY Kek, TK Chan, Wavelet Scattering Based Gated Recurrent Units for Binaural Acoustic Scenes Classification, <u>IEEE International Conference on Internet of Things and Intelligent Applications</u>, Zhenjiang, China, 27-29 Nov, 2020.

12) CS Chin, XY Kek, TK Chan, Scattering Transform of Averaged Data Augmentation for Ensemble Random Subspace Discriminant Classifiers in Audio Recognition, <u>7th International Conference on Advanced Computing and Communication Systems</u>, Coimbatore, India, 19 – 20 March, 2021.

13) Chan TK, Chin CS, Li Y. Non-Negative Matrix Factorization-Convolution Neural Network (NMF-CNN) for Sound Event Detection. IEEE Workshop on Detection and Classification of Acoustic Scenes and Events 2019 (<u>DCASE 2019</u>), New York, USA, 2019.

14) XY. Kek, C. S. Chin, Y Li, Acoustic Scene Classification using Bilinear Pooling on Time-Liked and Frequency-Liked Convolution Neural Network, <u>IEEE Symposium Series on Computational Intelligence</u>, Xiamen, China, 6-9 Dec, 2019.

# Thank you!