



## **Assignment #3 Speech Emotion Recognition (Total 120 Points)**

### **Problem Statement**

Speech is the most natural way of expressing ourselves as humans. It is only natural then to extend this communication medium to computer applications. We define speech emotion recognition (SER) systems as a collection of methodologies that process and classify speech signals to detect the embedded emotions. Below we will show the needed steps to achieve the goal of the assignment.

#### **1. Download the Dataset and Understand the Format (10 Points)**

- We will use CREMA dataset that is available at the following link:  
<https://www.kaggle.com/dmitrybabko/speech-emotion-recognition-en>
- Write your own function that loads an audio and listen to each of the classes you have and plot the waveform of the audio.

#### **2. Create the Feature Space (30 Points)**

We will create two feature spaces from the audio.

- You can work on time domain, or you can work in frequency domain. There are multiple of features that can help improving the model:
  - Zero crossing rate: The rate of sign-changes of the signal during the duration of a particular frame.
  - Energy: The sum of squares of the signal values, normalized by the respective frame length.
- Convert the audio waveform to mel spectrogram and use this as the feature space.

#### **3. Building the Model (40 Points)**

Don't forget to set seed 42 and stratify to true!!

- Split the data into 70% training and validation and 30% testing.
- Use 5% of the training and validation data for validation.
- CNN Model:**  
Build a CNN architecture of your own, a simple example is shown in the figure. For the time domain or frequency domain feature space, the feature space



will be 1 dimensional, therefore in the architecture we will be using 1D convolutions. While in melspectrogram feature space, the audio is represented as an image, therefore we will be using 2D convolutions. **N.B. Do not use the built-in models, but you can implement them from scratch.**

#### 4. Big Picture (20 Points)

Compare between the performance of the learned models (Different features, different learning models) by realizing the following.

- a. Compute the accuracy and F-Score for each model.
- b. Plot the confusion matrices and find the most confusing classes.

#### 5. Bonus (20 Points)

- Best result in the class will get 20 bonus points.
- Second best result will get 10 bonus points.

#### 6. Submission Notes

- a. Work in groups of 3 students.
- b. **[20 Points]** You are required to submit a clear and detailed report [in PDF format] illustrating every step in the assignment.

#### 7. References

Akçay, M. B., & Oğuz, K. (2020). Speech emotion recognition: Emotional models, databases, features, preprocessing methods, supporting modalities, and classifiers. *Speech Communication*, 116, 56-76.

[\[https://www.sciencedirect.com/science/article/abs/pii/S0167639319302262\]](https://www.sciencedirect.com/science/article/abs/pii/S0167639319302262)

*Good Luck*