# Video Call Application (WebRTC) with audio driven talking head generation

This the first iteration, where we are focused on installing python 3.X, setting up conda environment, installing required libraries for the project execution. Apart from that we experimented and learned about WebRTC + Streamlit. The decision is made on the ML model to implement the chosen audio driven approach.

## Installing python 3.X:

All the models were implemented using Python 3.6, hence as to implement the same models we decided on installing python 3.6

 The list of libraries needed to implement:

- ffmpeg-python: This is used to trim, concatenate and overlay the audio files that we record.its command-line interface.
- opencv-python: This can be utilized when dealing with image processing. In our case this library comes handy when we overlap caller's face with talking head image.
- face_alignment: It is used to detect 2D and 3D face alignment.
- scikit-learn: library for machine learning in Python. This provides us a selection of efficient tools for machine learning where we can implement regression, clustering and dimensionality reduction via a consistence interface.
- pydub: This library is used to handle the .wav files which are generated after recording the audio.
- pynormalize: Using this we save the metadata of the original audio file,which can be used for processing.
- soundfile: Read and write operations of audio file are performed using this module.
- librosa: This is used to normalize the audio file and to extract features of audio file.
- pysptk: It's a python wrapper library for speech signal processing.
- pyworld: This is used for speech breakdown in our project.
- resemblyzer: It is utilized to implement deep learning part of our project.

## Setting up conda environment:

We created a new environment in conda and installed the above-mentioned dependencies(libraries) to further execute ML model. There are basically two steps involved one creation and other activation.

Makeittalk is the new environment that we tried to create.

**Fig: Makeittalk environment is created**



**Fig: Environment activation**

## WebRTC+Streamlit:

We are using Streamlit to develop the UI of our project and WebRTC is used to provide the real-time video communication. The two libraries for this purpose are streamlit and stream-webrtc.

The commands to install python libraries required to have this setup are:

- Pip install Streamlit
- Pip install streamlit-webrtc

After studying about streamlit and webrtc protocol, we locally tried multiple bits and pieces of code to get hold of the structure and development environment. Finally we were able to develop a streamlit webpage that consists of two windows which are enabled with camera and microphone, and these can in future be modified to create a real-time call scenario. The webpage also consists of an option where user can chose the talking head.

Instructions to run the application:

- Navigate to the folder where the application resides.
- To run the application, use the syntax: "streamlit run filename.py"

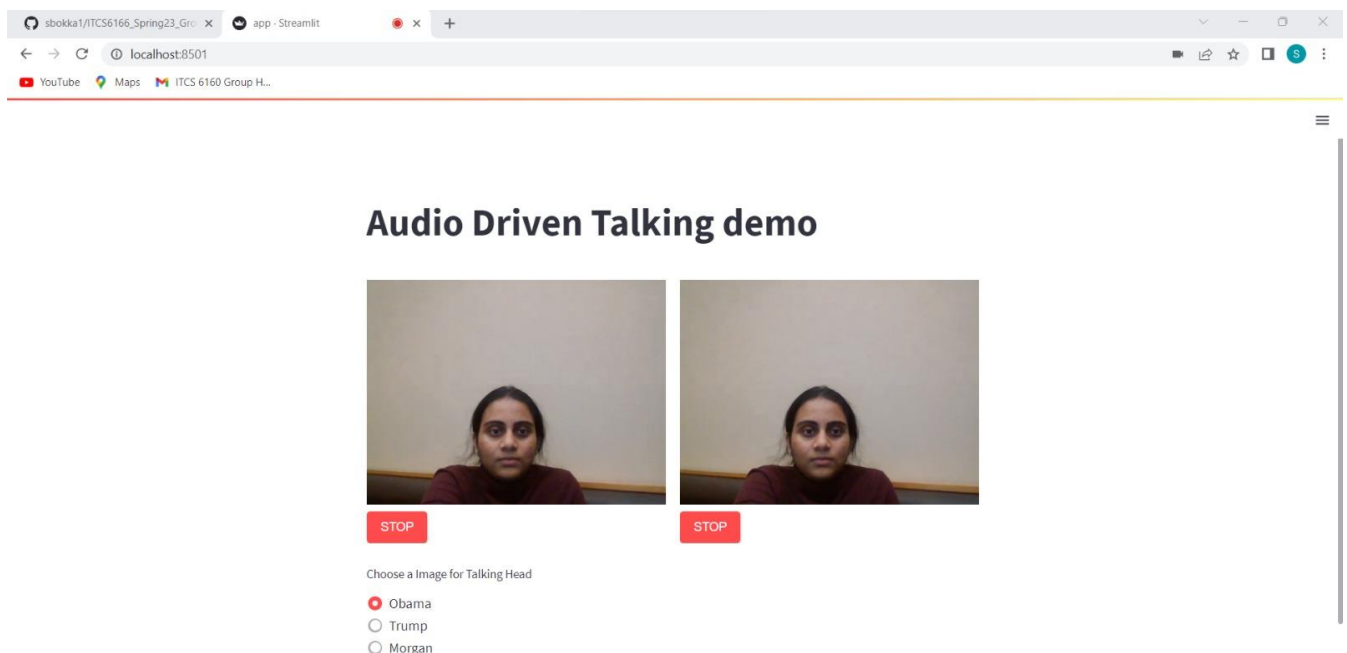**Fig: Screenshot representing the run instructions and run results.**



**Fig: Streamlit page**

The above shown image is our initial attempt towards the UI of the project and it will be further developed in coming iterations. The code file corresponding to the above page is attached in the GitHub.

## Decision on the Model:

After closely examining both models i.e Makeittalk and livespeech portraits. We finally decided on proceeding with Makeittalk model. For the following reasons: Both models are almost similar. Makeittalk model is additionally using libraries like face_alignment to get the precise dimensions of the user face in-order to ensure perfect overlap. The makeittalk model also utilizes libraries that can better process the audio files when compared to libraries in livespeech portraits.

### aioRTC Experiment:

We used aiortc repo: https://github.com/aiortc/aiortc to learn to understand how WebRTC works and tinker with its internals. Since the repository is quite old, newer version of python didn't work. So I figured python3.9 worked well.

We used the server example from the above repository. This example illustrates establishing audio, video and a data channel with a browser. It also performs some image processing on the video frames using OpenCV.
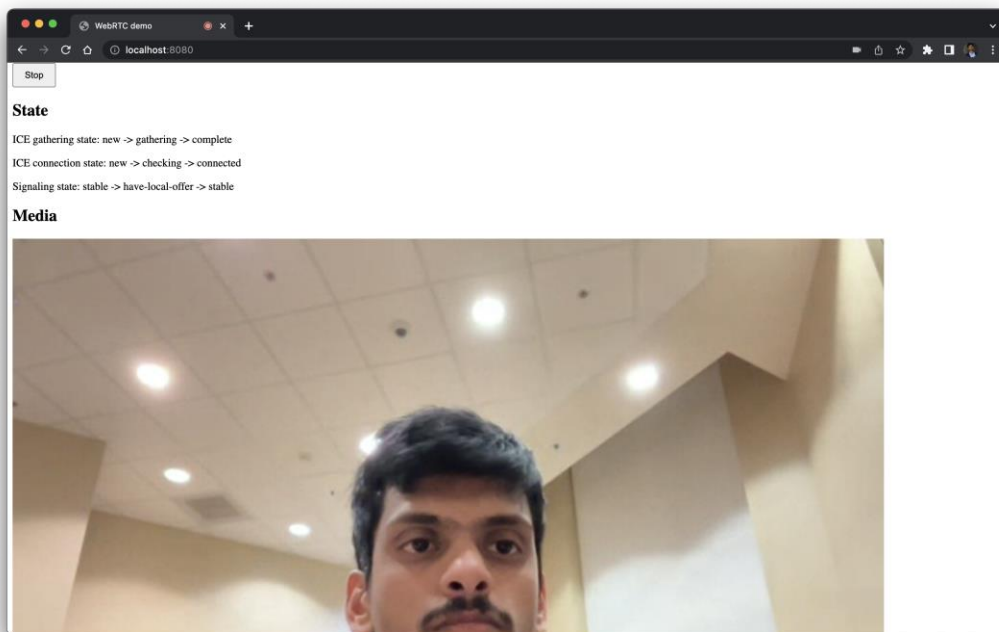


**Fig: Server example**

We completed all the planned tasks of iteration-1 and besides that we tried to test the ML models, there were few bottle necks in testing. We are positive about overcoming them in further iterations.

The test results and the codes of implementation of webrtc , streamlit and ML model executions can be found in the respective individual folders.