



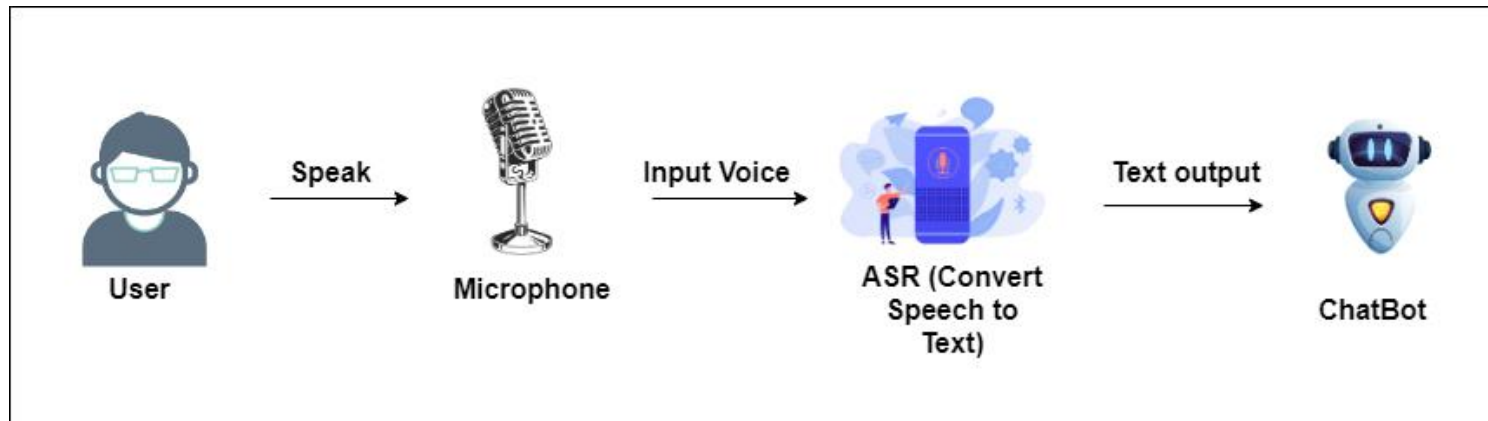
# IT20237554 | RATHNAWEERA R.P.W.G

B.Sc. (Hons) Degree in Information Technology Specialized in Data Science

# Specific and Sub Objectives

## Main Objective

Development of an Automatic Speech Recognition system for feeding queries to Chatbot



# Specific and Sub Objectives

- Sub Objective

- Dataset creation

- Developed custom data scrappers to gather data for the recommendation system

- Google Maps Scraping

- Developed a Google Maps Scraper to scrape the computer repair centers for recommendation

- Reviews automation

- Suggest the video link and a summary of a YouTube review regarding a recommended device to the user

# ASR – CTC Approach(Based on Deep Speech 2)

- Connectionist temporal classification used where

It's a model with RNN and CNN and trained on LJ speech dataset

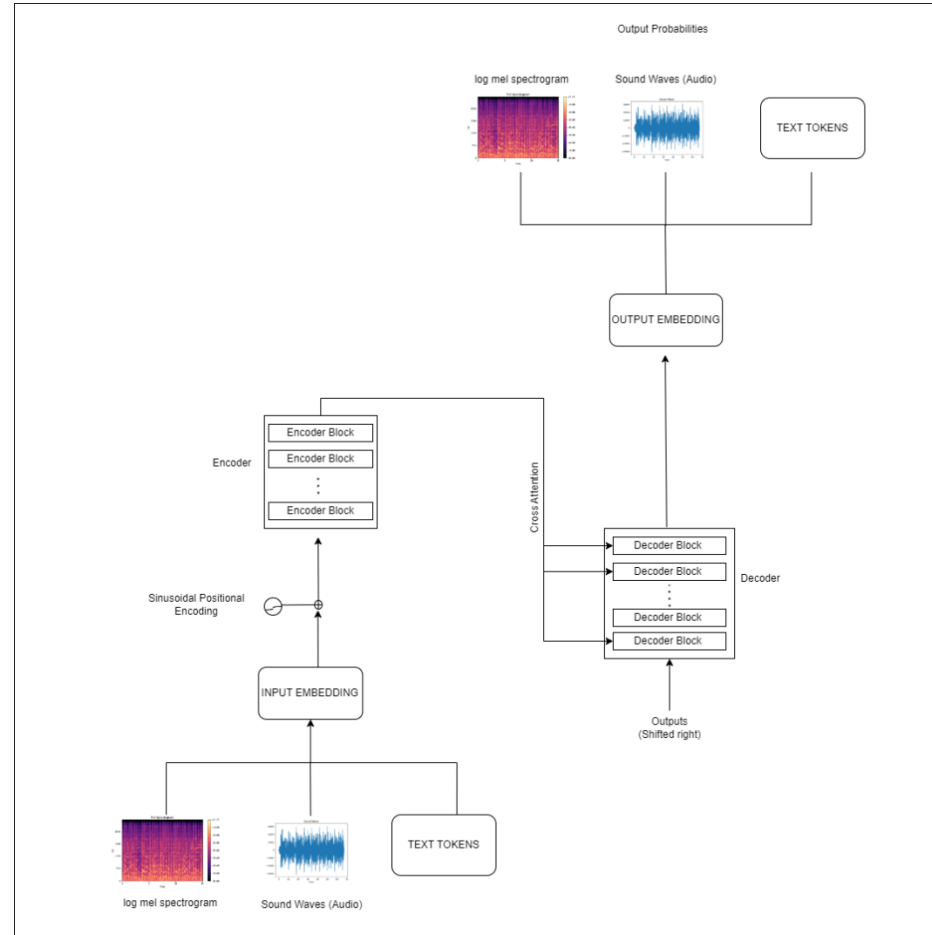
- Downsides of this model was it requires high computational power and less accurate even after running 20 epochs

```
1 def build_model(input_dim, output_dim, rnn_layers=5, rnn_units=128):
2     # Model's input
3     input_spectrogram = layers.Input((None, input_dim), name="input")
4     # Expand the dimension to use 2D CNN
5     x = layers.Reshape((-1, input_dim, 1), name="expand_dim")(input_spectrogram)
6     # Convolution layer 1
7     x = layers.Conv2D(
8         filters=32,
9         kernel_size=[11, 41],
10        strides=[2, 2],
11        padding="same",
12        use_bias=False,
13        name="conv_1",
14    )(x)
15    x = layers.BatchNormalization(name="conv_1_bn")(x)
16    x = layers.ReLU(name="conv_1_relu")(x)
17    # Convolution layer 2
18    x = layers.Conv2D(
19        filters=32,
20        kernel_size=[11, 21],
21        strides=[1, 2],
22        padding="same",
23        use_bias=False,
24        name="conv_2",
25    )(x)
26    x = layers.BatchNormalization(name="conv_2_bn")(x)
27    x = layers.ReLU(name="conv_2_relu")(x)
28    # Reshape the resulted volume to feed the RNNs layers
29    x = layers.Reshape((-1, x.shape[-2] * x.shape[-1]))(x)
30    # RNN layers
31    for i in range(1, rnn_layers + 1):
32        recurrent = layers.GRU(
33            units=rnn_units,
34            activation="tanh",
35            recurrent_activation="sigmoid",
36            use_bias=True,
37            return_sequences=True,
38            reset_after=True,
39            name=f"gru_{i}",
40        )
41    x = layers.Bidirectional(
42        recurrent, name=f"bidirectional_{i}", merge_mode="concat"
43    )(x)
44    if i < rnn_layers:
45        x = layers.Dropout(rate=0.5)(x)
46    # Dense layer
47    x = layers.Dense(units=rnn_units * 2, name="dense_1")(x)
48    x = layers.ReLU(name="dense_1_relu")(x)
49    x = layers.Dropout(rate=0.5)(x)
50    # Classification layer
51    output = layers.Dense(units=output_dim + 1, activation="softmax")(x)
52    # Model
53    model = keras.Model(input_spectrogram, output, name="DeepSpeech_2")
54    # Optimizer
55    opt = keras.optimizers.Adam(learning_rate=1e-4)
56    # Compile the model and return
57    model.compile(optimizer=opt, loss=CTCLoss)
58    return model
59
60
61
62 # Get the model
63 model = build_model(
64     input_dim=fft_length // 2 + 1,
65     output_dim=char_to_num.vocabulary_size(),
66     rnn_units=512,
67 )
```

```
1/1 [=====] - 1s 685ms/step
1/1 [=====] - 1s 728ms/step
1/1 [=====] - 1s 644ms/step
1/1 [=====] - 1s 750ms/step
1/1 [=====] - 1s 680ms/step
1/1 [=====] - 1s 675ms/step
1/1 [=====] - 1s 700ms/step
1/1 [=====] - 1s 739ms/step
1/1 [=====] - 1s 607ms/step
1/1 [=====] - 1s 621ms/step
-----
Word Error Rate: 0.2681
-----
Target : and you have to establish adequate circulation
Prediction: and you havet establish adequate circulation
-----
Target : and most of the records which have come from his time speak chiefly of his deeds of piety
Prediction: and most of the records which ove come from his time up chiefly of his deds of peity
-----
369/369 [=====] - 1322s 4s/step - loss: 38.1977 - val_loss: 44.8476
Epoch 20/20
1/1 [=====] - 1s 719ms/step loss: 28.91
1/1 [=====] - 1s 724ms/step
1/1 [=====] - 1s 739ms/step
1/1 [=====] - 1s 648ms/step
1/1 [=====] - 1s 645ms/step
1/1 [=====] - 1s 665ms/step
1/1 [=====] - 1s 713ms/step
1/1 [=====] - 1s 672ms/step
1/1 [=====] - 1s 671ms/step
1/1 [=====] - 1s 724ms/step
1/1 [=====] - 1s 685ms/step
1/1 [=====] - 1s 672ms/step
1/1 [=====] - 1s 650ms/step
1/1 [=====] - 1s 669ms/step
1/1 [=====] - 1s 686ms/step
1/1 [=====] - 1s 682ms/step
1/1 [=====] - 1s 663ms/step
1/1 [=====] - 1s 715ms/step
1/1 [=====] - 1s 624ms/step
1/1 [=====] - 1s 734ms/step
1/1 [=====] - 1s 639ms/step
1/1 [=====] - 1s 658ms/step
1/1 [=====] - 1s 670ms/step
1/1 [=====] - 1s 637ms/step
1/1 [=====] - 1s 655ms/step
1/1 [=====] - 1s 740ms/step
1/1 [=====] - 1s 686ms/step
1/1 [=====] - 1s 706ms/step
1/1 [=====] - 1s 732ms/step
1/1 [=====] - 1s 677ms/step
1/1 [=====] - 1s 632ms/step
1/1 [=====] - 1s 683ms/step
1/1 [=====] - 1s 617ms/step
1/1 [=====] - 1s 744ms/step
1/1 [=====] - 1s 678ms/step
1/1 [=====] - 1s 679ms/step
1/1 [=====] - 1s 670ms/step
1/1 [=====] - 1s 735ms/step
1/1 [=====] - 1s 686ms/step
1/1 [=====] - 1s 628ms/step
-----
Word Error Rate: 0.2689
-----
Target : all the misdemeanants whatever their offense were lodged in this chapel ward
Prediction: all the misdemeanants whatever ther offense were lodge in this chapel ward
-----
Target : to inquire into and report upon the several jails and houses of correction in the counties cities and corporate towns within england and wales
Prediction: to equire into an report upon the several jails and houses of cecrection in the counties cities and corporate bounds within england and wall
-----
369/369 [=====] - 1376s 4s/step - loss: 28.9184 - val_loss: 45.6976
```

# ASR – Seq2Seq Approach

- Used Whisper model which was released in end of September 2022 where it has been trained for 680,000 hours of audio data and fined tuned this whisper model for my requirements with “Commen Accent” dataset.



# Best Approach

CTC Approach	Seq2Seq Approach
High usage of computational power	Comparatively low usage of computational power
High rate of spelling errors	Low rate of Spelling errors
Only use Encoder	Use both decoder and encoder ( Can enhance the ASR for multiple languages)
Word Accuracy is Low	Word Accuracy is High
Difficult to handle High robustness, Noise inputs	High robustness, Noise inputs can be handled

Due to the above comparison, I had selected the seq2seq architecture model to fit for the chatbot

# Model Results

## CTC – Model based on Deep Speech2

WER = 0.26 (26%)

WAcc =  $1 - 0.26 = 0.74$  (74 %)



```
4/1 [=====] - 45 07ms/step
1/1 [=====] - 1s 632ms/step
1/1 [=====] - 1s 683ms/step
1/1 [=====] - 1s 617ms/step
1/1 [=====] - 1s 744ms/step
1/1 [=====] - 1s 678ms/step
1/1 [=====] - 1s 679ms/step
1/1 [=====] - 1s 670ms/step
1/1 [=====] - 1s 735ms/step
1/1 [=====] - 1s 606ms/step
1/1 [=====] - 1s 628ms/step
-----
Word Error Rate: 0.2609
-----
Target   : all the misdemeanants whatever their offense were lodged in this chapel ward
Prediction: all the misdemeanants whatever ther offense were lodge in this chapel ward
-----
Target   : to inquire into and report upon the several jails and houses of correction in the counties cities and corporate towns within england and wales
Prediction: to equire into an report upon the several jails and houses of ceration in the counties cities and corporate tounds wthan england and wail
-----
369/369 [=====] - 1376s 4s/step - loss: 28.9184 - val_loss: 45.6976
```

## Seq2Seq Whisper finetuned model

WER = 0.13 (13%)

Wacc =  $1 - 0.13 = 0.87$  (87 %)



2500/2500 3:35:20, Epoch 5/6]

Step	Training Loss	Validation Loss	Wer Ortho	Wer
500	0.101200	0.321511	16.378369	11.594108
1000	0.034500	0.348289	16.649578	11.845021
1500	0.018000	0.382895	17.162229	12.470673
2000	0.007500	0.406941	17.866711	13.011601
2500	0.005900	0.423361	17.922937	13.060480

TrainOutput(global\_step=2500, training\_loss=0.13272297571897507, metrics={'train\_runtime': 12935.8975, 'train\_samples\_per\_second': 3.092, 'train\_steps\_per\_second': 0.193, 'total\_flos': 1.15318725967872e+19, 'train\_loss': 0.13272297571897507, 'epoch': 5.71})

# Seq2Seq - Model Testing

- 2 models were built for the Seq2Seq approach.

<b>Model 1</b> <b>Model with 500 steps up to 8 epochs</b>	<b>Model 2</b> <b>Model with 2500 steps up to 6 epochs</b>
Training loss - 0.0019 Epoch - 7.94 Step - 500 Validation Loss - 0.6059 WER Ortho - 22.9616 WER - 19.0970	Training loss - 0.0059 Epoch - 5.71 Step - 2500 Validation Loss - 0.4234 WER Ortho - 17.9229 WER - 13.0605

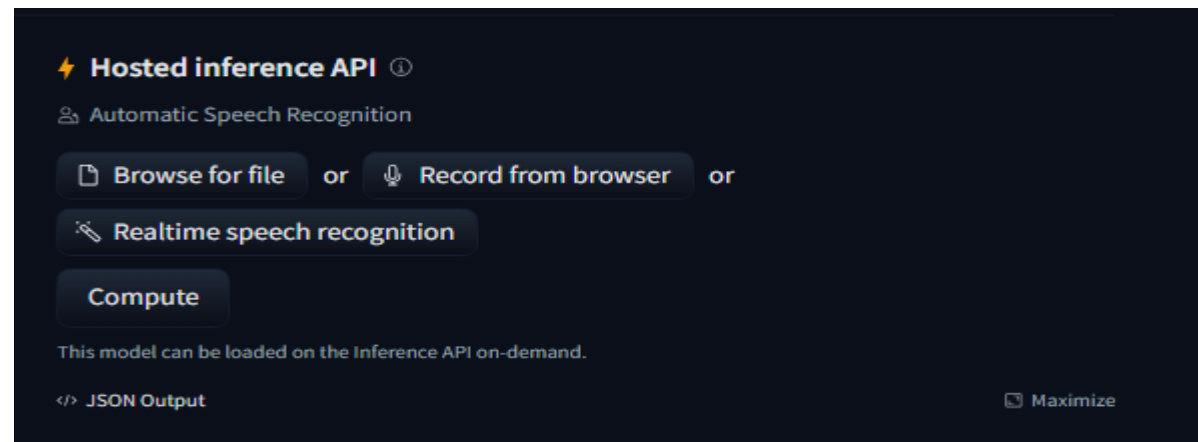


# Seq2Seq - Model Testing

- Created custom Audio clips to test the model.

Audio	Transcript
	I need a laptop to do my personal things and it should have a RTX 3060 VGA and 12GB RAM so I can do my stuff without any lag

- Since both models were deployed in the Hugging Face API, tested those models by uploading the transcript for the API interface.



# Seq2Seq Approach Model Results

## When uploading an audio clip

- Model 1 Results

Hosted inference API ⓘ

Automatic Speech Recognition

Browse for file or Record from browser or

Realtime speech recognition

Test Recording.wav

0:00 / 0:14

Compute

Computation time on Intel Xeon 3rd Gen Scalable cpu: 3.080 s

I need a laptop to do my personal things and it should have a RTX 3060 VGA and a 12GB RAM to do the things smoothly without having any lagging issues.

</> JSON Output Maximize

- Model 2 Results

Hosted inference API ⓘ

Automatic Speech Recognition

Browse for file or Record from browser or

Realtime speech recognition

Test Recording.wav

0:00 / 0:14

Compute

Computation time on Intel Xeon 3rd Gen Scalable cpu: 2.815 s

I need a laptop to do my personal things, and it should have a RTX 3060 VGA and a 12GB RAM to do the things smoothly without having any lagging issues.

</> JSON Output Maximize

# Seq2Seq Approach Model Results

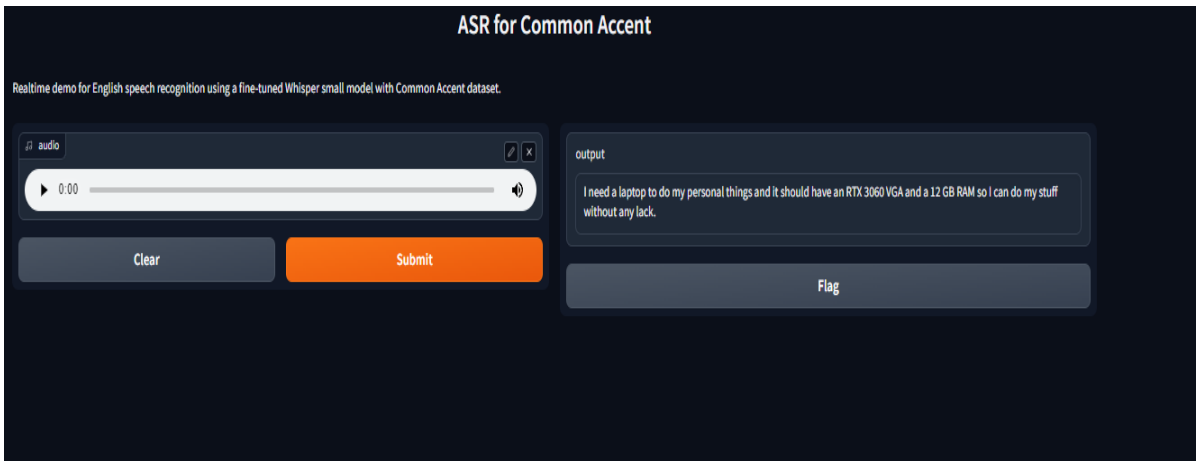
## For Real-time Voice recognition

- Loading model 1

```
# Load model directly
from transformers import AutoProcessor, AutoModelForSpeechSeq2Seq
#wishwa98/CommonAccent_TuneMore
processor = AutoProcessor.from_pretrained("wishwa98/CommonAccent_TuneMore")
model = AutoModelForSpeechSeq2Seq.from_pretrained("wishwa98/CommonAccent_TuneMore")

Downloading (...)processor_config.json: 100% 339/339 [00:00<00:00, 9.27kB/s]
Downloading (...)tokenizer_config.json: 100% 805/805 [00:00<00:00, 23.4kB/s]
```

- Model 1 Results



- Loading model 2

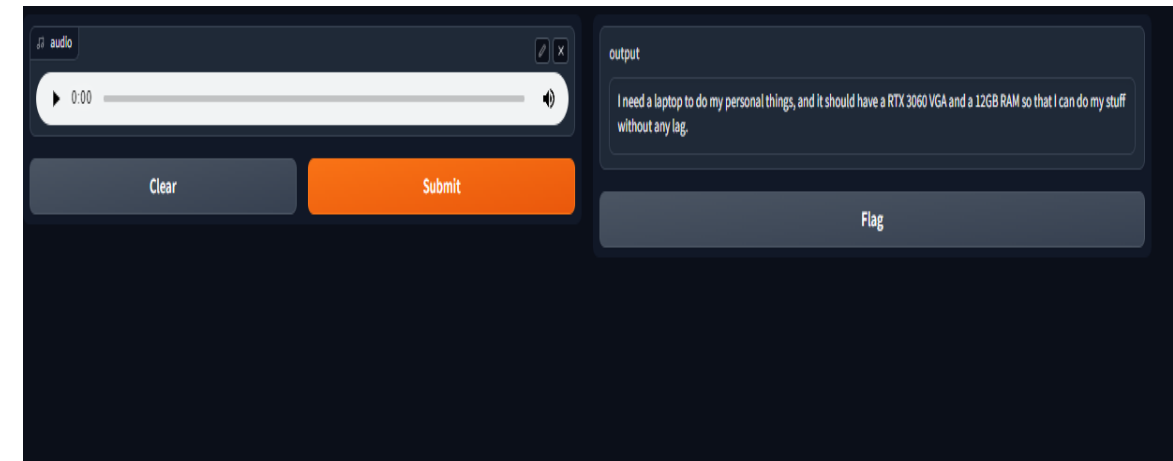
```
Model 2

# Load model directly
from transformers import AutoProcessor, AutoModelForSpeechSeq2Seq

processor = AutoProcessor.from_pretrained("wishwa98/ASRForCommonVoice")
model = AutoModelForSpeechSeq2Seq.from_pretrained("wishwa98/ASRForCommonVoice")

Downloading (...)processor_config.json: 100% 339/339 [00:00<00:00, 14.0kB/s]
Downloading (...)tokenizer_config.json: 100% 805/805 [00:00<00:00, 50.1kB/s]
```

- Model 2 Results



# Evaluation Metrics

There are different word error categories as Substitution, Insertion Deletion in both word level and character level Here we used the Word level, which calculates the substitutions insertions, and deletions on a word level These errors are annotated on a word-by-word basis,

S=substitute

D=Delete

I=Insertion

N= Total number of words

$$WER = \frac{S + I + D}{N}$$

Word Error Rate is an error matrix that can be convert into accuracy matrix

$$W_{Acc} = 1 - WER$$

# Seq2Seq Approach Model Results

- Expected Text – “ I need a laptop to do my personal things, and it should have a RTX 3060 VGA and a 12GB RAM so **that** I can do my stuff without any **lag** “

Predicted Text (Model 1)	Predicted Text (Model 2)
“I need a laptop to do my personal things and it should have an RTX 3060 VGA and a 12GB RAM so I can do my stuff without any <b>lack</b> .”	“I need a laptop to do my personal things, and it should have an RTX 3060 VGA and a 12GB RAM so <b>that</b> I can do my stuff without any <b>lag</b> .”

- Can clearly see in the model 1 word “That” has been deleted and the word “lag” has been substituted with word “lack”
- Due to the model evaluation results and the Realtime testing results , Selected the model 2 to implement In the chatbot interface.