# Automated Speech Recognition

1st Prayash Das
*Department of Computer Science and Engineering*
*Amity University Kolkata*
Kolkata, India
prayash.das@student.amity.edu

2nd Pronaya Bhattacharya
*Department of Computer Science and Engineering*
*Amity University Kolkata*
Kolkata, India
pbhattacharya@kol.amity.edu

*Abstract*—We all are habituated with the use of virtual assistants nowadays. The use of personal voice controlled assistants in mobile phones such as Siri of Apple, Google Assistant of Google, Cortana of Microsoft, Alexa of Amazon and many more has made our lives much easier in managing our daily tasks. The use of features like using hands-free call, sending text messages through voice dictation, asking weather without opening the app, setting up reminders and alarms has made our lives simpler and hassle-free and helped us to organise our daily tasks in proper arrangement. But with the benefits of these virtual assistants also comes the cons. I will discuss one of the major drawbacks of these AI assistants, that is the speech to text feature. Whenever it comes to the accuracy of this feature, in comparison to typed transcripts by human they are lagging significantly behind. No matter how clearly and fluently we dictate our voices to these voice based assistants, they still are unable to understand it accurately and make grammatical mistakes while conversion from speech to text. Even the tech giants in the AI Industry have an error rate of more than 10 percentage. A number of factors can impact word error rate, such as pronunciation, accent, pitch, volume, and background noise. The main goal of these Speech to Text Software Systems is to reach the error level as that of two humans speaking to each other.

*Index Terms*—Automated Speech Recognition, Artificial Intelligence, Speech to Text Software Systems, Word Error Rate, Algorithms, Emission Probability, Language Weighting, Speaker Labelling, Acoustics Training, Profanity Filtering, Hidden Markov Model, Borel Sets, N-grams, Speaker Dialisation, Language Model Score, Cross Entropy, Bits-per-Character

## I. INTRODUCTION

Speech to text is a speech recognition software that enables the recognition and translation of spoken language into text through computational languages. It is also known as speech recognition or computer speech recognition.It is an attribute of the Natural Language Processing which is a branch of Artificial Intelligence. With the introduction of virtual assistants, Speech Recognition became one of the fundamental attributes of Virtual Assistants. Speech Recognition became an integral part of our life, assisting us in accomplishing our tasks without the need to operate electronic gadgets/devices physically. In corporate world as well, these became an essential part handling and completing client's requests/queries. Speech Recognition plays an important roles in Automotive, Healthcare and Security sector as well. But these virtual

assistants cannot completely decipher our speech and makes grammatical errors and even generate meaningless words. No matter how much fluently we dictate our sentences and lines fluently and correctly( w.r.t grammar) to these voice based assistants makes error while deciphering our speech. In many cases, these virtual assistants are completely unable to interpret the client's requests and unable to complete the client's demands completely. The accuracy of human typed transcripts is nearly 99If we are to completely rely on these virtual assistants for communication purposes through voice, they will give out some error to the receiver who will be unable to interpret the meaning completely. The industry standard to measure Speech-to-Text accuracy is Word Error Rate (WER). WER counts the number of incorrect words identified during recognition, then divides by the total number of words provided in correct transcript. WER is calculated by the formula:- S + I + D)/N

S stands for substitutions I stands for insertions D stands for deletions

N stands for number of words in the references. The statistics below is shown of a survey conducted by Lionel Sujay Vailshery, Research expert covering the consumer electronics industry, in comparing speech-to-text error rate among leading companies worldwide in 2021. Based on these surveys, Rev.ai had the least error rate of 14.22

### A. Methodologies Used

a) Language Weighting- Weightage are given to those words that are often spoke frequently to improve precision, by going beyond terms that are already in the base vocabulary.

b) Speaker Labelling- These are tags/words used for identifying a person speaking. These labels can be either a speaker's name, role or other identifying attributes.

c) Acoustics Training- Training the system to an acoustic environment ( like ambient noise in a call centre) and speaker characteristics(voice, pitch).

d) Profanity filtering- Training the system to filter out profanity words and sanitise the speech output.

## II. MAJOR ALGORITHMS IMPLEMENTED

Various algorithms are implemented for the speech to text feature and to improve transcription accuracy. I have enlisted some of the major algorithms used:-

(a) Hidden Markov Model(HMM)- The Hidden Markov Model is a statistical model first proposed by Baum L.E. (Baum and Petrie) in 1966. HMM is a branch of Applied Probability and Statistics area. It is based on the Markov chain model, Markov chain is useful for observable events, while Hidden Markov Model is useful for hidden events This allows to incorporate hidden events such as part-of-speech labels into a probabilistic model.This arranges the labels in sequential manner i.e., words, syllables, sentences. This helps to create a mapping which helps us to determine correct labelling sequence.

Let Xt and Yt be discontinuous stochastic time processes and n 1

The pair (Xt, Yt) is a Hidden Markov model if:-

• The behavioural pattern of Xt is hidden( not observable) and Xt is a Markov process.

• Emission Probability of (Yt A X1 = x1,...,Xt = xt) is similar to Emission Probability of (Yt A Xt = xt),

For every n 1, x1,...,Xt, and every Borel set A.

Let Xa and Ya be continuous stochastic time processes. The pair (Xa, Ya) is Hidden Markov model if:-

• Xa is a Markov model and the behaviour of Xa is hidden.

• Emission Probability of(Ya A Xa Ba(t t0) is similar to Emission Probability of (Ya0 A Xa0 Ba0),

For every a0, every Borel set A, and every family of Borel setsBa.

( Borel set is a set that can be from open spaces through the operations of countable union, countable intersection, and complement relativity).

(b) N-grams- These are the simplest language models used in Speech Recognition Technique. An N-gram is a sequence of N- words. Example, Ram is a boy is a 4-gram. Recognition and accuracy is improved by grammar and determining the probability of certain words sequences.

(c) Neural Networks-These are the building blocks in Deep Learning Techniques. Neural Networks process training data with the help of layers of nodes.

Node is made up of inputs, weights, a threshold and an output. If the value of output exceeds the value of threshold, it activates that particular node passing the information to the next

layer. They are based on Supervised Learning. Although neural networks are more accurate and can be trained with more data, they are slower to train in comparison to conventional language models.

(d) Speaker Dialisation(SD)- SD identifies and segments speech through speaker identity. These helps programs to distinguish individuals in a conversation. Each speaker is separated by their unique audio characteristics and their utterances are bucketed together.

### A. Abbreviations and Acronyms

WER, BPC

### B. Equations

$$WER = S + I + D/N \qquad (1)$$

### C. Authors and Affiliations

**Prayash Das, Department of Computer Science and Engineering, Amity University Kolkata, Dr. Pronaya Bhattacharya**.
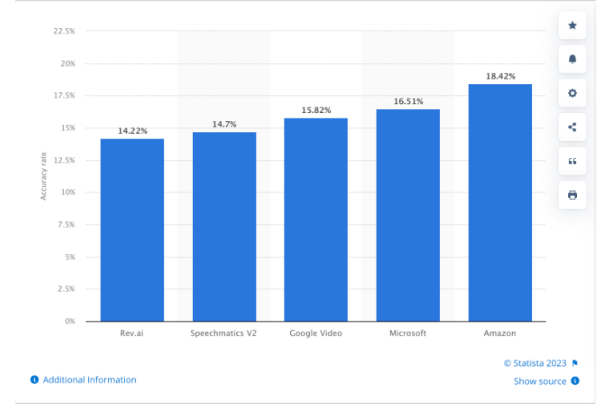
### D. Figures and Tables



Fig. 1. Speech to Text Transcript Error Rate Among Leading Companies Worldwide 2021

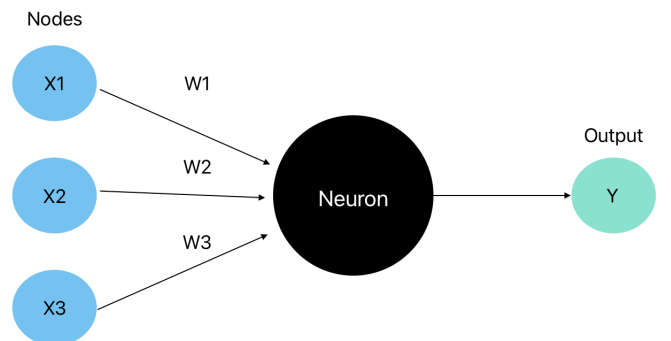### ACKNOWLEDGMENT

### REFERENCES

[**?**]. [**?**].

Fig. 2. Workflow of Neuron

## REFERENCES

[1] Frederick Jelinek, "Statistical Methods for Speech Recognition,", 2022.

[2] Kamil Ekštein, František Pártl, Miloslav Konopík, " Text, Speech, and Dialogue, " 24th International Conference, TSD 2021, Olomouc, Czech Republic, September 6–9, 2021, Proceedings

[3] Prof Philip M. Parker Ph.D, " The 2022 Report on Mobile Handset Speech Recognition ", World Market Segmentation by City, 2021.