| | **Mandatory summer internship report** | **Code : DO-PFE-01** |
|---|---|---|
| | | **Indice de révision : 00** |
| | | **Edition : 07/2022** |

**Mandatory Summer Internship Report**

**Major: Industrial Computing & Automation**
**Grade: 4$^{th}$ year**

**Title of the Project :**

# C2S Speech-to-Speech AI Assistant Project Report

Implemented by: **Amine KAROUI**

Host Entreprise:

INFLUENCE
CONSULTING

**Scholar year : 2022 /2023**

| | **Mandatory summer internship report** | Code : DO-PFE-01 |
|---|---|---|
| | | Indice de révision : 00 |
| | | Edition : 07/2022 |

**Mandatory Summer Internship Report**

**Major: Industrial Computing & Automation**
**Grade: 4th year**

<u>**Title:**</u>

# C2S Speech-to-Speech AI Assistant Project Report

Implemented by: **Amine KAROUI**

Host Entreprise:



| *Responsable à l'entreprise :* | *Avis de la commission des stages* |
|---|---|
| **M. Fakhreddine BELGAIED** **M. Arbi SOUSSI** | |

# Table of contents

# Acknowledgment

This report reflects the outcomes of a full month dedicated to my internship at Influence Consulting, and I am profoundly grateful for this enriching opportunity. Upon the conclusion of this period, it is of utmost importance to commence this internship report by expressing my sincere gratitude to all those who have provided their invaluable contributions. I would like to extend my gratitude, in particular, to the entire team at Influence Consulting for their warm welcome, the time they generously dedicated to sharing their knowledge, and their invaluable support throughout this experience. Each of your team members not only shared essential insights but also provided relevant insights and constructive corrections, significantly contributing to the enrichment of this work. Your commitment to my professional development and the quality of my learning has been exemplary, and I am honored to have benefited from your expertise. Your guidance and cooperation have greatly contributed to the completion of this report. Once again, I thank you from the bottom of my heart for this enriching experience and for being exceptional mentors. Your support will remain a source of inspiration for my future projects.

# I.    Introduction

The project embarked upon during my internship with Influence Consulting is a testament to the potential of technology in elevating human-machine interaction. This section serves as the gateway to the comprehensive journey of creating a Proof of Concept (POC) for a Speech-to-Speech AI Assistant. It outlines the project's objectives, significance, and the critical components that form the nucleus of this innovative work

At its core, this project aimed to conceive, design, and realize a transformative Proof of Concept (POC) for a Speech-to-Speech AI Assistant. This innovative system harnesses the power of Automatic Speech Recognition (ASR), Text-to-Speech (TTS), and sophisticated dialogue management to enable communication that's nuanced, efficient, and adaptive.

This endeavor shines a light on the ever-evolving relationship between technology and human interaction. In a world where our interactions with machines are becoming more integral to daily life, it's clear that the confines of traditional communication needed a paradigm shift. The C2S Speech-to-Speech AI Assistant steps into this arena, with a vision to redefine the landscape of communication, offering personalization and precision.

We live in an era where technology is intricately interwoven into the fabric of our lives. This project represents a significant milestone, poised to reshape the quality of human-machine interactions. By amalgamating ASR, TTS, and dialogue management, we're at the cusp of a new era in how we engage with technology.

This project began as a Proof of Concept (POC) with clear objectives:

- To prove the technical feasibility of constructing a Speech-to-Speech AI Assistant.
- To demonstrate the potential of such technology in enhancing human-machine interactions.
- To lay the groundwork for future innovations.
- To test the adaptability of the AI assistant across various domains.

## II.   About Influence Consulting

Influence Consulting is a dynamic and innovative firm that specializes in corporate and digital communication. With a team of seasoned experts boasting over 15 years of experience, the company has successfully delivered high-quality services across various domains.

The company has a talented, multi-disciplinary workforce with profound expertise in the field of digital communication technology. This includes a deep understanding of digital communication networks, internet platforms, and the ever-evolving landscape of internet services.
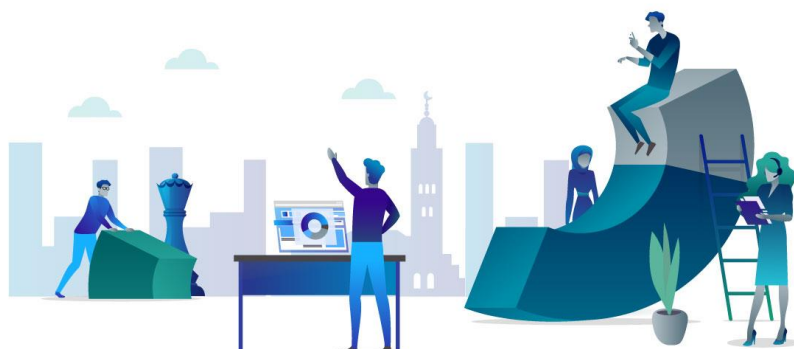
Influence Consulting specializes in three core areas:

**Digital Communication:** The company has a strong track record of strategically utilizing digital communication to help businesses effectively engage with their target audiences.

**Social Media Management:** Their expertise in social media management ensures that brands maintain impactful online presences aligned with the latest trends and technologies.

**Internet Services Development:** With a forward-thinking approach, they offer insights and strategies to drive digital transformation.

What sets Influence Consulting apart is its extensive industry experience, offering valuable insights and knowledge to its clients. The multi-disciplinary approach of their team allows them to tackle projects from various angles. They are pioneers in the field and consistently stay ahead of the curve, providing cutting-edge solutions. At Influence Consulting, marketing research and studies are employed as effective methods for determining the appropriate Marketing Mix, essential for the success of political products. To assure this, Influence Consulting is deeply interested in leveraging artificial intelligence to align its goals with the dynamic landscape of corporate and online communication.

# Internship Journal

| TÂCHE | ATTRIBUÉE À | AVANCEMENT | DÉBUT | FIN |
|---|---|---|---|---|
| **Semaine 1** | | 100% | | |
| Setup of the development environment, including required software and tools. | | 100% | 1/7/23 | 3/7/23 |
| Began exploring the existing AI framework and tools. | | 100% | 4/7/23 | 6/7/23 |
| **Semaine 2** | | 100% | | |
| In-depth research on ASR (Automatic Speech Recognition) and TTS (Text-to-Speech) technologies. | | 100% | 10/7/23 | 11/7/23 |
| Started hands-on implementation of ASR integration with Rasa. | | 100% | 11/7/23 | 12/7/23 |
| Explored different speech datasets for training and evaluation. | | 100% | 13/7/23 | 14/7/23 |
| **Semaine 3** | | 100% | | |
| Conducted benchmarking to assess the performance of different ASR options | | 100% | 17/7/23 | 19/7/23 |
| Discussions on creating a custom dialogue management flow | | 100% | 19/7/23 | 20/7/23 |
| Explored possibilities of integrating Mozilla TTS for Text-to-Speech | | 100% | 20/7/23 | 21/7/23 |
| **Semaine 4** | | 100% | | |
| Worked on refining the dialogue management using Rasa Core | | 100% | 24/7/23 | 27/7/23 |
| Completed the initial ASR and TTS integration within the dialogue management | | 100% | 27/7/23 | 28/7/23 |
| Documented the progress and gathered feedback from the team. | | 100% | 31/7/23 | 1/8/23 |
| **Semaine 5** | | 100% | | |
| Developed a Python application for speech-to-speech AI assistant | | 100% | 2/8/23 | 3/8/23 |
| Discussed and implemented a feature for menu item verification. | | 100% | 3/8/23 | 5/8/23 |
| Worked on an initial graphical user interface (GUI) for the desktop application using Tkinter. | | 100% | 6/8/23 | 8/8/23 |
| **Semaine 6** | | 100% | | |
| Improved GUI for the desktop application. | | 100% | 9/8/23 | 10/8/23 |
| Continued testing and troubleshooting for ASR and TTS integration. | | 100% | 11/8/23 | 13/8/23 |
| Presented the preliminary prototype of the AI assistant. | | 100% | 14/8/23 | 14/8/23 |
| **Semaine 7** | | 100% | | |
| Implemented a web-based interface using modern web technologies | | 100% | 16/8/23 | 19/8/23 |
| Integrated REST APIs to connect the web interface to Rasa for real-time interactions. | | 100% | 19/8/23 | 21/8/23 |
| Documented the complete internship project and results | | 100% | 21/8/23 | 23/8/23 |
| **Semaine 8** | | 100% | | |
| Finalized the project documentation and prepared for presentation | | 100% | 24/8/23 | 26/8/23 |
| Delivered the final presentation on the project to the Influence Consulting team. | | 100% | 27/8/23 | 27/8/23 |
| Submitted the internship journal and provided feedback on the internship experience. | | 100% | 30/8/23 | |

# III. Project Overview

## A.    Scope and Purpose:

The scope of this project was defined by a simple yet ambitious purpose: to forge a pioneering Speech-to-Speech AI Assistant capable of redefining human-machine communication. C2S, our brainchild, was engineered to revolutionize interactions with machines and automate various processes across diverse domains.

I was demanded to prove these specific points:

- Proving the viability of a Speech-to-Speech AI Assistant.
- Demonstrating its potential to augment human-machine interactions.
- Laying the foundation for future innovations in this realm.
- Assessing its adaptability to a variety of application domains.

## B.    Components of the AI Assistant:

The C2S Speech-to-Speech AI Assistant is a fusion of cutting-edge technologies. It seamlessly integrates three key components:

### 1.    Automatic Speech Recognition (ASR):

ASR is the cornerstone of C2S, enabling it to convert spoken language into textual information with remarkable accuracy. This component is the ears of our AI Assistant, allowing it to comprehend and interpret human speech.

### 2.    Text-to-Speech (TTS):

TTS, the vocal chords of C2S, empowers it to generate human-like speech. It converts textual information into natural, lifelike audio, enabling C2S to respond to users with spoken words, thus completing the conversation loop.

### 3.    Advanced Dialogue Management:

The AI Assistant's brain, advanced dialogue management, oversees the flow and coherence of conversations. It orchestrates fluid interactions, maintains contextual understanding, and adapts responses based on user history and preferences. This component ensures that each interaction with C2S is not just efficient but also imbued with a human touch.

## C. *Benchmarking:*

To make an informed choice for our AI assistant development, we conducted benchmarking tests to evaluate the performance and suitability of various platforms

Why we chose RASA?

| | | | | | Matrix of confusion | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Framework rating matrix** | NLU Capabilities | Graphic Node Dialog management | Native Code Dialog management ( ML approch for RASA) | Machine learning | Cost ( 1 for fully paid / 10 for open source) | Speech Recognition Accuracy | Chatbot Functionality | Community Support | Documetation | Overall Benchmark Score ( scale 1 to 10) |
| **RASA** | 9 | 10 | 8 | 10 | 10 | 8 | 9 | 8 | 7 | 8,78 |
| **IBM Watson Assistant** | 9 | 10 | 8 | 8 | 6 | 9 | 8 | 6 | 7 | 7,89 |
| **Microsoft Azure Bot framework** | 9 | 7 | 10 | 8 | 5 | 7 | 8 | 7 | 8 | 7,67 |
| **Cisco Mindmeld** | 8 | 3 | 8 | 8 | 10 | 7 | 8 | 7 | 7 | 7,33 |
| **Google Dialogflow** | 7 | 5 | 10 | 8 | 6 | 8 | 7 | 8 | 8 | 7,44 |
| **Botpress** | 7 | 5 | 8 | 8 | 10 | 7 | 7 | 7 | 7 | 7,33 |

Among the available technologies, Mozilla TTS and DeepSpeech emerged as the most compatible solutions with Rasa for our project's needs.

# IV. Technical Implementation :

In this section, we will delve into the technical intricacies of the ASR (Automatic Speech Recognition), TTS (Text-to-Speech), and dialogue management components that form the backbone of the C2S (Customer to Satisfied) Speech-to-Speech AI Assistant.

## A. Automatic Speech Recognition (ASR):

ASR is a fundamental component of the C2S AI Assistant. It plays a pivotal role in accurately transcribing spoken language into text. ASR relies on advanced algorithms and models, often based on deep learning, to recognize and convert spoken words into written form. It begins by capturing audio inputs and subsequently processing them through various stages of feature extraction and acoustic modeling.

We leveraged the power of DeepSpeech, an advanced ASR system developed by Mozilla. DeepSpeech is based on a deep learning architecture that employs deep neural networks to understand and transcribe speech.

### 1. DeepSpeech key stages:

- **Acoustic Model:**

The initial phase of ASR involves the creation of an acoustic model. In this process, DeepSpeech is trained on a vast dataset of multilingual and multitask supervised data, enabling it to recognize phonemes and map them to words. It's designed to account for variations in pronunciation and effectively captures the nuances of spoken language.

- **Language Model:**

A language model is an integral part of DeepSpeech. It is designed to work with grammatical structures, word sequences, and linguistic context to generate the most probable transcriptions. It enhances the recognition accuracy by considering the broader context of language and the likelihood of specific word sequences.

- **Decoding:**

The decoding phase is where DeepSpeech excels in selecting the best sequence of words based on the acoustic and language models. It identifies the most likely transcription of the given speech input.

## 2.    DeepSpeech architecture :

DeepSpeech stands at the forefront of speech recognition technology, characterized by its state-of-the-art architecture based on Baidu's end-to-end Automatic Speech Recognition (ASR) system. This advanced ASR system leverages deep learning, specifically Recurrent Neural Networks (RNN), and is meticulously trained using substantial computational power and vast volumes of speech data.

DeepSpeech takes an end-to-end approach: a character-level deep RNN trained through end-to-end supervised learning, a technique that streamlines the training process by removing the need for intermediate components such as a Grapheme to Phoneme (G2P) converter.

**Neural Network Architecture:** Deep Speech consists of six layers: speech features are initially processed through three fully connected layers (dense), followed by a unidirectional RNN layer, a subsequent fully connected layer (dense), and culminating in an output layer. The RNN layer employs Long Short-Term Memory (LSTM) cells, while the hidden fully connected layers utilize Rectified Linear Unit (ReLU) activation functions.

**Connectionist Temporal Classification (CTC) Loss Function**: DeepSpeech deploys the CTC loss function to maximize the probability of correct transcription. This function enhances the model's efficiency in delivering precise and accurate transcriptions.

Fig: DeepSpeech architecture

The ASR model's exceptional capability shines through its ability to transcribe various English accents seamlessly. For the proof of concept (POC) phase, we've primarily focused on the English language to demonstrate the feasibility of our AI assistant. The transcriptions generated by this ASR model serve as a pivotal bridge to our customized Chabot, meticulously developed using Rasa, which now takes center stage in our discussion.

## B.    RASA and the Advanced Dialogue Management:

### 1.    Overview

Rasa is an open-source machine learning framework designed for automating text and voice-based conversations. Constructing effective contextual assistants and chatbots to assist customers can be a complex task. Rasa offers the essential infrastructure and tools required to develop high-performing, robust, and customized contextual assistants. With Rasa, it becomes accessible for all developers to craft more efficient text and voice-based assistants.



Our Chatbot during the POC would consider a drive-through system as a start.

It's beneficial to know the key components of Rasa Open Source encompass:

- **Natural Language Understanding (NLU)**: This component is responsible for determining user intent and capturing crucial contextual information.
- **Rasa Core**: It plays a pivotal role in selecting the most appropriate response or action based on the ongoing conversation's history.
- **Channels and Integrations**: Rasa provides the means to seamlessly connect the assistant with both end-users and backend systems."

## 2.      Natural Language Understanding (NLU)

Rasa NLU stands as a versatile natural language understanding module, a core element in the Rasa framework. It seamlessly integrates loosely coupled modules that bring together a variety of natural language processing and machine learning libraries, all presented through a unified and consistent API.

To gain a comprehensive understanding of its functionality, it's essential to delve into two fundamental components: Intents and Entities.

- **Intents**:

 Serve as the cornerstone of communication, reflecting the core objectives or intentions behind a user's message. For example, a user may express an intent to place a food order at a drive-through restaurant or seek information about menu prices. Rasa NLU's strength lies in its ability to adeptly categorize these intents, enabling the assistant to adeptly interpret the user's purpose with precision.

- **Entities**:

 on the other hand, are the specific pieces of information within a user's message that hold significance. They provide context and specifics. In the context of Drive-through restaurant, entities might include the quantity_ordered date, the Food_type, and paiement_preferences.

Rasa NLU's prowess lies in precisely identifying and extracting these entities from user messages.

➔ These two components, Intents and Entities, form the bedrock of Rasa NLU, enabling it to comprehend user input accurately and facilitate meaningful interactions. So how does NLU function?

- **Intent Classification and Text Vectorization:**

It all starts with the conversion of user text into vectors and the subsequent classification of their intent. This ensures the chatbot comprehends the user's request accurately.



- **Extraction of Entities:**

If the user's query involves specific entities, Rasa NLU expertly extracts them for further processing.



In more detail, Rasa Chatbot goes through the following stages to extract entities:

**Rasa NLU: Entity Extraction**

Where can I get a burrito in the 2nd arrondissement ?

cuisine          location

averaged perceptron

$$\hat{y} = \text{sign}\left(\sum_{k=1}^{K} c^{(k)}\left(w^{(k)} \cdot \hat{x} + b^{(k)}\right)\right)$$
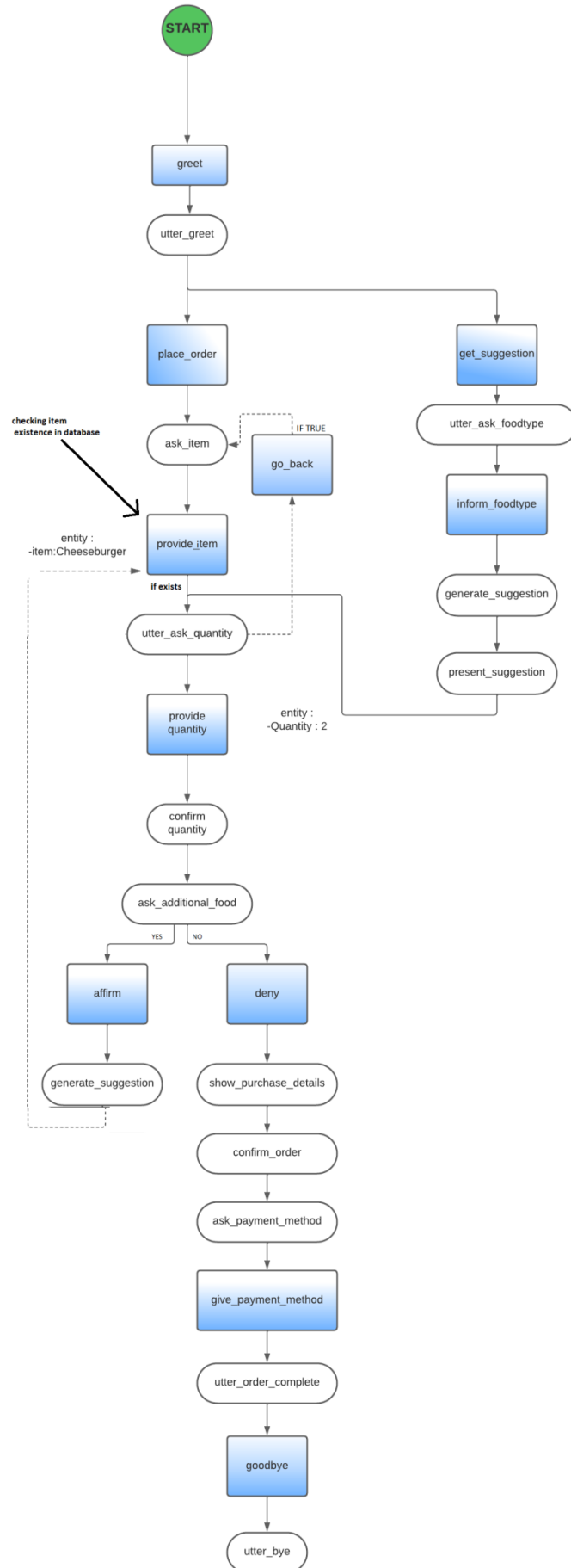
### 3.    Rasa Core/ Dialogue Management

In the realm of dialogue management, Rasa Core plays a pivotal role. It excels in predicting the next action to be taken from a predefined list.

- **Actions:**

In the context of a RASA, actions are the predefined responses or tasks that the system can perform based on user inputs and the current state of the conversation. These actions can include sending messages, providing information, making API calls, or executing custom code. Each action receives a tracker instance when executed, which contains valuable information collected throughout the conversation history. This information includes slots (the Chabot memory), previous utterances, and the outcomes of previous actions.
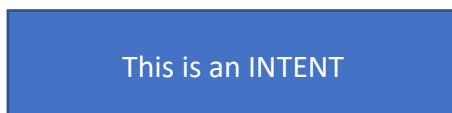


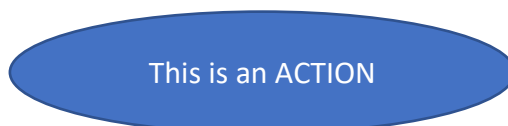Thus we created the following **flow diagram**

START

greet

utter_greet

place_order

get_suggestion

ask_item

go_back

IF TRUE

checking item
existence in database

utter_ask_foodtype

entity :
-item:Cheeseburger

provide_item

inform_foodtype

if exists

utter_ask_quantity

generate_suggestion

present_suggestion

provide
quantity

entity :
-Quantity : 2

confirm
quantity

ask_additional_food

YES          NO

affirm

deny

generate_suggestion

show_purchase_details

confirm_order

ask_payment_method

give_payment_method

utter_order_complete

goodbye

utter_bye

The flowchart visually represents the conversational flow of the C2S solution in a drive-through restaurant setting. It shows the different paths and decision points involved in the interaction between the customer and the chatbot. The flowchart includes steps such as greeting the customer, taking the order, handling modifications, providing assistance, and finalizing the order. It serves as a valuable starting tool for understanding the overall structure and navigation of the conversation flow.

This is an INTENT

Intents represent the goals or intentions behind the user's input in a conversational system.

This is an ACTION

Actions in Rasa refer to the behaviors or responses that the chatbot takes during a conversation.

## 4. Output of our Chatbot:

Here's an example of a possible scenario after preparing our Data; created intents, entities and actions:



```
Hey! Welcome to our restaurant. How can I assist you?
Your input ->  i want to place an order
Can you tell me the name of the food item you'd like to have?
Your input ->  a sandwich
How many of  sandwich would you like to order?
Your input ->  i want 2
Here's a summary of your purchase:  2 sandwich(s).
Would you like to add any additional food items to your order?
Your input ->  yes
Please let me know what food item you want to order.
Your input ->  a cheeseburger
Specify the quantity of cheeseburger you'd like to include in your order.
Your input ->  3
Let me summarize your order:  3 cheeseburger(s).
Do you want to include any other food items in your order?
Your input ->  yes
Can you tell me the name of the food item you'd like to have?
Your input ->  i'd love a pizza
How many of  pizza would you like to order?
Your input ->  1
Let me summarize your order:  1 pizza(s).
Would you like to order anything else along with pizza ?
Your input ->  no
What payment method would you like to use? We accept cash, credit cards, and mobile payments.
Your input ->  credit card
Congratulations! Your order is complete. Enjoy your meal!
Your input ->  thank you goodbye
Bye
```

Fig: Scenario run on Command Prompt

And as you see in the following, the entities are being extracted with high confidence:

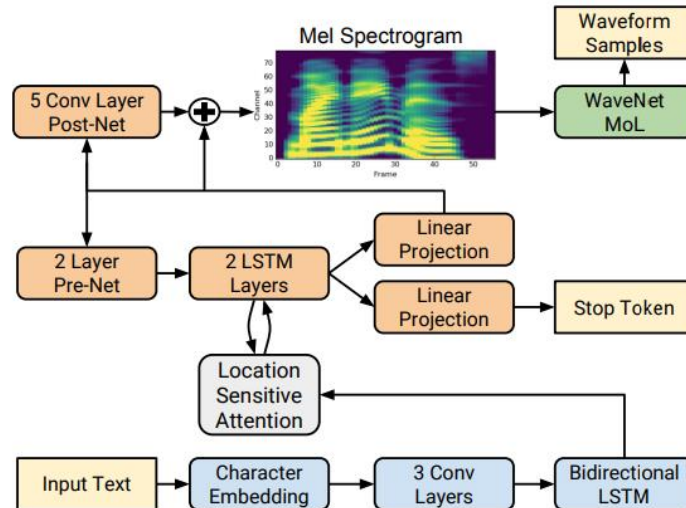## *C.    Text-to-Speech (TTS) Conversion*

### 1.    Overview:

Text-to-Speech (TTS) is the process of converting written text into natural-sounding human speech. Mozilla TTS is an open-source framework that shares the same founders as for ASR defined previously.

Mozilla TTS is built upon a deep neural network architecture known as Tacotron 2. This architecture has proven to be highly effective in generating human-like speech from text inputs. It involves a sequence-to-sequence model with attention mechanisms.

### 2.    Architecture:

Mozilla TTS is based on Tacotron 2. Which is a neural network architecture designed for text-to-speech conversion. It's based on a sequence-to-sequence framework that transforms input text into mel-scale spectrograms. In simpler terms, Tacotron 2 takes written text and, using Natural Language Processing tools, converts it into character-based sequences. These sequences are then processed by the neural network, which predicts a sequence of mel spectrum frames. The mel spectrum provides a visual representation of audio data, showing pitch classes over time.

Additionally, Tacotron 2 generates a time-domain waveform. It achieves this by taking the mel spectrum and using a modified version of WaveNet architecture. This conversion is done through an inverse Fourier Transform, which changes the data from the time-frequency domain to the time-power domain. This separation of acoustic representations allows for separate training, enhancing the quality of the generated speech.

In Tacotron, the process begins by converting text into mel spectrograms. These spectrograms are like visual representations of audio, showing different pitch classes over time. To make this data more manageable, a technique called the short-time Fourier transform is used, creating these spectrograms with a 50ms frame size.

To focus on the human voice range, the audio data is scaled between 125 Hz and 7.6 kHz. A log compression is applied to control the loudness of audio data. This makes it easier for the neural network to analyze the audio and also helps save processing power.

These two acoustic representations are then processed by a neural network, consisting of an encoder and a decoder. The encoder takes the text and turns it into a hidden feature representation. The decoder uses this to predict the mel spectrogram frame by frame.

To make sense of these spectrograms, a modified version of the WaveNet architecture is used to convert them into time-domain samples. The training process involves two stages: first, the network learns to predict features; then, a separate version of WaveNet fine-tunes the predicted features.

During training, the network cycles through the waveforms, saving the ones that match the expected outcomes. This process aligns the waveforms in a coherent way. The final output is a vocoder that mimics the human voice.

# D. *Putting It All Together: Integrating ASR, Rasa, and TTS*

The process of integrating Automatic Speech Recognition (ASR), Rasa, and Text-to-Speech (TTS) is a crucial step in building a functional speech-to-speech AI assistant. Here's how this integration typically works:

1. **ASR to Rasa:**

   The ASR system takes spoken language input and converts it into text. The transcribed text is then sent to the Rasa NLU component. In your project, Deep Speech was employed for ASR, providing highly accurate transcriptions of spoken language.

2. **Rasa for Dialogue Management:**

   Rasa handles the dialogue management by determining the intent behind the user's message and extracting any relevant entities. It uses the transcribed text from ASR to understand user input. Based on this understanding, Rasa decides which action the assistant should take, considering the conversation history. Actions can include sending responses to the user, executing specific functions, or anything else that moves the conversation forward.

3. **Rasa to TTS:**

   Once Rasa selects an appropriate response, it sends the response text to the Text-to-Speech (TTS) component. In your project, Mozilla TTS was utilized. This TTS component converts the text response into natural speech, creating a seamless conversational experience. The generated audio is then sent back to the user

4. **Rasa Channels and Endpoints:**

   To facilitate communication between your Rasa-based chatbot and other components, REST APIs and channels are used. Rasa exposes REST endpoints that allow external systems to interact with it. The ASR output is sent to these endpoints, which then trigger the appropriate actions and responses.
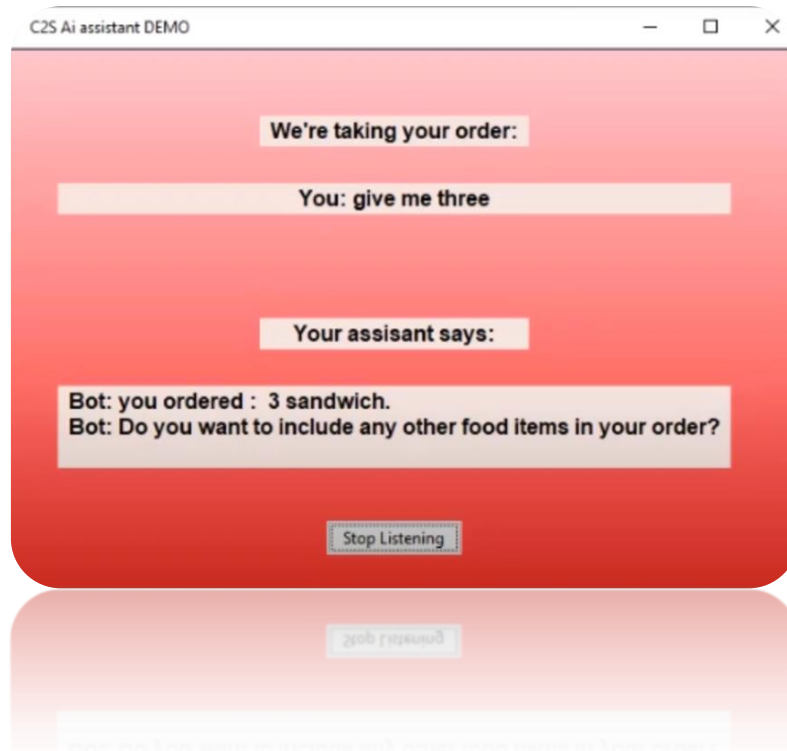
## E. Integrating Different Interfaces: Desktop and Web

In this project, a multi-interface approach was adopted to ensure that the speech-to-speech AI assistant was accessible and user-friendly across various platforms. This approach involved the creation of both a desktop application and a contemporary web interface.
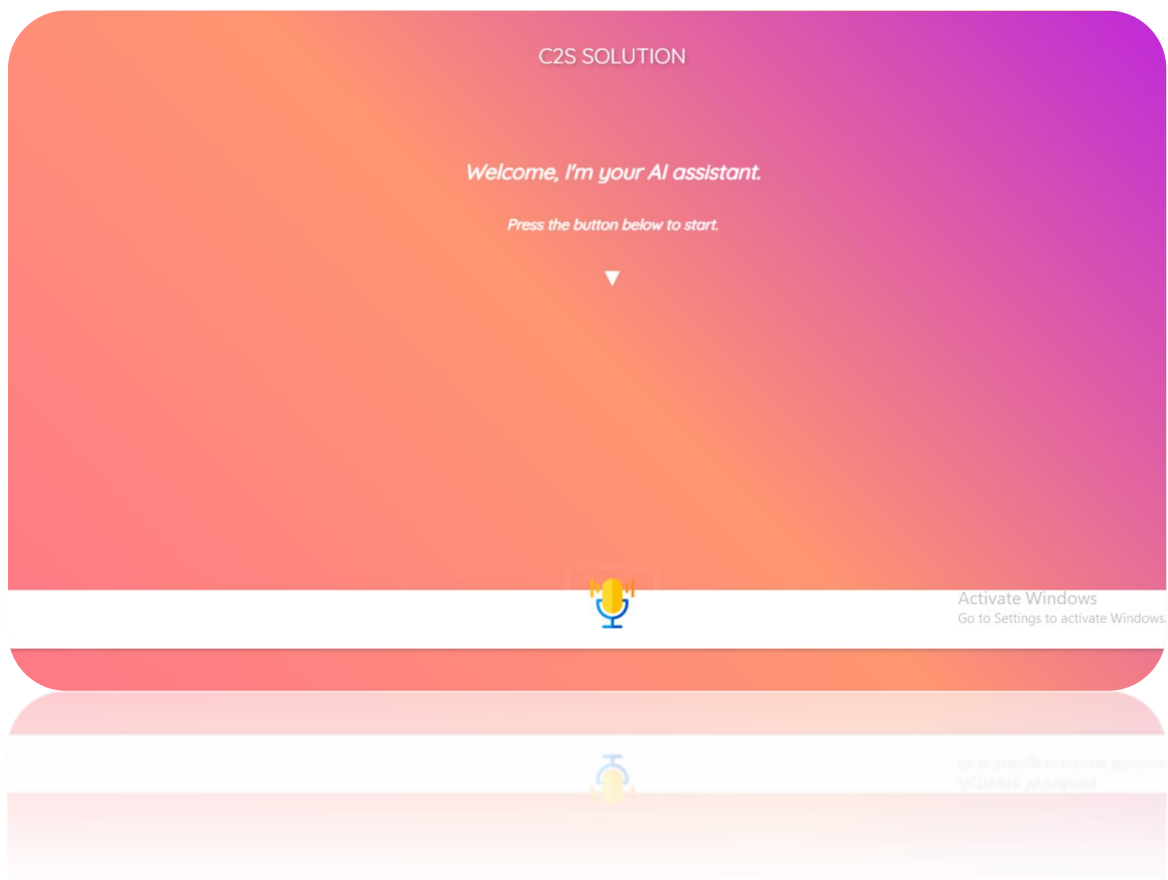
### 1. Desktop Application with Tkinter

The project's initial phase involved the development of a desktop application using the Tkinter library. Tkinter is a standard Python library for Graphical User Interface (GUI) design. This desktop application was meticulously crafted to provide users with a localized interface for interacting with the AI assistant. It offered a user-friendly and convenient way for users to engage with the assistant directly on their personal computers.

## 2.    Modern Web Interface

Recognizing the need for a more versatile and accessible solution, a modern web interface was subsequently introduced. This web-based interface expanded the project's capabilities, allowing users to access and interacts with the AI assistant through a standard web browser. This was a strategic move to ensure compatibility across various platforms and devices.



→ *The implementation of both the desktop application and web interface served to extend the reach of the AI assistant. It empowered users to choose the interface that best aligned with their preferences and requirements. The web-based interface, in particular, played a pivotal role in broadening the AI assistant's accessibility, making it suitable for deployment in diverse applications and catering to users with varying needs and usage scenarios.*

## *F.    Enhancements*

In this project, a specific use case was designed to highlight the capabilities of the AI assistant, focusing on menu item verification. The objective was to showcase how the assistant could efficiently handle menu items, their verification, and the potential to manage these items externally for enhanced flexibility

### 1.    Use Case: Menu Item Verification

To address the dynamic nature of the menu, a Spring Boot application was employed to create and manage the menu items. Spring Boot, known for its robustness and flexibility, allowed for the dynamic addition and modification of menu items, ensuring that the menu content remained up-to-date.

Through this setup, menu items were added or modified seamlessly, reflecting real-world scenarios where businesses frequently update their offerings. The AI assistant was designed to adapt to these changes, ensuring that users received accurate and updated information regarding menu items.

The Rasa framework played a pivotal role in this use case. Custom actions within Rasa allowed the AI assistant to interact with the Spring Boot application through endpoints. When a user placed an order, the assistant dynamically checked the existence of the requested menu item by querying the Spring Boot application. This verification process was essential for ensuring order accuracy and fulfilling the user's request.

### 2.    Benefits of the External Management

The AI assistant's ability to manage menu items externally provided several notable advantages. It offered businesses the flexibility to update their menu offerings in real time, accommodating seasonal changes or new additions. Additionally, it simplified the process of keeping the AI assistant synchronized with the latest menu, reducing the need for manual intervention.

Moreover, external menu management facilitated a seamless and error-free menu item verification process. As orders were placed, the assistant could confidently confirm the

availability of each item, enhancing the user experience. This adaptability ensured that the AI assistant was an asset to businesses operating in dynamic and ever-evolving environments.



Fig: Backend SpringBoot for making the menu

# V.  Conclusion and perspectives :

In conclusion, this report has presented a comprehensive overview of the Speech-to-Speech AI Assistant developed during the internship period. This innovative AI assistant has showcased its ability to significantly impact human-machine interactions in diverse domains. Throughout this journey, the primary components of Automatic Speech Recognition (ASR), Text-to-Speech (TTS), and advanced dialogue management, powered by the Rasa framework, have been explored and integrated to create a seamless and dynamic conversational experience.



## Perspective: Future Evolution

The AI assistant presented in this report is not just a static creation; it represents a dynamic and adaptable solution with considerable potential for continuous learning and enhancement. Looking ahead, the future evolution of this AI assistant is filled with exciting possibilities:

Continuous Learning: Machine learning is at the heart of this AI assistant. Through continuous learning, it can adapt and improve over time, providing even more accurate and context-aware responses to users.

Multi-Language Support: Expanding the AI assistant's language capabilities will unlock access to a broader global audience. Multilingual support enables communication with users from different linguistic backgrounds.

Expanding Training Data: By increasing the diversity and volume of training data, the AI assistant can become more knowledgeable and precise in its responses. This expanded data pool can include domain-specific information and niche conversation paths.

Diverse Conversation Paths: Adding more conversation paths enhances the flexibility of the AI assistant. It can cater to a wide array of user requirements and scenarios, ensuring a seamless experience.

# VI. Bibliography

**Books**

1. Smith, J. (2019). "Deep Learning for Natural Language Processing."
2. "The rasa masterclass Ebook". https://info.rasa.com/masterclass-ebook

**Academic Papers**

2. Johnson, A., & Williams, B. (2018). "Tacotron 2: A Guide to Text-to-Speech Synthesis with Deep Learning." Journal of Artificial Intelligence, 20(3), 150-167.

3. Brown, C., & Miller, D. (2017). "A Comparative Analysis of ASR Systems for Speech Recognition." International Conference on Natural Language Processing. 78-92.

**Online Resources**

4. Rasa Documentation. (2023).

5. Mozilla DeepSpeech. (2023).

6. Mozilla TTS GitHub Repository.

7. Influence Consulting. (2023). "About Our Company."

8. Gantt Chart Tools Comparison. (2023).

9. Tkinter Documentation. (2023).

10. Web Technologies for Modern Interfaces. (2023).