

Plan for Sprint 2 & 3

- Introduction
- Sprint 2 Plan
 - Objective
 - Requirements to Develop
 - Data Collection and Preparation
 - Rationale
 - Infrastructure to Deploy
 - Technology to Use
- Sprint 3 Plan
 - Objective
 - Requirements to Develop
 - Q&A Platform Using GPT-4
 - Furhat Integration
 - Initial Testing and Feedback Gathering
 - Rationale
 - Infrastructure to deploy
 - Technology to Use

Introduction

This document provides a comprehensive and detailed plan for Sprints 2 and 3 of our project to develop a sophisticated Q&A platform. This platform will utilize a Furhat robot, augmented by data derived from specifically chosen websites, leveraging the advanced capabilities of GPT-4 for dynamic and insightful user interactions.

Sprint 2 Plan

Objective

Establish a robust foundational codebase, commence the meticulous process of data scraping and preprocessing, and set up a scalable and secure infrastructure to support development and deployment phases.

Requirements to Develop

Requirement	User Story	Estimation	Priority	Jira Issue
<u>Domain-Specific Language Model function</u> Must provide accurate information of Melbourne Connect to user when asked.	U3.1: As a user interested in services offered at Melbourne Connect, I want the robot to provide detailed information (such as room's information, location and provided services) from websites of Melbourne Connect, so that I can clearly understand the overview of Melbourne Connect without navigating through those websites.	LARGE	MUST HAVE	<input checked="" type="checkbox"/> C2QK-3: Must provide accurate information of Melbourne Connect to user when asked. 待办

<p><u>Furhat Robot interaction</u></p> <p>Must be able to filter and sort information based on user-defined criteria such as location, size, and function of the room.</p>	<p>U2.2: As a user, I want the robot to filter and sort information based on my provided criteria (e.g., location, size, function of the room), so that I can find what I'm looking for more efficiently.</p>	<p>MEDIUM</p>	<p>SHOULD HAVE</p>	<p><input checked="" type="checkbox"/> C2QK-4:</p> <p>Must be able to filter and sort information based on user-defined criteria such as location, size, and function of the room. 待办</p>
<p><u>Furhat Robot interaction</u></p> <p>Must be able to describe the key information based on the website summary to visually impaired user.</p>	<p>U2.3: As a visually impaired user, I want the robot to describe images and visual content from websites, so that I can comprehend visual information that I cannot see.</p>	<p>LARGE</p>	<p>COULD HAVE</p>	<p><input checked="" type="checkbox"/> C2QK-6:</p> <p>Must be able to describe the key information based on the website summary to visually impaired user. 待办</p>
<p><u>UI interactive interface</u></p> <p>Must protect user's private information and provide clear privacy policies.</p>	<p>U1.4: As a user, I want the Q&A platform to protect my personal information and provide clear privacy policies, so that I can ensure the security and confidentiality of my data.</p>	<p>MEDIUM</p>	<p>MUST HAVE</p>	<p><input checked="" type="checkbox"/> C2QK-7:</p> <p>Must protect user's private information and provide clear privacy policies. 待办</p>
<p><u>Domain-Specific Language Model function</u></p> <p>Must be able to summary the information of the website clearly and quickly.</p>	<p>U3.2: As a user, I want the chatbot to automatically extract and summarize the main content of a website I provide, so that I can quickly grasp what the website is about without reading all the content.</p>	<p>LARGE</p>	<p>SHOULD HAVE</p>	<p><input checked="" type="checkbox"/> C2QK-8:</p> <p>Must be able to summarize the information of the website clearly and quickly. 待办</p>
<p><u>Domain-Specific Language Model function</u></p> <p>Must be able to recommend information based on user's interests and website summary.</p>	<p>U3.3: As a user looking for entertainment or news, I want the chatbot to recommend and summarize articles, blog posts, stories, or provide me with the latest headlines based on my interests, so that I can easily find engaging or informative content.</p>	<p>LARGE</p>	<p>MUST HAVE</p>	<p><input checked="" type="checkbox"/> C2QK-9:</p> <p>Must be able to recommend information based on user's interests and website summary. 待办</p>

<u>Domain-Specific Language Model function</u> Must be able to extract information from various careers websites and provide a clear summary job list.	U3.4: As a job seeker, I want the chatbot to compile job listings from various websites, so that I can find opportunities that match my skills and preferences more easily.	LARGE	MUST HAVE	<input checked="" type="checkbox"/> C2QK-1 6: Must be able to extract information from various careers websites and provide a clear summary job list. 待办
<u>Domain-Specific Language Model function</u> Must be able to extract academic information from scientific websites.	U3.5: As a user doing research, I want to ask the chatbot specific questions about content found on academic or scientific websites, so that I can gather information efficiently for my studies or work.	LARGE	MUST HAVE	<input checked="" type="checkbox"/> C2QK-1 1: Must be able to extract academic information from scientific websites. 待办
<u>Furhat Robot interaction</u> Must be able to avoid answering unrelated questions with the website.	U2.5: As a user, I want the robot to avoid answering questions that are unrelated to website navigation content or involve sensitive information, and to inform me when my queries are outside the scope of available information, to maintain professionalism and set clear expectations.	HIGH	MUST HAVE	<input checked="" type="checkbox"/> C2QK-1 3: Must be able to avoid answering unrelated questions with the website. 待办

Data Collection and Preparation [🔗](#)

- Implement advanced web scraping techniques to extract relevant textual and multimedia information from targeted websites.
- Employ sophisticated data cleaning methodologies to ensure data quality and relevance.
- Structure the cleaned data into a format readily usable for training the LLM, ensuring compatibility with GPT-4 input requirements.

Rationale [🔗](#)

We need to train Q/A agent with large amount of multimedia data scraped from chosen websites so that it can successfully provide the domain-specific conversational service to its users. Therefore, the first thing to do is to prepare these data in the format which can be efficiently used by our chosen LLMs.

Infrastructure to Deploy [🔗](#)

- Choose suitable websites for extracting multi-type data including text, images, maps.
- Set up python development, staging, and production environments.

Technology to Use [🔗](#)

- Web Scraping: Python (BeautifulSoup, Scrapy)
- Data Cleaning and Preprocessing: Python (Pandas, NumPy)
- RAG (Retrieval-Augmented Generation): Python

Sprint 3 Plan

Objective

To complete the LLM training, integrate the Q&A agent with the Furhat robot, and design the conversational interface.

Requirements to Develop

Requirement	User Story	Estimation	Priority	Justification	Jira Issue
<u>UI interactive interface</u> Support text interaction in the robot interface.	U1.1: As a user, I want to be able to support text interaction in the robot interface, so I can type and express what I need to query.	SMALL	MUST H...	This is a must have because it is the basic function of user interaction. It may involve only the front end and is expected to be small.	https://comp90082-2024-qa-koala.atlassian.net/browse/C2QK-14
<u>UI interactive interface</u> Generate an livechat box when asking the question.	U1.2: As a user, I want the robot to generate an livechat box when I ask my question, so I can gain the information and answer I need directly.	SMALL	MUST H...	This is a must have because it is the basic function of user interaction. It may involve only the front end and is expected to be small.	https://comp90082-2024-qa-koala.atlassian.net/browse/C2QK-15
<u>Furhat Robot interaction</u> Must summary the information quickly.	U2.1: As a user, I want to interact with the Furhat robot in a conversational manner to obtain information directly, so that I can save time by not having to search and filter information on the web myself.	LARGE	MUST H...	This is a must have because the robot needs to feedback the information needed by the user, and involves the model, UI and robot interaction, which is expected to be a large project.	https://comp90082-2024-qa-koala.atlassian.net/browse/C2QK-2
<u>Domain-Specific Language Model function</u> Must provide accurate information of Melbourne Connect to user when asked.	U3.1: As a user interested in services offered at Melbourne Connect, I want the robot to provide detailed information (such as room's information, location and provided services)	LARGE	MUST H...	This is a must have because the language model needs to analyze user problems and obtain relevant information of Melbourne Connect, the project is	https://comp90082-2024-qa-koala.atlassian.net/browse/C2QK-3

	from websites of Melbourne Connect, so that I can clearly understand the overview of Melbourne Connect without navigating through those websites.			expected to be large due to the model involved.	
<u>Furhat Robot interaction</u> Must be able to filter and sort information based on user-defined criteria such as location, size, and function of the room.	U2.2: As a user, I want the robot to filter and sort information based on my provided criteria (e.g., location, size, function of the room), so that I can find what I'm looking for more efficiently.	MEDIUM	SHOULD...	This is a should have because filters and sorting functions can better present information. It involves robot interaction with UI, and the project is expected to be medium.	https://comp90082-2024-qa-koala.atlassian.net/browse/C2QK-4
<u>UI interactive interface</u> Must be able to guide user on how to using furhat robot.	U1.3: As a new user, I want the robot to offer an introduction on how to use it, so that I can quickly understand and start utilizing the platform services.	SMALL	SHOULD...	This is a should have because the introduction function allows users to better understand the role of QA robots. The front-end is involved, and the engineering quantity is expected to be small.	https://comp90082-2024-qa-koala.atlassian.net/browse/C2QK-5
<u>Furhat Robot interaction</u> Must be able to describe the key information based on the website summary to visually impaired user.	U2.3: As a visually impaired user, I want the robot to describe images and visual content from websites, so that I can comprehend visual information that I cannot see.	LARGE	COULD ...	This is a could have because this feature is an additional optimization feature to visually impaired users and does not affect the core functionality of the project. Due to the development of images, models, etc., a large amount of work is expected.	https://comp90082-2024-qa-koala.atlassian.net/browse/C2QK-6

<u>Domain-Specific Language Model function</u> Must be able to summary the information of the website clearly and quickly.	U3.2: As a user, I want the chatbot to automatically extract and summarize the main content of a website I provide, so that I can quickly grasp what the website is about without reading all the content.	LARGE	SHOULD...	This is a should have because it is based on the language model to complement the functionality. The project is expected to be large due to the model involved.	https://comp90082-2024-qa-koala.atlassian.net/browse/C2QK-8
<u>Domain-Specific Language Model function</u> Must be able to recommend information based on user's interests and website summary.	U3.3: As a user looking for entertainment or news, I want the chatbot to recommend and summarize articles, blog posts, stories, or provide me with the latest headlines based on my interests, so that I can easily find engaging or informative content.	LARGE	MUST H...	This is a must have because it is the key function with language model to provide the correct and necessary answer. The project is expected to be large due to the model involved.	https://comp90082-2024-qa-koala.atlassian.net/browse/C2QK-9
<u>Domain-Specific Language Model function</u> Must be able to extract information from various careers websites and provide a clear summary job list.	U3.4: As a job seeker, I want the chatbot to compile job listings from various websites, so that I can find opportunities that match my skills and preferences more easily.	LARGE	MUST H...	This is a must have because it is the key function with language model to provide the correct and necessary answer. The project is expected to be large due to the model involved.	https://comp90082-2024-qa-koala.atlassian.net/browse/C2QK-16
<u>Domain-Specific Language Model function</u> Must be able to extract academic information from scientific websites.	U3.5: As a user doing research, I want to ask the chatbot specific questions about content found on academic or scientific websites, so that I can gather information efficiently for my studies or work.	LARGE	MUST H...	This is a must have because it is the key function with language model to provide the correct and necessary answer. The project is expected to be large due to the model involved.	https://comp90082-2024-qa-koala.atlassian.net/browse/C2QK-11

<u>Furhat Robot interaction</u> Must be able to translate web content from various of languages.	U2.4: As a user learning a new language, I want the chatbot to translate content from websites in foreign languages(e.g. from Chinese to English), so that I can understand the content without being fluent in the language.	LARGE	COULD ...	This is a could have because it is based on the language model to complement the functionality. The project is expected to be large due to the model involved.	https://comp90082-2024-qa-koala.atlassian.net/browse/C2QK-12
<u>Furhat Robot interaction</u> Must be able to avoid answering unrelated questions with the website.	U2.5: As a user, I want the robot to avoid answering questions that are unrelated to website navigation content or involve sensitive information, and to inform me when my queries are outside the scope of available information, to maintain professionalism and set clear expectations.	HIGH	MUST H...	This is a must have because it is the key function with language model to provide the correct and necessary answer. The project is expected to be large due to the model involved.	https://comp90082-2024-qa-koala.atlassian.net/browse/C2QK-13

Q&A Platform Using GPT-4 [↗](#)

- Configure GPT-4 settings and parameters for optimal learning outcomes using the structured knowledge base.
- Execute RAG (Retrieval-Augmented Generation) sessions, monitoring for accuracy and response relevance.
- Iteratively refine the model based on testing feedback and performance assessments.

Furhat Integration [↗](#)

- Implement the Q&A agent integration on the Furhat robot.

Initial Testing and Feedback Gathering [↗](#)

- Conduct initial testing on functionality.

Rationale [↗](#)

The question processing process is mainly completed by the large language model, therefore we have to train the model and fine-tune it to ensure our core Q/A functions. We then need to integrate it into the Furhat robot and equip it with a suitable interface to help users better use the agent.

Infrastructure to deploy [↗](#)

- GPT-4 development environment
- Furhat robot SDK

Technology to Use [↗](#)

- LLM Integration: GPT-4 APIs and Python
- Furhat robot API