

Leveraging Retrieval-Augmented Generation (RAG) and Large Language Models (LLMs) for Chatbot Development

Introduction

Problem

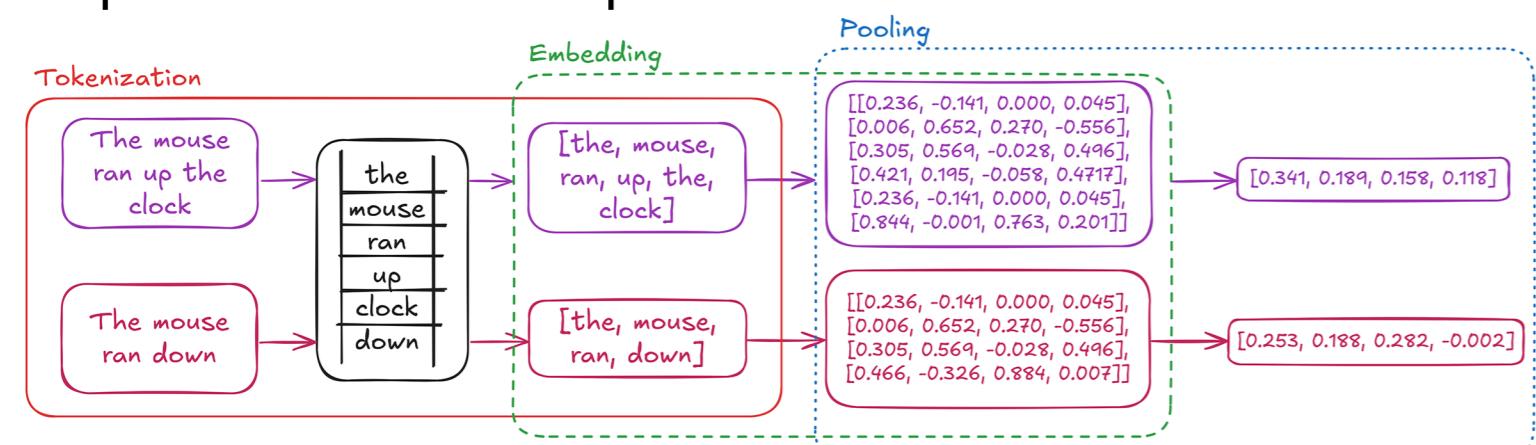
Universities face challenges in providing instant responses to students' administrative and academic queries. Traditional search methods are inefficient, requiring multiple clicks and manual navigation.

Methodology

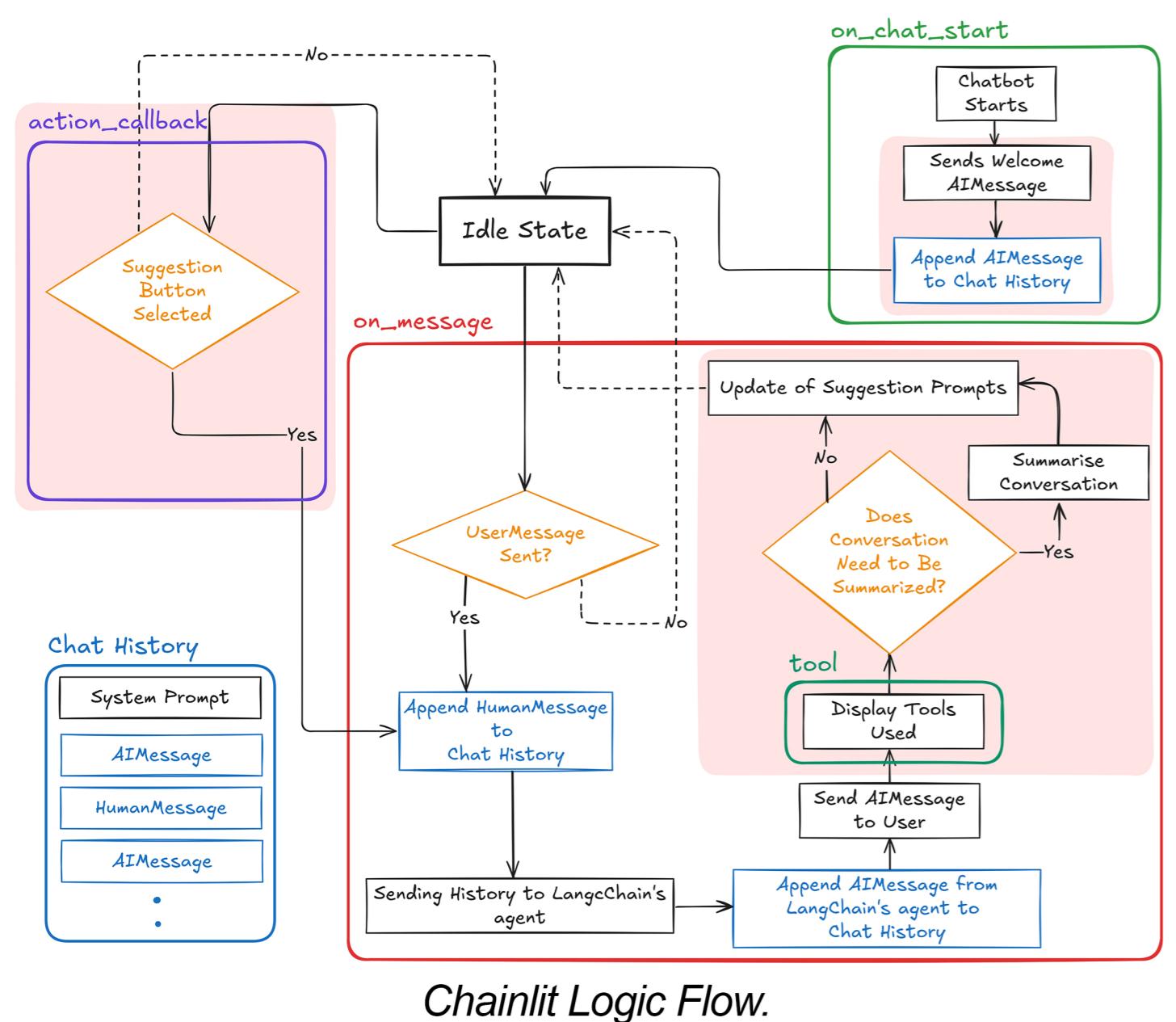
1. Sentence Embeddings

Sentences are converted into numerical representations called embeddings. This process involves three steps:

- Tokenization – Breaking the sentence into smaller units (tokens).
 - Embedding – Converting each token into a numerical vector.
 - Pooling – Combining token embeddings into a single representation for comparison.



Text Processing in RAG: Tokenization, Embedding, and Pooling.



User Study & Results

User Study

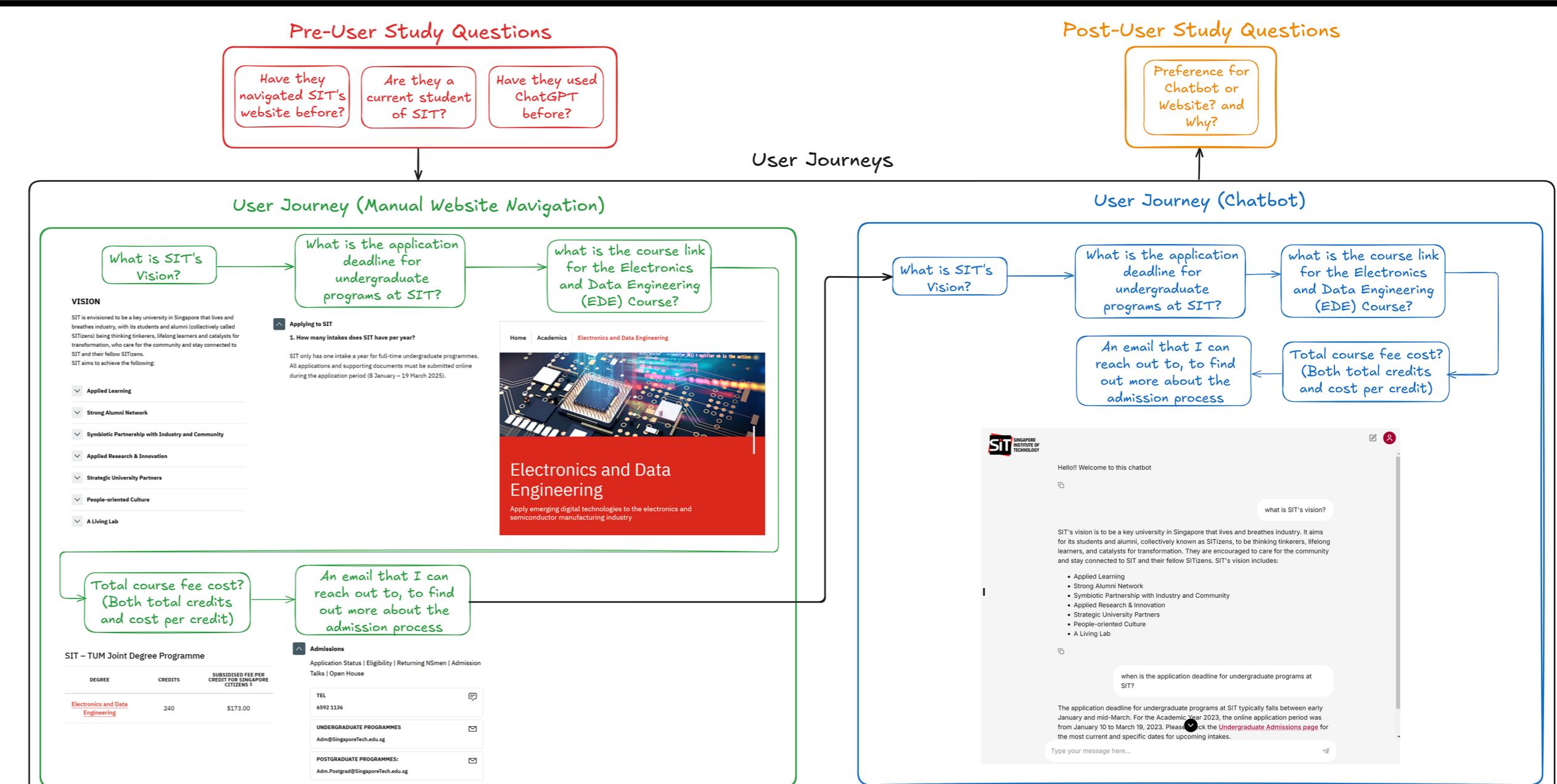
The User Study involved 23 participants comparing chatbot interactions with manual website navigation. Each participant attempted to find answers to five questions, simulating a student's typical user journey. This process was conducted twice – once using the chatbot and once through manual website navigation. Response times, user feedback, number of prompts sent to the chatbot and the number of pages navigated were evaluated to assess efficiency and usability.

Results

Results showed that 82.6% (19 participants) preferred the chatbot, which demonstrated up to a 300% increase in speed compared to manual website navigation. Qualitative analysis was conducted by visualizing user feedback through word clouds and categorizing responses to understand the reasoning behind their preferences, leading to key improvements in the chatbot's functionality.

	Website Avg. Time(s)	Chatbot Avg. Time(s)
Question 1	14.65	19.57
Question 2	112.83	31.70
Question 3	31.22	22.78
Question 3	98.00	25.90
Question 4	93.26	28.82
Overall Average	70.00	25.75

Comparison of Average Response Times



Overall Structure of User Study



Word Cloud of Website Preference



Word Cloud of Chatbot Preference.