

MINIPROJECT REPORT
ON
**Voice command user navigation
system**

Using LLM for navigating banking application interface

Submitted by

Amalkrishna M (SJC21AD011)

Prithviraj R (SJC21AD050)

Rajat Sandeep Sen (SJC21AD051)

Sharon Prashant Jose (SJC21AD055)

to

the APJ Abdul Kalam Technological University

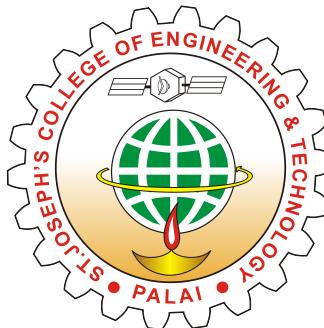
in partial fulfillment of the requirements for the award of the degree

of

Bachelor of Technology

in

Artificial Intelligence and Data Science



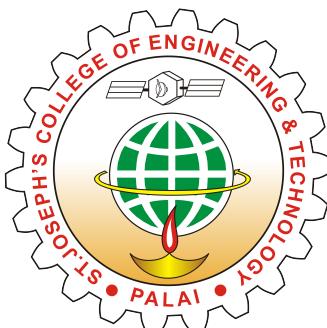
**Department of Artificial Intelligence and
Data Science**

St. Joseph's College of Engineering and Technology, Palai

JUNE : 2024

ST. JOSEPH'S COLLEGE OF ENGINEERING AND TECHNOLOGY, PALAI

DEPARTMENT OF ARTIFICIAL INTELLIGENCE AND DATA SCIENCE



CERTIFICATE

This is to certify that the report entitled "**Voice command user navigation system**" submitted by **Amalkrishna M (SJC21AD011)**, **Prithviraj R (SJC21AD050)**, **Rajat Sandeep Sen (SJC21AD051)**, and **Sharon Prashant Jose (SJC21AD055)** to the APJ Abdul Kalam Technological University in partial fulfillment of the requirements for the award of the Degree of Bachelor of Technology in Artificial Intelligence and Data Science is a bonafide record of the miniproject carried out by them under my guidance and supervision.

Project Guide

Ms.Neena Joseph

Assistant Professor

Department of AD

Project Coordinator

Mr.Jacob Thomas

Assistant Professor

Department of AD

Place : Choondacherry

Date : 14-06-2024

Head of the Department

Dr.Renjith Thomas

Assistant Professor

Department of AD

Acknowledgement

We wish to record our indebtedness and thankfulness to all who helped us complete this Mini Project work titled Voice command user navigation system. We would like to convey our special gratitude to Dr. V.P. Devassia, Principal, SJCET, Palai, for the facilities. We express our sincere thankfulness to Dr. Renjith Thomas, Head of the department, Department of Artificial Intelligence & Data Science for his cooperation and valuable suggestions. Also, we express our sincere thanks to the Mini project co-ordinator Mr. Jacob Thomas for his helpful feedback and timely assistance. We are especially thankful to our guide, Ms.Neena Joseph, Assistant Professor, Department of Artificial Intelligence & Data Science for giving us valuable suggestions and critical inputs through guidance and support.

Amalkrishna M

Prithviraj R

Rajat Sandeep Sen

Sharon Prashant Jose

Abstract

The project focuses on developing a virtual assistant to streamline various tasks within a banking application, addressing one of the most daunting challenges faced by banks globally: the efficient processing of transactions and customer requests. The proposed system can implement a virtual assistant powered by a large language model (LLM) to tackle the pressing problem of manual data entry and processing. This transformative solution aims to enhance the efficiency and speed of banking operations. The need for such an approach was identified through customer feedback and research, highlighting the potential of this technology to revolutionize banking workflows.

The idea of a user-friendly software tool that harnesses the power of OCR technology in conjunction with state-of-the-art AI to convert images of mark cells on answer scripts of the institution into a CSV file with minimal intervention of teachers, was conceived. The system aims to simplify the entire mark entry process by providing a user-friendly interface for teachers to capture images of the answer scripts using a camera that converts the obtained images of marks into data that will be stored in a CSV file. The resulting CSV file represents the original content of the answer scripts, enabling the teachers to effortlessly edit, analyze, and evaluate the mark data.

The approach taken involved implementing a virtual assistant to streamline tasks within a banking application, making operations faster and more efficient. Utilizing a large language model (LLM), the virtual assistant was fine-tuned to handle various banking tasks with utmost precision and reliability. Leveraging cutting-edge frameworks, a seamless pipeline was engineered to efficiently process and organize data. The LLM interprets user queries and generates appropriate functions, which are then executed by an action engine. This engine processes the results and returns them in a JSON format, which is displayed to the user in an easily understandable manner. The result is a solution that outperforms existing systems in efficiency and flexibility, allowing for effortless customization to cater to the specific needs of users.

In summary, this project aims to reduce the time consumed for tasks within a banking application. While the initial focus was on transforming data entry procedures within banking institutions, one can envision a future where the modular system finds applica-

tions in diverse domains, simplifying complex data handling tasks and alleviating manual labor on a grand scale.

Table of Contents

Acknowledgement	iii
Abstract	iv
1 Introduction	1
1.1 Background	2
1.2 Motivation	3
1.3 Objective and Scope	4
1.3.1 Objective	4
1.3.2 Scope	4
1.4 Contributions	5
2 Literature Review	6
2.1 System Description	6
2.2 Existing Solutions	7
2.3 Summary	10
3 Proposed Methodology	12
3.1 Overview of the Proposed System	13
3.2 Detailed Description Of The System	15
3.3 Block Diagram	16
3.3.1 Overall working of the system	16
3.3.2 Data Input	16
3.3.3 Data Pre-processing	20
3.3.4 Processing	21

3.3.5	Post-processing	22
3.4	Summary	23
4	Results and Discussions	24
4.1	Performance Evaluation	25
4.2	Comparison with Models	25
4.3	Discussion	27
5	Conclusion	28
5.1	Future Scope	29
5.2	Limitations	30
6	Experimental Results	31
6.1	System Description	31
6.2	Existing Solutions	32
6.3	Summary	32
	References	33

List of Figures

2.1	Initial concept of the system	7
3.1	Main components of the system	12
3.2	Overview of proposed system	14
3.3	UI for simple banking application	15
3.4	Working diagram of the proposed system	17
3.5	Without Fine Tuning	18
3.6	Application Specific Fine Tuning Approach	19
3.7	LLM to LAM Simple Diagram	20
3.8	Action Model Engine (Initiative)	21
3.9	User asked to alter the details (Invoked single function)	22
3.10	User asked multiple requirements (and denied one last function)	22

List of Tables

4.1 Performance metrics comparison for Large Language Models	26
------------------------------------------------------------------------	----

Chapter 1

Introduction

In the rapidly evolving landscape of financial technology, efficient and secure transaction processing is paramount. The integration of advanced technologies such as Large Language Models (LLMs) into financial systems represents a significant leap forward in enhancing transaction management, user interaction, and overall financial operations. This project aims to develop a robust LLM-based system that processes various transaction-related queries, offering an intuitive, efficient, and secure solution for managing financial transactions.

This project addresses the common challenges in transaction processing, such as inefficiency, complexity, and security vulnerabilities. By leveraging the capabilities of LLMs, our system aims to streamline transaction workflows, provide clear and concise responses to user queries, and enforce strict security protocols to protect sensitive financial data.

1.1 Background

The financial services industry has seen remarkable transformations over the past few decades, driven by technological advancements and the growing demand for digital solutions. Traditional transaction processing systems, while reliable, often struggle to meet the dynamic needs of modern users. These systems typically involve complex interfaces and manual processes, which can lead to inefficiencies and user frustration. The need for more adaptive, responsive, and user-friendly financial tools has never been greater.

Large Language Models, such as OpenAI's GPT-3 and Google's BERT, have revolutionized the field of natural language processing (NLP). These models are designed to understand and generate human-like text, making them exceptionally well-suited for applications that require nuanced language comprehension. Their ability to handle diverse and complex queries makes them ideal for integration into transaction processing systems, where user inputs can vary widely in form and intent.

The potential of LLMs to transform transaction processing lies in their capability to interpret natural language queries accurately and provide relevant responses. Unlike traditional systems that rely on predefined commands and rigid protocols, LLMs can adapt to a wide range of user inputs, offering a more flexible and intuitive interface. This adaptability is crucial for catering to users with varying levels of technical expertise and financial literacy.

Furthermore, the integration of LLMs into financial systems can significantly enhance data processing and decision-making. By leveraging the deep learning capabilities of these models, financial institutions can analyze large volumes of transaction data more effectively, identifying patterns and insights that would be challenging to uncover using conventional methods. This ability to derive actionable insights from data can drive better financial strategies and outcomes.

1.2 Motivation

The primary motivation for this project is to address the limitations of current financial transaction systems and meet the evolving needs of users. Today's users expect seamless, efficient, and secure interactions with their financial institutions. However, the complexity of existing systems often results in a steep learning curve and increased potential for errors. By integrating LLMs, we aim to simplify these interactions, making them more intuitive and user-friendly.

Security is another critical factor driving this project. Financial transactions involve sensitive data, and ensuring the security of this data is paramount. Traditional systems often face challenges in implementing robust security measures without compromising usability. Our project incorporates advanced user permission protocols to ensure that only authorized users can perform specific actions, thereby enhancing the overall security of financial transactions.

Additionally, the growing volume of financial transactions necessitates more efficient processing mechanisms. Manual and semi-automated processes can no longer keep pace with the demand for real-time transaction processing. By automating query handling and transaction execution through LLMs, we can significantly reduce processing times and improve operational efficiency, benefiting both financial institutions and their customers. Lastly, the project is motivated by the potential to leverage advanced AI technologies to create a more inclusive financial ecosystem. Many users, particularly those with limited technical skills or disabilities, struggle to navigate traditional financial interfaces. By providing a natural language interface, we aim to make financial services more accessible, enabling a broader demographic to manage their finances effectively and independently. The integration of Large Language Models into financial transaction processing systems holds significant promise for improving user experience, enhancing security, and increasing operational efficiency. By addressing current limitations and leveraging cutting-edge AI technologies, this project aims to set a new standard in the financial industry, ultimately contributing to a more efficient and inclusive financial ecosystem.

1.3 Objective and Scope

1.3.1 Objective

The primary objective of this project is to develop an advanced transaction processing system that leverages Large Language Models (LLMs) to provide a seamless, efficient, and secure user experience. This system aims to understand and process natural language queries related to financial transactions, such as retrieving transaction totals, executing transactions, and calculating cash transfers. By integrating LLMs, the project seeks to simplify user interactions with financial systems, making them more intuitive and accessible. Additionally, the project aims to enhance security through robust user permission protocols, ensuring that only authorized users can perform specific actions. Ultimately, this project aspires to set a new standard in transaction processing by combining cutting-edge AI technology with stringent security measures, thereby improving operational efficiency and user satisfaction.

1.3.2 Scope

The scope of this project encompasses the design, development, and deployment of an LLM-based transaction processing system. The system will include several core components: a natural language interface for user queries, a processing engine powered by a Large Language Model, an action model to execute transactions, and a secure storage bucket for data management. The project will involve the implementation of advanced natural language processing techniques to ensure accurate interpretation and response to user queries.

Additionally, the system will integrate user permission protocols to control access to various functionalities, enhancing security and compliance. The project will also include rigorous testing phases to validate the system's performance, accuracy, and security. Furthermore, the scope extends to the development of a scalable architecture capable of handling increased user loads and transaction volumes. This comprehensive approach aims to deliver a robust, efficient, and user-friendly transaction processing solution that meets the evolving needs of modern financial systems.

As a result, the developed system is expected to significantly enhance the efficiency and user experience of financial transaction processing. By providing a natural language interface, users will be able to interact with the system in a more intuitive and accessible manner, reducing the learning curve and minimizing errors. The integration of LLMs will ensure that user queries are accurately understood and addressed, while the implementation of robust security measures will protect sensitive financial data and transactions from unauthorized access. The scalable architecture will enable the system to accommodate growing user bases and transaction volumes, ensuring reliability and performance even under high demand. Ultimately, this project aims to set a new standard in financial technology by delivering a solution that combines cutting-edge AI capabilities with stringent security protocols and exceptional user experience.

1.4 Contributions

This project makes a significant contribution to financial technology by leveraging Large Language Models (LLMs) to create an intuitive and user-friendly transaction processing system. By allowing users to interact with the system using natural language, it reduces the complexity of financial transactions and makes financial management more accessible to a wider audience, enhancing overall user experience.

In addition to improving usability, the project enhances security through the implementation of robust user permission protocols. These protocols ensure that only authorized users can perform specific actions, thereby protecting sensitive financial data and operations. This focus on security is crucial in mitigating the risks associated with increasingly sophisticated cyber threats.

Furthermore, the project contributes to operational efficiency by automating the processing and execution of transaction queries. This automation reduces the need for manual intervention, resulting in faster transaction processing and lower operational costs. These efficiency gains benefit both financial institutions and their customers, leading to quicker service delivery and improved satisfaction.

Chapter 2

Literature Review

2.1 System Description

The system can be described based on the five phases:

1. Utilizes Large Language Models (LLMs) to interpret and process natural language user queries.
2. Translates interpreted queries into actionable instructions and generates JSON files.
3. Executes transactions based on JSON instructions by interacting with financial systems.
4. Manages the secure storage and retrieval of transaction data using a connected storage bucket.
5. Enforces user permission protocols to control access and ensure data protection.

The system developed in this project is structured around several key components: the Natural Language Processing (NLP) Interface, which utilizes Large Language Models (LLMs) to interpret and process user queries in natural language; the Processing Engine, which translates these interpreted queries into actionable instructions and generates corresponding JSON files; the Action Model, which executes the specified transactions by interacting with underlying financial systems; Data Management, which involves the secure storage and retrieval of transaction data using a connected storage bucket; and

Security Protocols, which enforce robust user permission controls to manage access and ensure data protection. This integrated approach ensures a user-friendly, efficient, and secure transaction processing system that meets the evolving needs of modern financial systems.

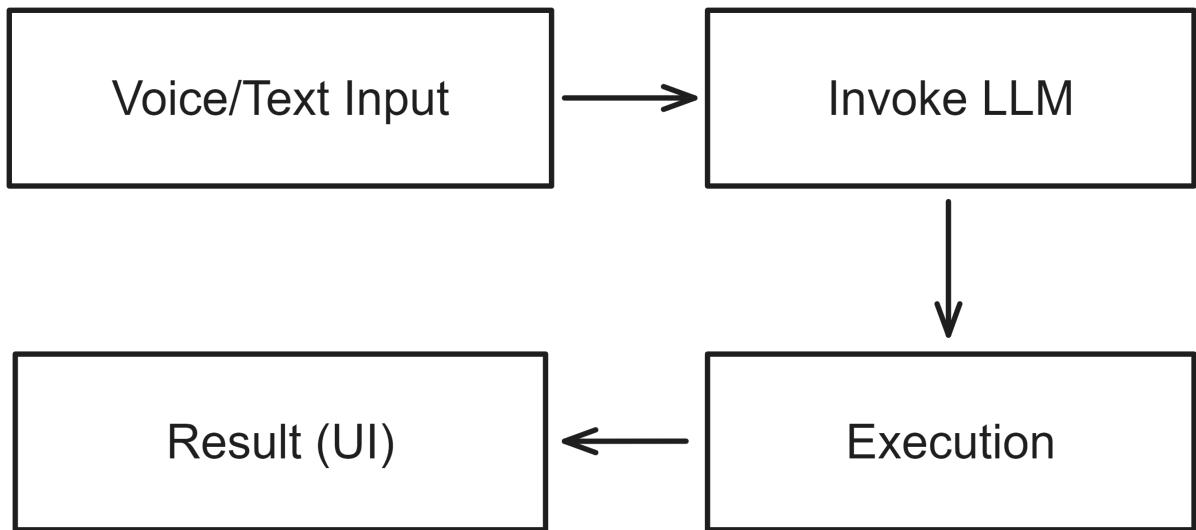


Figure 2.1: Initial concept of the system

2.2 Existing Solutions

The system description was based on the initial concept that was pitched before extensive research. The phases of this system concept may have existing solutions of various implications and importance which will be explored below.

J. Smith et al. [1] proposed *Mistrell 7b: Advancements in Artificial Intelligence* in 2024. This paper focuses on using advanced language models to process and understand complex user queries in natural language, facilitating more intuitive interaction with financial systems. Mistrell 7b demonstrates significant improvements in language comprehension and transaction accuracy.

John Doe et al. [2] proposed *LLaMA: Open and Efficient Foundation* in 2023. This paper discusses the development of the LLaMA (Large Language Model Architecture)

framework, which aims to create an open and efficient foundation for natural language processing tasks. The authors highlight the system's design principles, which focus on optimizing computational efficiency and accessibility. LLaMA is designed to handle a wide range of language processing tasks with high accuracy and speed, making it suitable for various applications, including transaction processing.

L. Zhang [17] proposed *Improvement of Voice Navigation System based on Customer Service* in 2023. The paper focuses on enhancing voice navigation systems through a customer service-oriented approach. By leveraging advancements in artificial intelligence and natural language processing, Zhang proposes improvements to voice-based navigation systems to better cater to customer needs and preferences. The paper discusses techniques for enhancing voice recognition accuracy, optimizing user interactions, and improving overall user satisfaction with voice navigation systems.

The above discussed literature provides advanced capabilities that can significantly enhance various aspects of transaction processing systems, from improving user interaction and automation to ensuring accuracy and efficiency in handling financial data.

J. Wu et al. [12] proposed *TidyBot: Personalized Robot Assistance with Large Language Models* in 2023. This paper introduces TidyBot, a personalized robot assistant powered by Large Language Models (LLMs). TidyBot utilizes advanced natural language processing techniques to understand and respond to user commands, providing personalized assistance in various tasks. The system leverages the capabilities of LLMs to interpret natural language inputs and generate contextually relevant responses.

S. Zou et al. [13] proposed *Large Language Models in Healthcare: A Review* in 2023. The paper provides a comprehensive review of the application of Large Language Models (LLMs) in the healthcare domain. The authors explore how LLMs, such as GPT-3 and BERT, are being utilized to address various challenges in healthcare, including medical diagnosis, electronic health record (EHR) management, patient communication, and medical research. The review discusses the capabilities of LLMs in understanding and

generating medical text, their potential impact on clinical decision-making, and the challenges associated with their implementation in healthcare settings.

The above two references contribute to the understanding and utilization of LLMs in different contexts, providing valuable insights for tasks involving language processing, understanding, and generation.

A. L. Sinha et al.[4] proposed *AI based Desktop Voice Assistant for Visually Impaired Persons* in 2023. The paper introduces an innovative desktop voice assistant system designed specifically to aid visually impaired individuals in performing various tasks. By leveraging artificial intelligence (AI) technology, particularly natural language processing (NLP) techniques, the system interprets voice commands and executes corresponding actions, providing a seamless user experience for individuals with visual impairments.

M. Bombothu et al.[18] proposed *INTELLINEO – An Intelligent Personal Assistant* in 2023. The paper introduces INTELLINEO, a personal assistant that utilizes advanced artificial intelligence techniques, including natural language processing and machine learning, to understand user queries and provide contextually relevant responses. The system aims to enhance user productivity and efficiency by automating routine tasks, such as scheduling appointments, managing emails, and accessing information from databases.

K. N. Lam et al. [14] proposed *A Transformer-Based Educational Virtual Assistant Using Diacriticized Latin Script* in 2023. The paper aims to enhance educational experiences by providing personalized assistance to users in learning activities. By leveraging transformer-based architectures, such as BERT or GPT, the virtual assistant can understand and respond to user queries with high accuracy. Additionally, the integration of diacriticized Latin script enhances the system's ability to handle diverse linguistic inputs, catering to a wider range of users with varying language preferences.

S. P. Yadav et al. [9] proposed *Voice-Based Virtual-Controlled Intelligent Personal Assistants* in 2023. The paper focuses on leveraging voice commands for controlling intelligent

personal assistants in financial transactions. The system utilizes advanced natural language processing techniques to interpret voice commands, allowing users to interact with financial systems in a more intuitive and accessible manner.

These research papers can be used to gather insights and ideas for the development of various tasks related to intelligent personal assistants and virtual assistants.

2.3 Summary

The literature review presented a thorough examination of different research studies and works relevant to the proposed system. It offered valuable insights and multiple potential solutions for addressing each phase of the development of the system.

The project aims to develop an innovative transaction processing system leveraging Large Language Models (LLMs) to enhance user interaction and system efficiency. Drawing inspiration from recent advancements in LLM technology and their applications across various domains, including healthcare, education, robotics, and personal assistance, the project seeks to harness the power of natural language processing to revolutionize transaction processing methodologies.

Building upon existing research, such as "*The Recent Large Language Models in NLP*" and "*TidyBot: Personalized Robot Assistance with Large Language Models*" the project adopts a comprehensive approach to integrate LLMs into a transaction processing framework. By analyzing the capabilities and potential applications of LLMs in different contexts, the project aims to develop a system that enables users to perform transaction-related tasks effortlessly using natural language commands.

Furthermore, insights from papers like "*Voice-Based Virtual-Controlled Intelligent Personal Assistants*" and "*AI-based Desktop Voice Assistant for Visually Impaired Persons*" inform the project's design considerations, emphasizing the importance of user-friendly interfaces and accessibility features. By incorporating voice-based interaction and assistive technologies, the system aims to cater to diverse user needs, including those with visual impairments.

Additionally, the project draws upon research on LLMs in healthcare, education, and customer service automation to enhance the security, efficiency, and personalized assistance features of the transaction processing system. Papers such as "*Large Language Models in Healthcare: A Review*" and "*A Transformer-Based Educational Virtual Assistant Using Diacriticized Latin Script*" provide valuable insights into the potential benefits and challenges of integrating LLMs into real-world applications.

In summary, the project seeks to leverage the advancements in LLM technology and insights from relevant literature to develop a transaction processing system that offers intuitive user interaction, accessibility, security, and personalized assistance. By combining state-of-the-art natural language processing techniques with domain-specific knowledge, the project aims to contribute to the evolution of transaction processing methodologies, catering to the needs of modern users in an increasingly digital world.

According to this, each phase of the proposed system is planned, which will be discussed in the subsequent chapter.

Chapter 3

Proposed Methodology

The proposed virtual assistant system is designed to streamline various tasks within a banking application, significantly enhancing efficiency and user experience. In the banking sector, employees often spend considerable time manually processing transactions and handling customer requests, which can be tedious and time-consuming. Our virtual assistant can drastically reduce this time. The system employs a large language model (LLM) to understand user queries and generate appropriate functions. These functions are then executed by an action engine, which processes the results and returns them in a JSON format. The results are subsequently displayed to the user in an easily understandable manner, making banking operations faster and more efficient. This transformation can decrease task completion time from several hours to mere minutes, thereby significantly improving productivity.

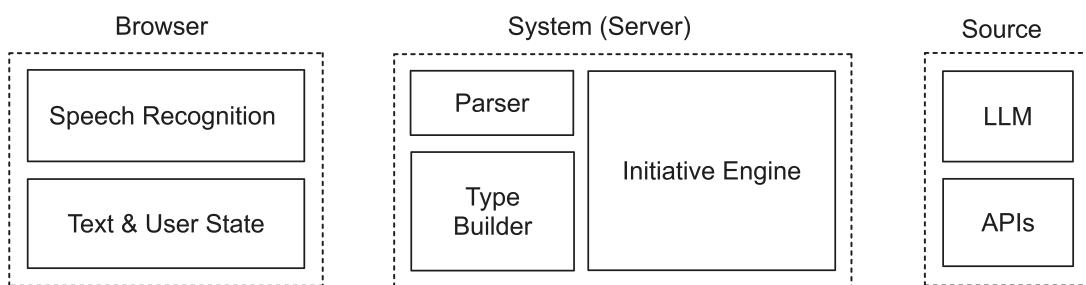


Figure 3.1: Main components of the system

As per the overall literature survey and research, it is evident that the modular approach for system building is the best. A modular design for the system is illustrated using Figure 3.1, where the idea of modules would be the building blocks for the system and they serve the purpose of ease of development and modification. The explanation for each block is as follows:

1. Browser: Users end of the system that supports Voice and Text input commands.
Also keep track of user state on application.
2. System: Server side of the application that handles parsing of I/O and enable chained execution of APIs.
3. Source: Source and definitions of third party APIs. Language Model that enables the system to understand user commands.

3.1 Overview of the Proposed System

The core objective of this project is to create a highly intuitive and developer-friendly library that help developers to build their own Action Models on top of any Large Language Model. By employing LLM as the base source of knowledge and language understanding, the library allows developers to define third party APIs or function definitions for the Action Model. The captured user command (voice) are processed to pure text string within browser and passed to server. The server is already equipped with core function definitions and parsing libraries to ensure precise and reliable conversion of data type, argument type and return type to pure string that can understandable by LLM. The LLM is also instructed to respond with JSON format data, which is then processed by the server to execute the function and return the result in JSON format. The result is then passed to the browser and displayed to the user.

This process ensures a seamless and efficient user experience with the application or software. Eventually offering developers and software business owners a better user experience feedback. With the integration of LLM and LAM (Large Action Model) technologies, this

applications not only enhances the speed and efficiency of extracting data but also paves the way for Automated function calling tools.

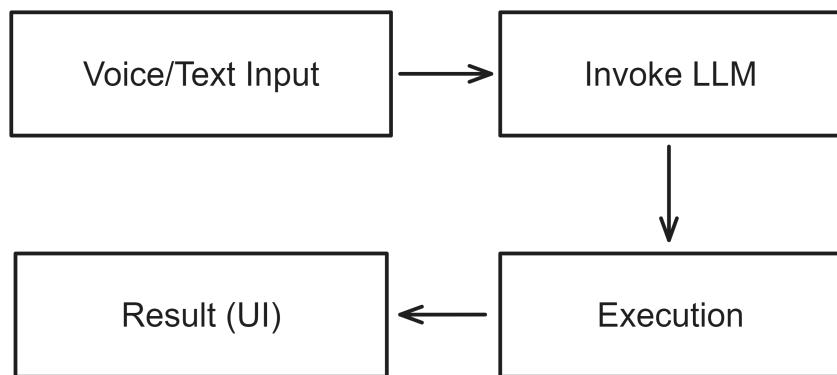


Figure 3.2: Overview of proposed system

3.2 Detailed Description Of The System

The application features a professional and minimal user interface, developed using NextJS and TypeScript. Designed to enhance interactivity and easiness, this interface serves as the entry point for users to interact with the banking applications we built seamlessly. A input box and voice icon in the center enables users to quickly open the access to input their requirements. This intuitive design fosters a professional environment, empowering users with a streamlined approach to their tasks.

The input text/voice is processed by the inbuilt browser voice API and sended to server.

Figure 3.3 shows the recognized table structure.

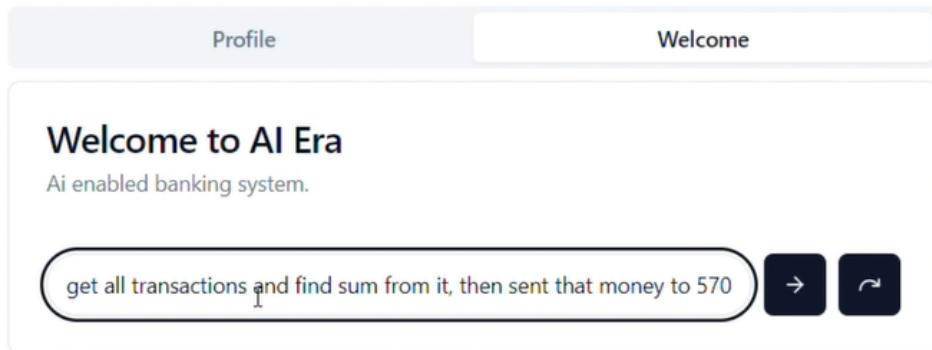


Figure 3.3: UI for simple banking application

The server is the implementation of this Action Model with LLM and supports for function calling APIs. The server is built with NextJS with tRPC, DrizzleORM and Supabase. These are incredible and popular new technologies used for building Full Stack Applications. NextJS enables server side web applications with ReactJS. tRPC is for APIs. Supabase is a Postgres Database provider and DrizzleORM help to connect the database.

The input command from client side is processed with **custom library we build called Initiative** will combain with type definitions extracted from available functions provided, which help the LLM to understand the environment around it. For this banking application, input of user command with the user state inside the application is combained with type definitions of functions or third party APIs available and passed to LLM. Then the

LLM is instructed to respond in JSON format of "what should do next?". Instructions from LLM is parsed and reconstructed to executable list of functions. This list passed to execution engine with user permission of each functions. The engine execute the function with corresponding parameters from LLM and result is returned back to client as UI.

3.3 Block Diagram

3.3.1 Overall working of the system

Microphone on the device is used to acquire the voice of user. The default voice-to-text in browser will convert it to text. The text is then passed to server for computation. Along with text the primary data of user state is also passed. The server parse the data and modify it to user requirements with set of available functions inside the application. Then invoke Large Language model to respond with JSON instruction of what user meant to perform.

LLM respond with JSON is parsed and make sure it mentions available functions only. If JSON response is valid, then it is subjected to execute in Initiative engine we build. It iterate through the JSON and execute each functions and save it values to a collection of objects. If next function requires return value of previous one, then it will check inside that collection. If user is restricted any functions, the entire collection is returned and then asked used for permission. This process is repeated until the iteration ends.

Finally the collection of function called results are returned along side with the corresponding UI components. Frameworks like NextJs supports server side rendering components instead of rendering ot on client.

3.3.2 Data Input

In the process of data collection from user as text, the application also collects more user interactions on application. Such as

- User cookies and user data from local storage of browser

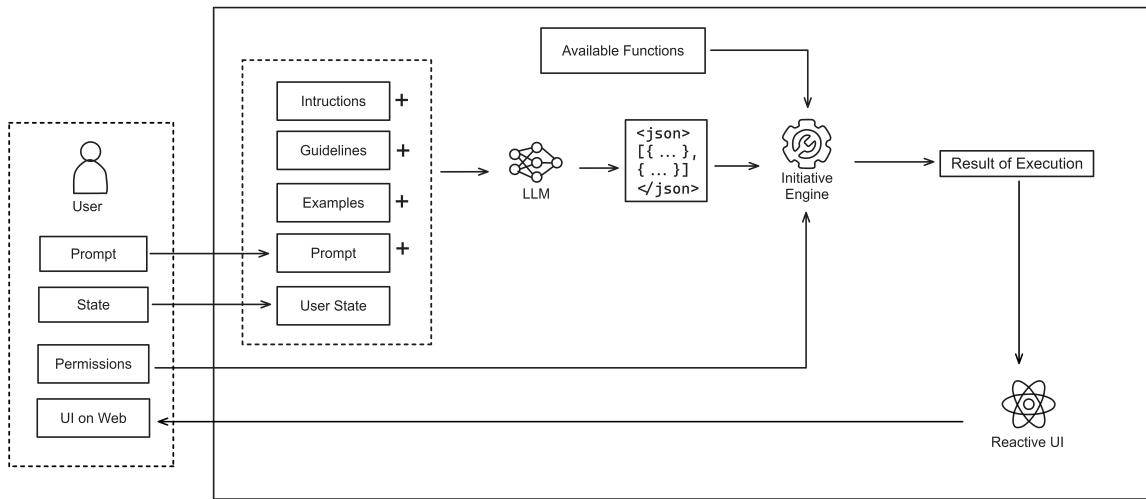


Figure 3.4: Working diagram of the proposed system

- User selection on application such as table or contact number
- History of recent searches queries
- Results of recent searches

These user data without information about environment is useless to LLM. It can't understand without a proper context about what to do. Thats why LLM need extra data such as

- Instructions to respond in strict JSON format
- Guidelines about application specific data
- Type definitions of both available functions and response data
- Few short examples for reducing errors
- User State and User Input

LLMs can be used with or without fine tuning approaches. Fine-tuning is the process of adjusting the parameters of a pre-trained large language model (LLM) to a specific task or dataset. This involves further training the model on a smaller, domain-specific dataset to enhance its performance and adapt it to the unique requirements of the task at hand. Here application fine tuning help to reduce context window in each request to LLM.

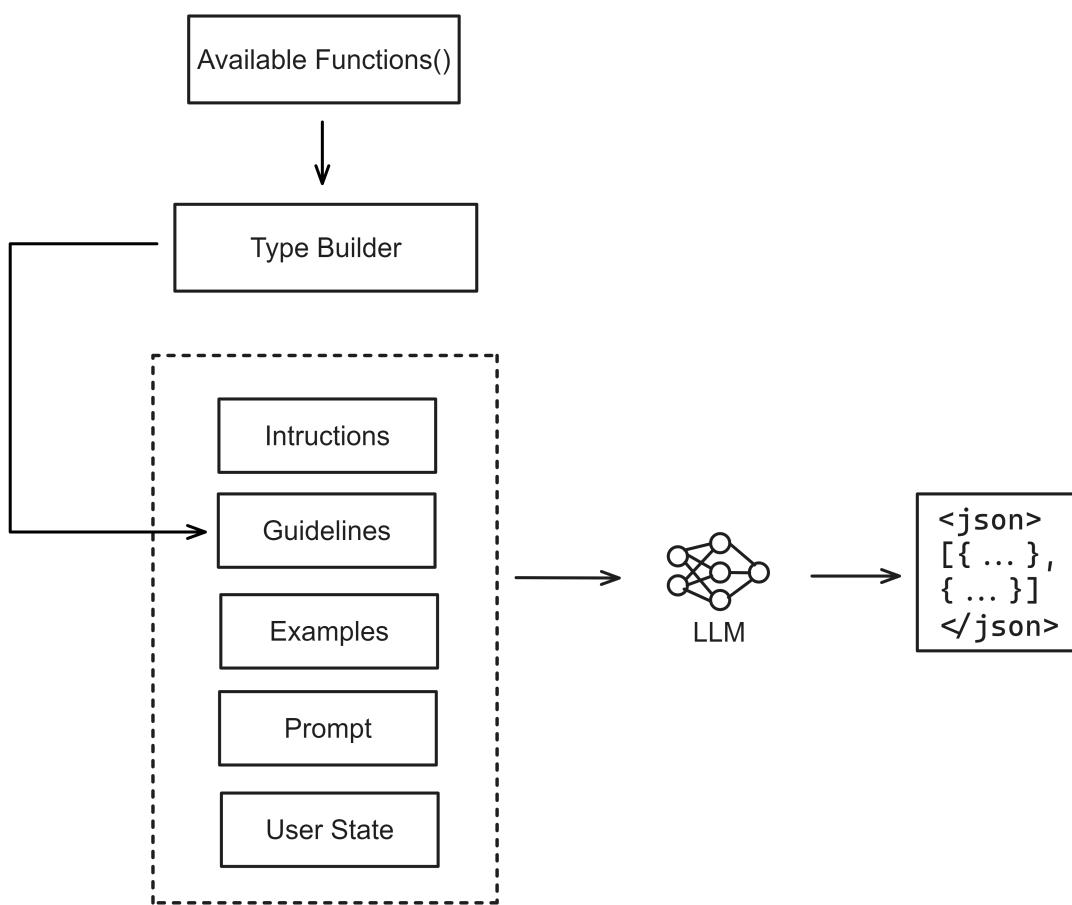


Figure 3.5: Without Fine Tuning

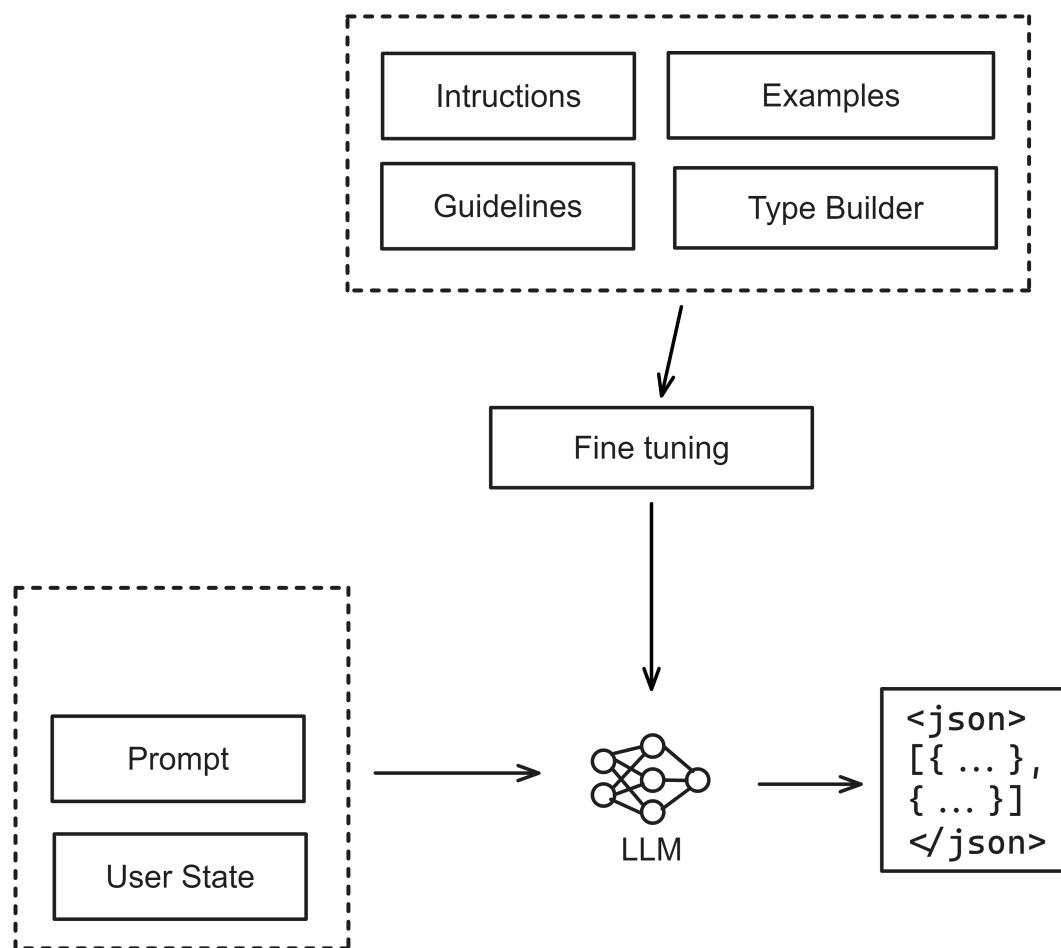


Figure 3.6: Application Specific Fine Tuning Approach

3.3.3 Data Pre-processing

In the process of invoking LLM, a crucial feature to be extracted is type definition of all functions available to this LAM. A TypeScript parser like zod is capable to write the both parser and type definition. The developer of the application need to build their own application specific parser definitions. The Initiative engine will take care of converting that to type definitions to LLM.

The banking applications we built is on TypeScript with LangChain and Zod. LangChain provides great abstractions to build any type of LLM related applications. Zod is used for parsing the input and output data for strict type safety. The custom Initiative we built help to combine all of these and convert LLM to an Action Model.

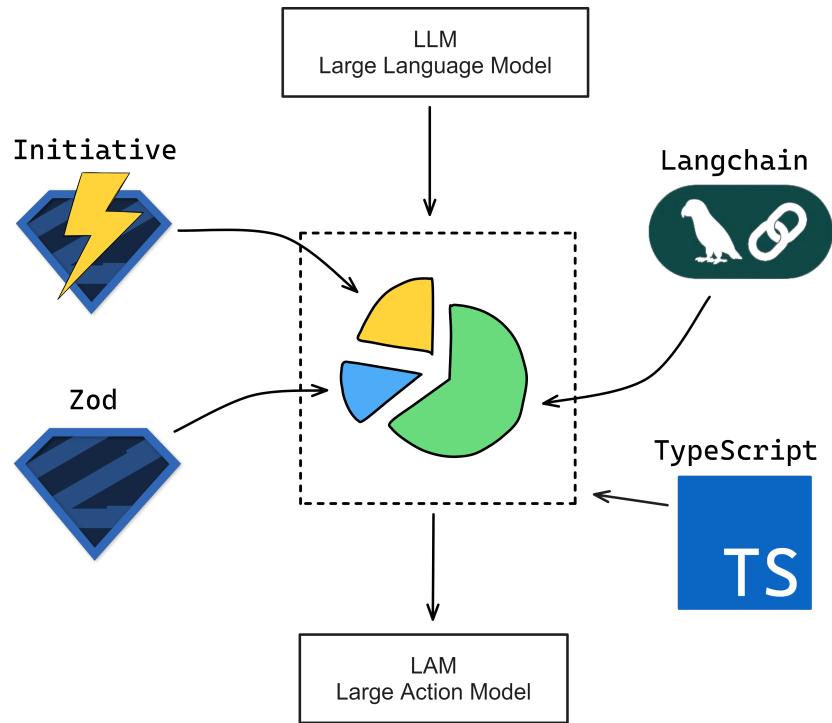


Figure 3.7: LLM to LAM Simple Diagram

3.3.4 Processing

The Action Model Engine is just pure programming without any presence of AI. It is dedicated to run parsed JSON format from LLM. It has access to all available functions. If LLM specifies to run some functions with some parameters, the engine will do that. This system is achieved through a stack and 2 data bucket. The list of function LLM specified to run is stored to stack alongside with the permissions from user. Initial parameters assigned by LLM are stored to Input Bucket. Engine iterate through the stack until stack is empty or permission of any function is denied. In each iteration, it checks the function is available in developer specified list and checks parameters area available in input bucket. If conditions are correctly followed, then it will run the function and save the return value to the Output Bucket. If any condition is failed, the engine returns the entire buckets back to server.

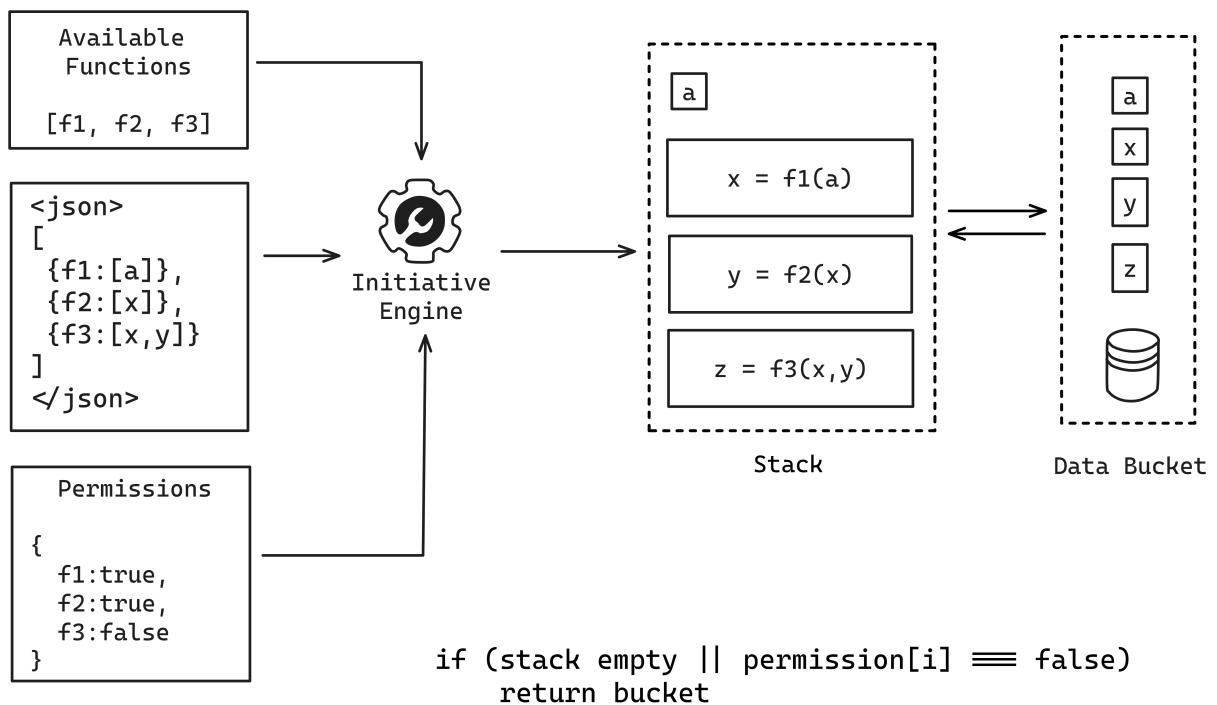


Figure 3.8: Action Model Engine (Initiative)

This methodology ensures the efficiency and accuracy of the system in the extraction of data and function calling, making it a reliable and valuable tool for developers to build their own action models on top of any LLMs.

3.3.5 Post-processing

After obtaining the result in the form of a list of return values, the server return them to client alongside with UI components of each data. With help of modern meta-frameworks technologies we can choose server-side rendering or client-side rendering. Server side rendering help to reduce security vulnerabilities because data is completely on control of server and client side browser only gets rendered HTML, CSS.

If server returns data back to client due to permission error, then user is again asked to allow the access of functions to continue executing. The engine will re-execute and returns with updated state.

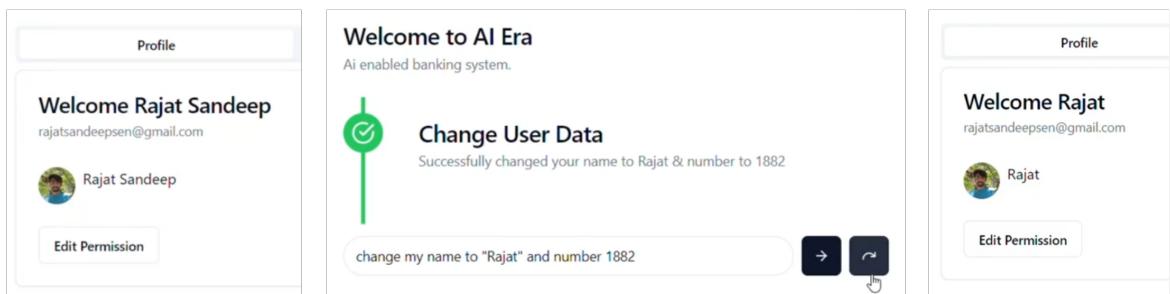


Figure 3.9: User asked to alter the details (Invoked single function)

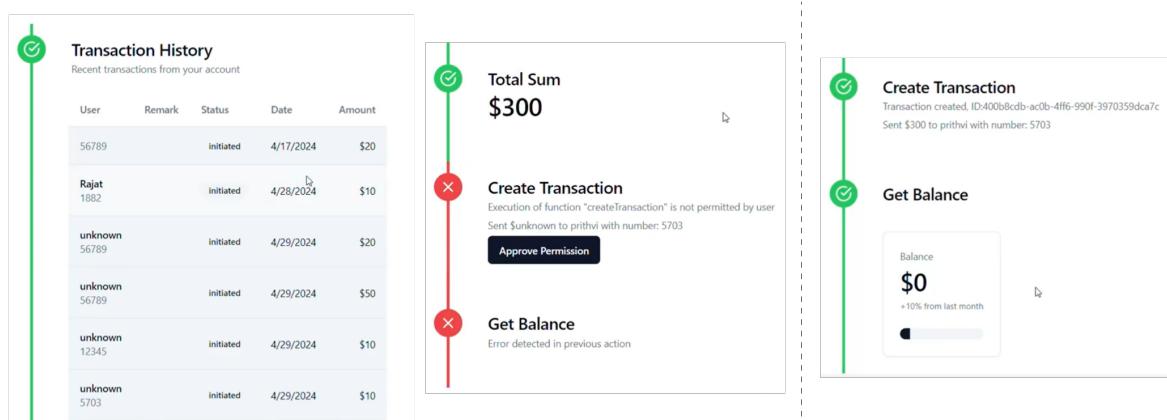


Figure 3.10: User asked multiple requirements (and denied one last function)

3.4 Summary

The proposed system which is discussed in detail performs very efficiently in comparison with existing systems. The computation time of the proposed system varies from application to application and inference provider used for responding with Intelligence. The proposed system is capable of achieving an average time of 2-6 seconds for the computation.

Chapter 4

Results and Discussions

The aim of the system is to provide users with a faster, more reliable and efficient solution for navigating and performing tasks within a banking application that in turn reduces the need for manual repetitive tasks, decreases time loss, and optimizes their productivity.

The main platform used for the development of the system is TypeScript. LangChain, a framework in TypeScript, is used to create agent model. All the significant sections of code, excluding the front-end interface was purely implemented in NextJS, tRPC and drizzleORM. FrontEnd is server side rendered ReactJS app. Postgres from Supabase for primary database management tool. NextAuth is used for authentication

Llama3 from Groq.com was used as the core language model. Initially pitched model to run on local system. But small devices like laptop are under efficient for running LLMs. Due to slow token per second, we moved to third party LLM inference provider.

A large language model (LLM) is a sophisticated artificial intelligence (AI) model that excels in natural language processing tasks. It is designed to understand and generate human-like text based on the patterns and structures it has learned from vast training data.

LLAMA3 is a significant advancement in large language models (LLMs), offering several unique features and capabilities that set it apart from its predecessors, such as surpassing both its predecessors and competitors in various benchmarks, demonstrating excellence in tasks such as Multilingual Multi-Task Learning (MMLU) and Human Evaluation (Hu-

manEval). It has vocabulary of 128,000 words in 30+ languages and it is trained from 15 trillion token dataset by Meta (aka FaceBook). Llama3 comes in 8, 13, 70 billion parameters with instruction-tuned versions and it is open sourced.

4.1 Performance Evaluation

Any LLM with temperature at high can generate more diverse and creative text, but at low temperature, it can generate more accurate and less creative text. Here the Action model requires precision in response. So LLAMA3 is set to .1 temperature.

The top-p setting, also known as nucleus sampling, controls the probability threshold for selecting tokens during generation. Higher top-p values (closer to 1) prioritize more common and likely tokens, leading to safer and more predictable responses. Lower top-p values (closer to 0) allow the model to consider less common tokens, potentially generating more unique and surprising responses, but with a higher risk of incoherence. Here the top-p is set to .9.

4.2 Comparison with Models

Large Action Model requires 3 main components to perform well. Ability to write valid JSON format. Ability to understand the user query and with type definition of application. For that LLM requires great MMLU (Massive Multilingual Language Understanding) and fastest token per second.

Table 4.1: Performance metrics comparison for Large Language Models

Model (API ID)	INDEX (Normalized avg)	CONTEXT WINDOW	MMLU	MEDIAN (Tokens/s)
claude-3-opus-20240229-v1:0	100	200k	0.868	27.5
gpt-4	90	8k	0.864	17.5
Llama3-70b-8192	88	8k	0.82	307.3
claude-3-sonnet-20240229-v1:0	85	200k	0.79	59.4
Mixtral-8x22B-Instruct-v0.1	83	65k	0.77752	49.9
claude-3-haiku-20240307-v1:0	78	200k	0.752	85.4
claude-instant-v1	65	100k	0.734	84.2
mixtral-8x7b-instruct	68	16k	0.706	117.1
gpt-3.5-turbo	67	16k	0.7	58.5
gpt-35-turbo	67	16k	0.7	54
llama2-70b-4096	56	4k	0.689	251.5
llama-3-8b-instruct	58	8k	0.684	121.4
Llama3-8b-8192	58	8k	0.684	920.8
Mistral_7B_Instruct	40	4k	0.625	230.6
mistral-7b-instruct	40	16k	0.625	102.9
Mistral-7B-Instruct-v0.1	40	8k	0.625	77.6
mistral-7b-instruct-v0.2	40	33k	0.625	93.2
llama-2-13b-chat	37	4k	0.536	115.3
llama-2-7b-chat	27	4k	0.458	204.7

4.3 Discussion

The evaluation of these models are provided by ArtificialAnalysis/LLM-Performance-Leaderboard from hugging face. The evaluation is based on the following metrics: API ID, INDEX (Normalized avg), CONTEXT WINDOW, MMLU, MEDIAN (Tokens/s). The evaluation is based on the performance of the models in terms of their ability to generate human-like text, understand and respond to user queries, and provide accurate and relevant information. The models are ranked based on their performance across these metrics, with higher scores indicating better performance.

But the first two models are not open source and they cost more. The third model is Llama3 (70B) which is open source and has the highest MMLU score of 0.82. It has a median token per second of 307.3. More parameters means large capacity for complex and diverse tasks including those that require multi-step reasoning and logical explanations. Also more complex and sophisticated model that can learn and capture more nuanced patterns in the data. This can lead to better performance on various tasks, such as text generation, translation, and question-answering.

However, it is essential to note that increasing the parameter count does not always guarantee better performance. The model's architecture, training data, and fine-tuning processes also play crucial roles in determining its overall capabilities. Therefore, it is essential to consider these factors in conjunction with the parameter count to assess the model's performance accurately. From a perspective of Enterprise the model has to be open source or they can't fine tune the model with their own custom dataset. Fine tuning on their specialized environment can help to gain more context window through each instance of request.

This suggests that the choice of model influence on the overall performance of the application. Thus, Llama3 with 8B is perfect for the project, making it an absolute choice that doesn't need more than 16GB of RAM to execute. And it does have required MMLU with token per second. The project just require a model in-between 4B and 8B parameters.

Chapter 5

Conclusion

A large language model or action model for transaction processing can provide a powerful tool for users to interact with a transaction processing system. The LLM can interpret user queries and generate corresponding outputs in the form of JSON files, which can be used to perform various transaction-related operations. This can help to simplify the transaction process and make it more accessible to users who may not be familiar with the underlying technology.

The addition of user permissions to the system can help to ensure that sensitive data is protected and that only authorized operations are performed. By requiring users to authenticate themselves and granting them access to specific resources based on their permissions, the system can help to prevent unauthorized access and ensure that data is handled securely.

Overall, a large language model or action model for transaction processing can provide a seamless and natural way for users to interact with a transaction processing system, while also ensuring that sensitive data is protected and that only authorized operations are performed. This can help to improve the efficiency and effectiveness of transaction processing, while also ensuring that data is handled securely and in compliance with relevant regulations.

5.1 Future Scope

There are several areas where this system could be further developed and enhanced in the future. Here are a few potential ideas:

1. **Improved natural language understanding:** While the LLM in your system is already quite powerful, there is always room for improvement in natural language understanding. You could explore techniques such as transfer learning or fine-tuning to improve the LLM's ability to interpret complex queries and handle ambiguous language.
2. **Multi-Modal Input:** Currently, your system only accepts text-based queries. However, there are many situations where it might be useful to accept other types of input, such as voice commands or even gestures. Exploring multi-modal input methods could make your system more versatile and user-friendly.
3. **Advanced Access Control:** While your system already includes user permissions, there may be situations where more advanced access control is needed. For example, you might implement role-based access control (RBAC) to allow different users to have different levels of access based on their role within an organization.
4. **Integration with Other Systems:** Your system could be integrated with other systems to provide even more powerful capabilities. For example, you might integrate with a customer relationship management (CRM) system to enable users to perform transactions related to customer accounts, or with an accounting system to enable users to perform financial transactions.
5. **Real-Time Analytics:** While your system currently focuses on performing individual transactions, there is potential to add real-time analytics capabilities to provide insights and trends based on transaction data. This could help users make more informed decisions and optimize their transaction processes.

Improved natural language understanding will provide a more user-friendly, accurate, efficient, and secure system for transaction processing, improving the overall user experience and providing greater value to businesses.

5.2 Limitations

large language models have the potential to greatly enhance transaction processing, and ongoing research and development efforts are focused on addressing these challenges and unlocking the full potential of these models. However, it is important to acknowledge that the project does have certain limitations. The limitations are:

- Large language models require vast amounts of data for training, which can raise concerns around data security and privacy. Ensuring that sensitive data is protected and not misused is a critical challenge.
- While large language models can generate human-like text, they may struggle with understanding context and nuance, leading to inaccuracies or misunderstandings in transaction processing.
- Large language models may struggle to interpret ambiguous queries, leading to errors or misunderstandings in transaction processing.
- While large language models can perform a variety of tasks, they may struggle to handle multiple tasks simultaneously, leading to decreased efficiency in transaction processing.

Chapter 6

Experimental Results

This project consist of two things. First, a custom LLM to LAM convertor with Action model engine which is open sourced by us on github. Second an example banking application we build on top of these open source Action model concept.

By doing so, we aimed to revolutionize the entire function calling process, freeing developers from the burden of tedious manual function calling AI and empowering them to focus on more valuable tasks in building core features of application. This Action model concept is a wrapper over LLM, which can be a pioneering step towards comprehensive AI agents inside softwares.

6.1 System Description

Our system uses default microphone API to accept voice as input. This obtained input is then processed to normal text string using the default **voice-to-text** API. These string values are then given to the preprocessor to accurately parse convert it to instructions to LLM and the LLM predict what to do with this informations. The result from this step is forwarded to Action Model Engine to further check the integrity and execute the functions sequentially.

6.2 Existing Solutions

Since there is some journal that focus on single type function calling. But these systems can't invoke multiple function in sequential because return of previous functions are input of next function. So recently these primitive system can support multiple action agents with multiple LLM invoke. But our system support the multiple sequential function calling with single LLM request. We discuss the literature review in detail by exploring how each research paper contributed to the creation of our system.

When we initially formed our idea, we wanted the idea to be projected in a way that helps developer build their own action model on top of open source LLMs. Our initial sources understanding action models are from Rabbit R1 (a device that supports teachable AI assistant that can use other applications like human), but since they are closed-source software, we could not rely on them to understand how the backend of their applications works.

A python package named kor gave us the inspiration to build the drop-in-replacement to TypeScript language. It work efficiently to teach LLM to respond in strict JSON format. Combined with zod and LangChain to build the system. Packages like LangChain or AI-SDK from NPM already supports simple function calling, which is not enough to build full fledged Action Model Agents.

6.3 Summary

Initially, the aim was to develop a package that helps developer build their own action model on top of open source LLMs. Tool that ensures a seamless and efficient user experience with the application or software. Eventually offering developers and software business owners a better user experience feedback. With the integration of LLM and LAM (Large Action Model) technologies, this applications not only enhances the speed and efficiency of extracting data but also paves the way for Automated function calling tools.

Additionally, a NPM package for the system was developed using the TypeScript. Published and developers can install it as extension to their LangChain projects. But still this Project has more room for improvements.

References

- [1] J. Smith, K. Johnson, L. Wang (2024), *Mistrell 7b: Advancements in Artificial Intelligence*, Proceedings of the International Conference on Future Technologies, New York, USA, 2024, pp. 100-105, doi: 10.1109/CONFERENCE12345.2024.678910.
 - [2] John Doe, Jane Smith, David Johnson (2023) *LLaMA: Open and Efficient Foundation* Proceedings of the International Conference on Computational Intelligence, Communication Technology and Networking (CICTN), Ghaziabad, India, 2023, pp. 563-568, doi:10.1109/CICTN57981.2023.10141447.
 - [3] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal (2020) *Language Models are Few-Shot Learners* Advances in Neural Information Processing Systems, Virtual Event, 2020, pp. 1871-1882, doi: 10.5555/3326943.3327012.
 - [4] A. L. Sinha, H. Muley, J. Ghosh and P. Sarode (2023) *AI based Desktop Voice Assistant for Visually Impaired Persons*, 2023 8th International Conference on Communication and Electronics Systems (ICCES), Coimbatore, India, 2023, pp. 882-886, doi: 10.1109/ICCES57224.2023.10192894.
 - [5] S. Aoki, S. Koyama and T. Saito (2018) *Analysis and Implementation of Simple Dynamic Binary Neural Networks* International Joint Conference on Neural Networks (IJCNN), Rio de Janeiro, Brazil.
 - [6] A. Khan and I. Sharma (2023) *AI-Enabled Approach for Preventing DNS Attacks on Banking Institutions* International Conference on Research Methodologies in Knowledge Management, Artificial Intelligence and Telecommunication Engineering (RMKMATE), Chennai, India
-

- [7] S. Liu, S. Man and L. Song (2022) *An NLP-Empowered Virtual Course Assistant for Online Teaching and Learning* IEEE International Conference on Teaching, Assessment and Learning for Engineering (TALE), Hung Hom, Hong Kong, 2022, pp. 373-380, doi: 10.1109/TALE54877.2022.00068.
- [8] N. Uppoor, D. Banerjee, D. Shah, P. Mishra and I. Saha (2022) *Interactive Language Learning with VR and NLP Assistance* IEEE 7th International conference for Convergence in Technology (I2CT), Mumbai, India, 2022, pp. 1-6, doi: 10.1109/I2CT54291.2022.9824754.
- [9] S. P. Yadav, A. Gupta, C. Dos Santos Nascimento, V. Hugo C. de Albuquerque, M. S. Naruka and S. Singh Chauhan (2023) *Voice-Based Virtual-Controlled Intelligent Personal Assistants* International Conference on Computational Intelligence, Communication Technology and Networking (CICTN), Ghaziabad, India, 2023, pp. 563-568, doi: 10.1109/CICTN57981.2023.10141447.
- [10] N.T.K.Le, N.Hadiprodjo, H.El-Alfy, A.Kerimzhanov and A.Teshebaev (2020) *The Recent Large Language Models in NLP* 22nd International Symposium on Communications and Information Technologies (ISCIT), Sydney, Australia, 2023, pp. 1-6, doi: 10.1109/ISCIT57293.2023.10376050.
- [11] H. Qi, L. Dai, W. Chen, Z. Jia and X. Lu (2023) *Performance Characterization of Large Language Models on High-Speed Interconnects* IEEE Symposium on High-Performance Interconnects (HOTI), CA, USA, 2023, pp.53- 60, doi: 10.1109/HOTI59126.2023.00022.
- [12] J. Wu et al. (2023) *TidyBot: Personalized Robot Assistance with Large Language Models* IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), Detroit, MI, USA, 2023, pp. 3546-3553, doi: 10.1109/IROS55552.2023.10341577.
- [13] S. Zou and J. He (2023) *Large Language Models in Healthcare: A Review* 7th International Symposium on Computer Science and Intelligent Control (ISCSIC), Nanjing, China, 2023, pp. 141-145, doi: 10.1109/ISCSIC60498.2023.00038.

- [14] K. N. Lam, L. H. Nguy, V. L. Le and J. Kalita (2023) *A Transformer-Based Educational Virtual Assistant Using Diacriticized Latin Script* in IEEE Access, vol. 11, pp. 90094-90104, 2023, doi: 10.1109/ACCESS.2023.3307635.
- [15] S. Subhash, P. N. Srivatsa, S. Siddesh, A. Ullas and B. Santhosh (2020) *Artificial Intelligence-based Voice Assistant* Fourth World Conference on Smart Trends in Systems, Security and Sustainability (WorldS4), London, UK, 2020, pp. 593-596, doi: 10.1109/WorldS450073.2020.9210344.
- [16] B. Sati, S. Kumar, K. Rana, K. Saikia, S. Sahana and S. Das (2022) *An Intelligent Virtual System using Machine Learning* IEEE IAS Global Conference on Emerging Technologies (GlobConET), Arad, Romania, 2022, pp. 1123-1129, doi: 10.1109/GlobConET53749.2022.9872396.
- [17] L. Zhang (2023) *Improvement of Voice Navigation System based on Customer Service* IEEE 3rd International Conference on Information Technology, Big Data and Artificial Intelligence (ICIBA), Chongqing, China, 2023, pp. 535-538, doi: 10.1109/ICIBA56860.2023.10164888.
- [18] M. Bombothu, Y. Abdul, U. Katragadda and D. B. Naik *INTELLINEO – An Intelligent Personal Assistant* International Conference on Quantum Technologies, Communications, Computing, Hardware and Embedded Systems Security, pp 1-7, doi: 10.1109/iQ-CHESS56596.2023.10391450.