

Treball Final de Màster

Estudi: Màster en Ciència de Dades

Títol: Ajustament d'un model generatiu de llenguatge per a la creació de xatbots personalitzats per administracions públiques

Document: Memòria

Alumne: Martí Mas Fullana

Tutor: Josep Suy Franch

Tutor: Miquel Tarragona Margarit

Departament: Departament d'Informàtica, Matemàtica Aplicada i Estadística

Àrea: Intel·ligència Artificial

Convocatòria (mes/any): Setembre 2024

TREBALL FINAL DE MÀSTER

Ajustament d'un model generatiu de llenguatge per a la creació de xatbots personalitzats per administracions públiques

Autor:

Martí MAS FULLANA

Setembre 2024

Màster en Ciència de Dades

Tutors:

Josep SUY FRANCH

Miquel TARRAGONA MARGARIT

Resum

Agraïments

Per començar vull agrair molt especialment a ...

Índex

1	Introduction	1
2	Background	3
2.1	Introduction to Dialogue Systems	3
2.2	Towards Language Models	3
2.3	GPT and its Contribution	3
2.4	Retrieval Augmented Generation (RAG)	4
2.5	Applications and Benefits of RAG-based Chatbots	4
3	Objectives	5
4	Methodology	7
4.1	Data Collection and Preparation	7
4.2	RAG Implementation	7
4.3	User Interface Design	7
5	Architecture	9
5.1	Hotword/Wakeword Detection	9
5.2	Azure Services	11
5.2.1	PostgreSQL Database	11
5.2.2	GPT Model	11
5.3	Backend	12
5.3.1	Conversation Flow	12
5.4	Vector Search API	12
5.4.1	Model Routing	12
5.4.2	Information Retrieval	13
6	Results	15
	Bibliografia	17
	Flowable Diagrams	19
	Our Contributions to Open-Source	23

Índex de figures

5.1	App Architecture	9
5.2	Mel spectrogram of the audio “Hello world” (image taken from [1])	10
5.3	Example of conversation messages being routed to different models according to their complexity	13
5.4	Small To Big Retrieval (STBR) embeddings hierarchy	14
1	Top level flow diagram for the conversation	19
2	Flow diagram for the Anti DAN stage	19
3	Flow diagram for treating a scenario	20
4	Flow diagram for the answering a question using RAG (if a question exists)	20
5	Flow diagram for the scenario where we need to discover the user’s situation	20
6	Flow diagram for inferring necessary variables for the current scenario based on the conversation	21

Índex de taules

Índex

Introduction

We present a chatbot system that uses GPT (Generative Pre-trained Transformer) technology and RAG (Retrieval Augmented Generation) to provide assistance with social rights and benefits in Catalonia. Our client is the department of social rights of the Generalitat de Catalunya (*Departament de Drets Socials* or DSO).

Background

2.1 Introduction to Dialogue Systems

Dialogue systems, also known as chatbots, have experienced a significant step-change in the last few years. Initially these systems were based on predefined rules and decision trees [2, 3], limiting their capacity for understanding and answering user queries in a natural and flexible manner. These rudimentary systems, commonly referenced as rule-based chatbots, might have been enough for simple tasks, but could not have managed the full complexity and variability of natural language.

2.2 Towards Language Models

As the first machine learning-based language models appeared, such as the Sequence-to-Sequence (Seq2Seq) model [4], and more recently the transformer-based models such as GPT (Generative Pre-trained Transformer) [5, 6], the capacity of chatbots to understand and generate natural language has improved significantly. These models are trained on large datasets of text, learning the complex patterns and structures of language, and are able to generate text that is coherent and contextually relevant.

2.3 GPT and its Contribution

The GPT model [6], developed by OpenAI, has been one of the most notable advances in this field. GPT uses the transformer architecture [5], which is a type of neural network that is particularly well-suited for processing sequences of data, such as text. Its capacity for generating coherent and contextually relevant responses has been leveraged in a wide range of applications, from virtual assistance to automated content generation.

2.4 Retrieval Augmented Generation (RAG)

One of the most recent advances in the integration of language models has been the use of retrieval augmented generation (RAG) [7]. RAG combines the strengths of information retrieval from databases with the generative capacity of language models. In this context, when a user query is received, the system first retrieves relevant information from a database, and then the language model generates a coherent and precise response based on this information. This approach has been shown to improve the accuracy and relevance of the responses generated by chatbots [7].

2.5 Applications and Benefits of RAG-based Chatbots

RAG-based chatbots offer a variety of benefits compared with more traditional systems. They are able to generate responses that are more coherent and contextually relevant. In this way users are both less frustrated and more satisfied. These systems also allow the chatbots to have access to newer, more up to date information than the data the model was originally trained on, as the data provided to the information retrieval component can be updated by simply adding new entries to the database. This makes the chatbot more adaptable and flexible, and allows it to provide more accurate and relevant information to users, reducing the necessity of performing full or partial retraining of the model, which can be prohibitively expensive.

CAPÍTOL 3

Objectives

The main goal of this project is to develop an advanced chatbot system that uses GPT (Generative Pre-trained Transformer) technology and RAG (Retrieval Augmented Generation) to provide responses to user queries based off of the content of a database. This general goal can be broken down into the following specific objectives:

1. Pick an appropriate GPT Model

- Choose a GPT model that is well-suited for the task of generating responses to user queries based on the content of a database.

2. Integrate RAG Technology

- **Information Retrieval:** Develop and implement a system for retrieving relevant information from a database based on user queries.
- **Combine Retrieval and Generation:** Integrate the information retrieval system with the GPT model to generate coherent and contextually relevant responses to user queries.

3. Facilitate User-Chatbot Interaction

- **UI Design** Develop a user interface that allows users to interact with the chatbot in a natural and intuitive way.
- **UX Design** Ensure that the user experience is smooth and seamless, and that users are able to easily access the information they need.

4. Accessibility

- **Multilingual Support** Implement support for multiple languages to make the chatbot accessible to a wider range of users.
- **Accessibility Features** The system must be designed to be accessible to users with visual or motor impairments. As such, it should support voice input. The voice input feature must be able to be activated through a voice command.

5. Evaluate and Validate the System

- **User Testing** Conduct user testing to assess the usability and effectiveness of the chatbot system.
- **Results Analysis** Analyze the results of the different tests to identify areas for improvement and optimization.

Methodology

4.1 Data Collection and Preparation

- **Data Collection:** We are given by the stakeholders a set of websites documenting the laws related to social rights in Catalonia and also documenting available social benefits. These websites contain a variety of information, including the text of the laws and the social benefits, what conditions are necessary to access them, and how to apply for them.
- **Data Preparation:** We will extract the text from the websites using web scraping and convert them into a format that can be used by the information retrieval system. This will involve cleaning the text and removing any extraneous characters. This will be done using custom web scraping scripts.

4.2 RAG Implementation

- **Information Retrieval:** Develop a system capable of retrieving relevant information from a database based on user queries. This will involve creating an index of the laws related to social rights and available social benefits in Catalonia, and implementing a search algorithm that can return the most relevant laws based on the user query. In practice this will be done using the LlamaIndex Python library, which chunks the text into smaller pieces and indexes them.
- **Combining Retrieval and Generation:** Integrate the retrieval system with the GPT model to generate coherent and contextually relevant responses to user queries. This will involve passing the retrieved information to the GPT model, which will generate a response based on this information.

4.3 User Interface Design

- **UI Design:** Develop a user interface that allows users to interact with the chatbot in a natural and intuitive way. This will involve creating a chat

interface that allows users to input queries and receive responses from the chatbot.

- **UX Design:** Ensure that the user experience is smooth and seamless, and that users are able to easily access the information they need. This will involve testing the user interface with a group of users to identify any areas for improvement, as well as demoing the system to the stakeholders.

CAPÍTOL 5

Architecture

We develop an Angular Frontend conversation UI (similar to other messaging apps) that communicates with a Backend that manages the calls to the language models and the database. The Backend also handles the management of the database and the indexing of the information. The diagram of the app's architecture is shown in Figure 5.1.

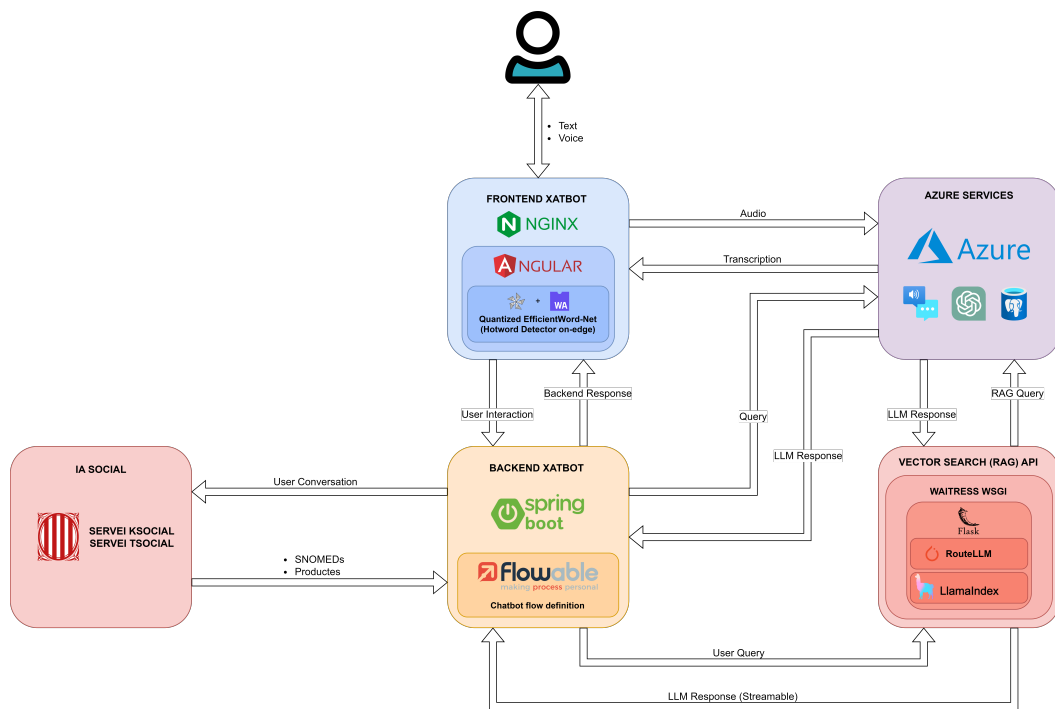


Figura 5.1: App Architecture

Over the following sections we will describe the different components of the system in more detail, and how a user's interaction is processed through the system.

5.1 Hotword/Wakeword Detection

On the frontend we implement a custom version of the EfficientWord-Net [1] hotword detector, that has been quantized [8]. To implement it we use ONNX

[9] and WebAssembly (Wasm) [10], which run directly in the browser, ensuring no data leaves the client’s machine inadvertently. The hotword detector listens to the microphone data in real-time, and when it hears the keyword “chat” or “chatbot”, it activates the microphone and starts recording. We will see how this data is processed in the 5.2 section.

The EfficientWord-Net model has been quantized to 8-bits, which has reduced the size of the model from 80MB to 20MB, and has improved the response time by 100%. This has been achieved without perceptibly sacrificing the accuracy of the hotword detection.

The hotword detector analyzes audio data in chunks of 1.5 seconds, overlapped by 0.75 seconds. The raw audio signal is first converted to a Mel spectrogram (Figure 5.2), which is then passed through a ResNet [11] model to generate semantic embedding vectors. These vectors are then compared to the embedding vectors of reference recordings of the hotword (which are prerecorded) using cosine similarity. If the similarity is above a certain threshold, the hotword is detected. We use the pretrained model as provided by the EfficientWord-Net library, with no further training.

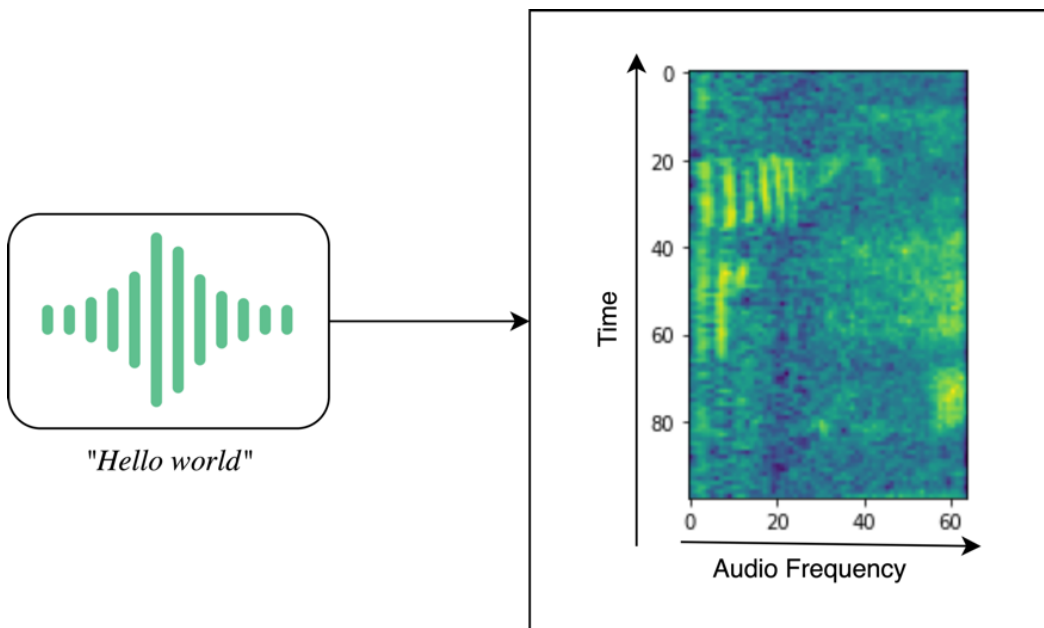


Figura 5.2: Mel spectrogram of the audio “Hello world” (image taken from [1])

Because this system needs to run on the browser, and there is no existing implementation of a Mel spectrogram converter, we implement our own converter directly in TypeScript. This converter is as close as possible to a direct translation –from Python to TypeScript– of the original code from the EfficientWord-Net library. By doing so we ensure data generated by our implementation is equiva-

lent to that generated by the original implementation. We evaluate that this is the case by making bitwise comparisons of the output of both implementations (subtracting one image from the other), and find there are no differences (all pixels in the resulting image are exactly 0).

On our reference system, which consists of a Dell Latitude 3440 laptop with a 13th Gen Intel(R) Core(TM) i5-1345U CPU, the hotword detector has a response time of around 80-100 milliseconds, which is well within the acceptable range for real-time applications.

5.2 Azure Services

5.2.0.1 Speech-to-Text

When the hotword detector activates the microphone, the audio data is sent to the Azure Speech-to-Text service. This service converts the audio data into text, which is then sent to the Backend for processing. The Azure Speech-to-Text service uses the Whisper [12] model to accurately transcribe speech into text.

5.2.1 PostgreSQL Database

The PostgreSQL database contains the text of the laws related to social rights and available social benefits in Catalonia. This data is indexed using the LlamaIndex Python library to chunk the text into smaller pieces, index them and help create semantic embeddings. When a user query is received, the information retrieval system searches the database for the most relevant laws and benefits based on the query, and returns this information to the GPT model for generation.

5.2.2 GPT Model

The GPT model is used to generate coherent and contextually relevant responses to user queries based on the information retrieved from the database. The GPT model is a transformer-based [5] language model that is trained on a large dataset of text to generate human-like responses to user queries.

We use the GPT-4o and GPT-4o Mini models, which are the latest versions of the GPT model developed by OpenAI at the time of writing.

5.3 Backend

5.3.1 Conversation Flow

The Backend implements, the conversation flow followed by the chatbot. It implements multiple stages and follows a state machine to manage the conversation. User queries at each stage are routed to the appropriate model.

Conversations flows are implemented using Flowable. Figures 1, 2, 3, 4, 5, and 6 show the flow diagrams for the different stages of the conversation.

At each stage the backend can decide to talk to one of the other components of the architecture. For example, during the RAG stage it uses our Vector Search API component to generate a response based on the information retrieved from the database. In another case, in the scenario where we are discovering the user's situation, the backend talks to the IA Social component to and asks it if, given the current conversation, there are any SNOMEDs that might apply to the user.

The backend also implements any error handling that might be necessary. The error handling is implemented directly on the flow diagram itself.

5.4 Vector Search API

As we needed to implement a few custom components all related to the RAG system, we decided to abstract them into a single separate component, the Vector Search API. This component is responsible for managing the information retrieval system, and for generating responses based on the information retrieved from the database. This component is abstract enough to be reused in chatbot systems other than the one we are developing for the DSO.

This API is implemented in Python using Flask to serve the API endpoints, and with Waitress as the WSGI server, as depicted in Figure 5.1.

The following sections describe the features that this API provides.

5.4.1 Model Routing

We use a slightly tweaked version of the fairly new RouteLLM Python library. This allows us to route a certain percentage of queries to a “stronger”, more expensive model and the rest to a “weaker”, less expensive model. With this system we are able to achieve X% of the performance of the stronger model at a fraction of the cost. The RouteLLM system uses a small BERT classifier model that decides which queries should be routed to the stronger model and which to the weaker model. This classifier was trained by the original library authors on

a dataset of human preferences augmented with synthetic data generated using GPT-4. They report good generalization performance, so we apply the system on a pair consisting of GPT-4o and GPT-4o Mini. We have also made the necessary changes to the library to make it compatible with the Azure OpenAI models, which didn't have official support.

We have translated the dataset the authors used to train the BERT classifier to Catalan, and are in the process of retraining the classifier on this new dataset, at the time of writing.

The translated training dataset consists of two parts: the translated texts and the embeddings of the texts. The translations were generated using the GPT-4o model and the embeddings using the *text-embedding-3-large* model. The total cost of generating the translations and embeddings was around 500 euros. The datasets can be found here [\[translated texts\]](#) and here [\[embeddings\]](#).

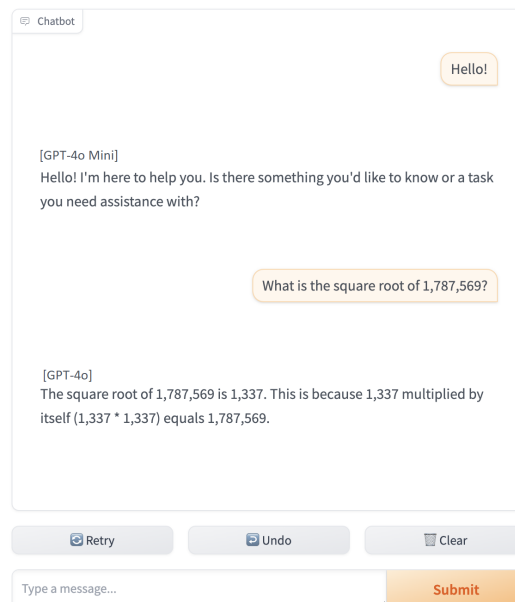


Figura 5.3: Example of conversation messages being routed to different models according to their complexity

5.4.2 Information Retrieval

We do RAG using the LlamaIndex python library.

We use two versions of RAG: normal RAG and Small To Big Retrieval (STBR).

- **Normal RAG:** This is the standard RAG algorithm. It simply creates an embedding of the user query and does cosine similarity with the embeddings of the chunks in the database to retrieve the N most relevant chunks.

In practice the database implements this with the Hierarchical Navigable Small World (HNSW) [13] algorithm, which is a fast approximate nearest neighbor search algorithm, in order to avoid having to compare the user query with all the entries in the database's table. However this implementation is invisible, as the database itself executes it.

- **Small To Big Retrieval (STBR):** This is a less common RAG algorithm that uses hierarchical embeddings of chunks, allowing us to capture both fine-grained and coarse-grained information. Figure 5.4 shows our hierarchy-of-embeddings setup. Each level's embeddings correspond to smaller and smaller chunks of the text. We have 3 levels of embeddings. This can be thought of as analogous to the way a CNN model captures both fine-grained and coarse-grained information, where deeper layers have a larger receptive field. We implement the algorithm using the `RecursiveRetriever` class from the LlamaIndex library.

The STBR algorithm is implemented in the Vector Search API and runs at the top level, so the embedding searches are still performed by the database itself, still using the HNSW algorithm. Therefore STBR is more expensive because it performs multiple database queries (one per hierarchy level), but it is also more accurate for retrieving the most relevant information.

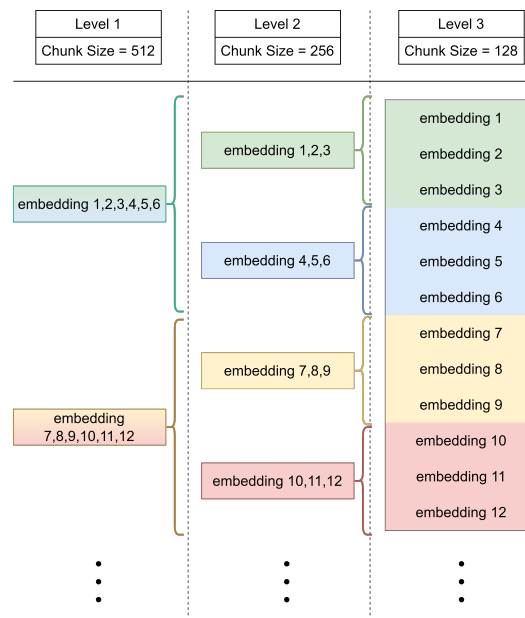


Figura 5.4: Small To Big Retrieval (STBR) embeddings hierarchy

CAPÍTOL 6

Results

RESULTSSSS

Bibliografia

- [1] R. Chidhambararajan, A. Rangapur, S. Sibi Chakkaravarthy, A. K. Cherukuri, M. V. Cruz, and S. S. Ilango, “EfficientWord-Net: An open source hotword detection engine based on few-shot learning,” *Journal of Information & Knowledge Management*, vol. 21, no. 04, p. 2250059, 2022. (Cited on pages [vii](#), [9](#) and [10](#).)
- [2] J. Weizenbaum, “ELIZA: A computer program for the study of natural language communication between man and machine,” *Commun. ACM*, vol. 9, p. 36–45, jan 1966. (Cited on page [3](#).)
- [3] B. Abushawar and E. Atwell, “ALICE chatbot: Trials and outputs,” *Computación y Sistemas*, vol. 19, 12 2015. (Cited on page [3](#).)
- [4] I. Sutskever, O. Vinyals, and Q. V. Le, “Sequence to sequence learning with neural networks,” 2014. (Cited on page [3](#).)
- [5] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” 2023. (Cited on pages [3](#) and [11](#).)
- [6] A. Radford and K. Narasimhan, “Improving language understanding by generative pre-training,” 2018. (Cited on page [3](#).)
- [7] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W. tau Yih, T. Rocktäschel, S. Riedel, and D. Kiela, “Retrieval-augmented generation for knowledge-intensive NLP tasks,” 2021. (Cited on page [4](#).)
- [8] J. Zhang, Y. Zhou, and R. Saab, “Post-training quantization for neural networks with provable guarantees,” 2023. (Cited on page [9](#).)
- [9] ONNX Community, “ONNX: Open neural network exchange.” <https://onnx.ai>, 2024. <https://onnx.ai>. (Cited on page [10](#).)
- [10] World Wide Web Consortium (W3C), “WebAssembly: A binary instruction format for a stack-based virtual machine.” <https://webassembly.org/>, 2024. <https://webassembly.org/>. (Cited on page [10](#).)
- [11] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” 2015. (Cited on page [10](#).)

- [12] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, “Robust speech recognition via large-scale weak supervision,” 2022. (Cited on page [11](#).)
- [13] Y. A. Malkov and D. A. Yashunin, “Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs,” 2018. (Cited on page [14](#).)

Flowable Diagrams

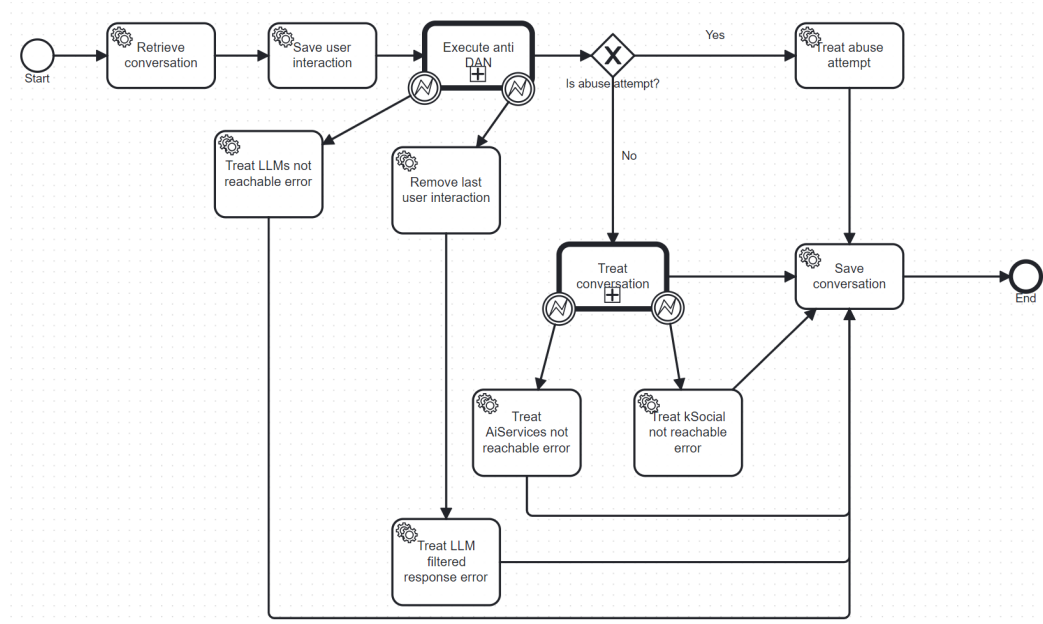


Figure 1: Top level flow diagram for the conversation

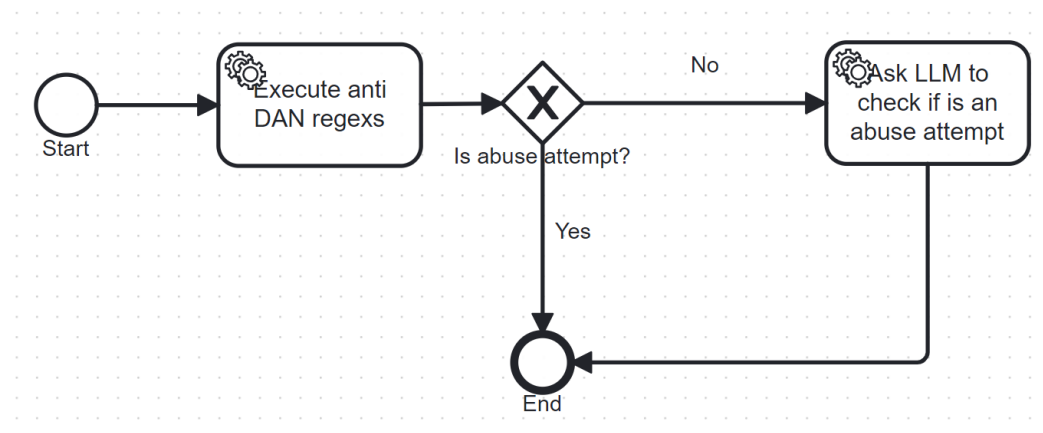


Figure 2: Flow diagram for the Anti DAN stage

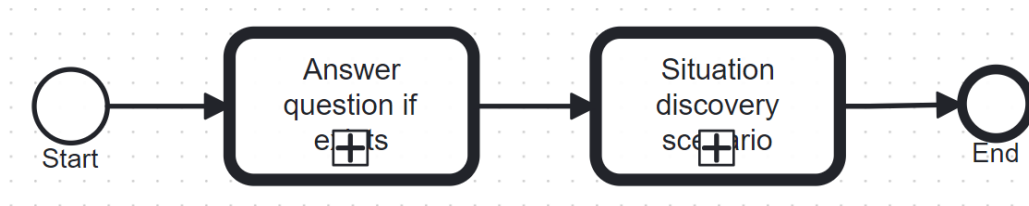


Figura 3: Flow diagram for treating a scenario

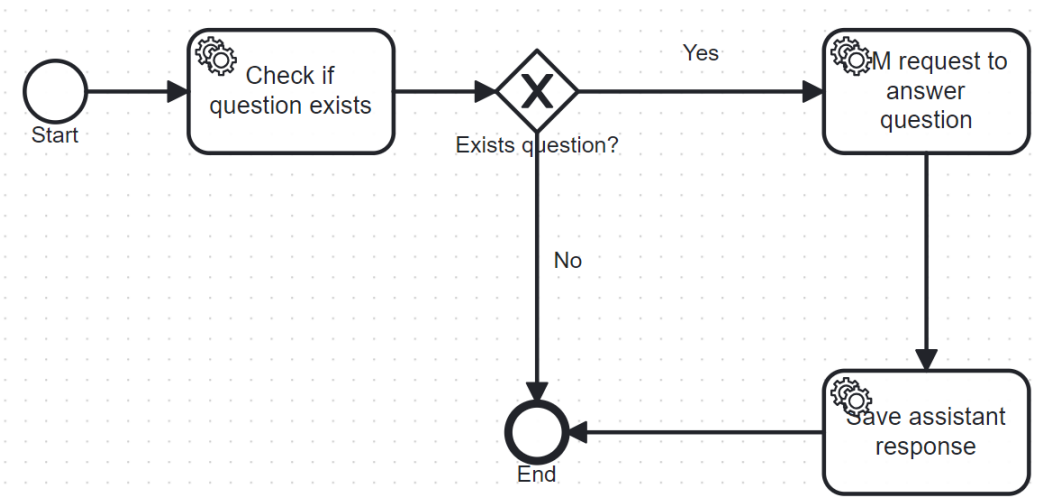


Figura 4: Flow diagram for the answering a question using RAG (if a question exists)

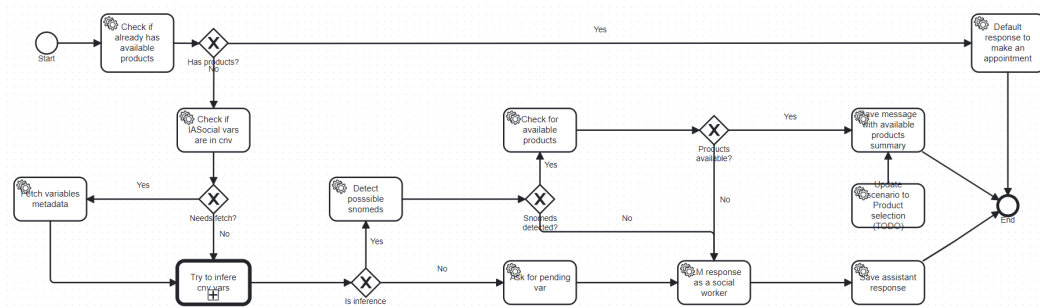


Figura 5: Flow diagram for the scenario where we need to discover the user's situation

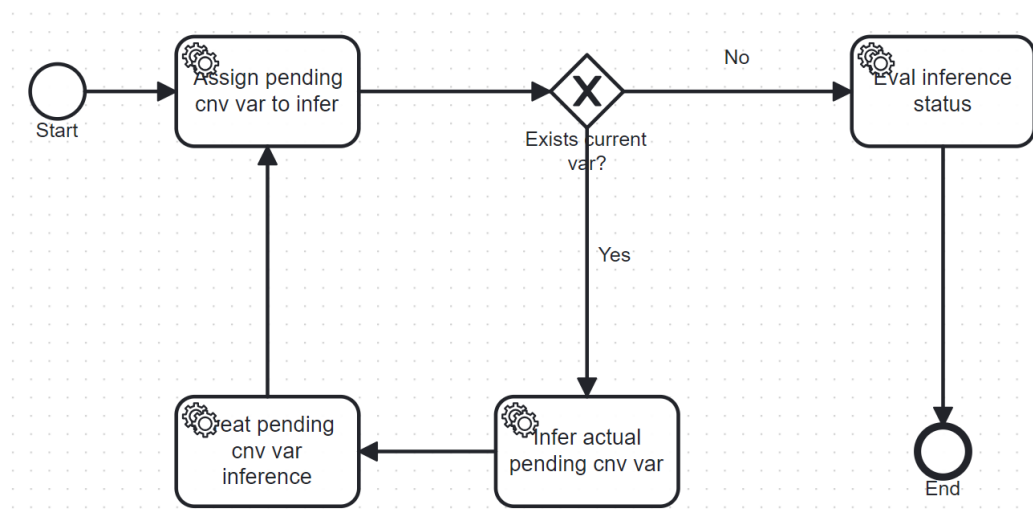


Figura 6: Flow diagram for inferring necessary variables for the current scenario based on the conversation

Our Contributions to Open-Source

During the development of this project we have iterated and made changes to a few existing tools and components. Here we link to our forks of the repositories.

- EfficientWord-Net: <https://github.com/SupremeLobster/EfficientWord-Net>
- Our fork of the EfficientWord-Net library, which has been modified to reduce response time, client machine requirements, and the size of downloaded files, through Quantization methods.
- RouteLLM: <https://github.com/SupremeLobster/RouteLLM> - Our fork of the RouteLLM library, which allows us to route a certain percentage of queries to a “stronger”, more expensive model and the rest to a “weaker”, less expensive model. Our changes are in the branch “feature/support-azure-openai-embeddings”. These changes made the library compatible with the Azure OpenAI models, which didn’t have official support.