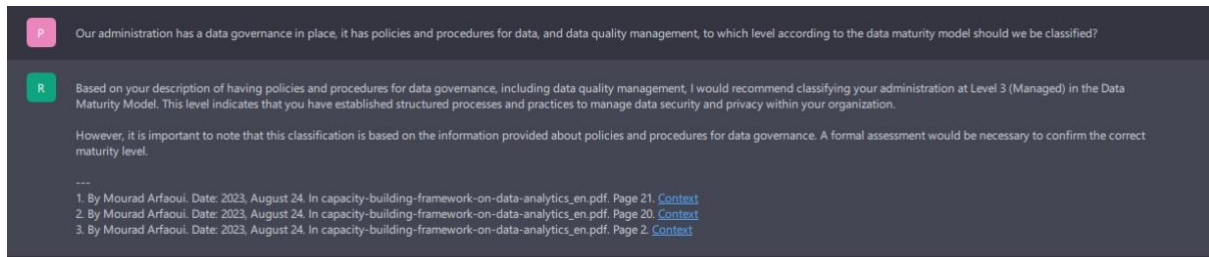# How to create a personal, secure Chatbot assistant on Customs Data

Are you interested in evaluating your organization's data analytics maturity based on the WCO's proposed model, and keen to enhance your capabilities? As demonstrated, the chatbot trained on WCO data can offer these insights and more, tailored specifically to your domain and requirements. For instance, it can conduct diagnostic assessments for you based on your specific queries.



ChatGPT has garnered significant attention for its wide range of functionalities, including answering questions, writing emails or essays, and serving as a personal or coding assistant. The tool rapidly gained popularity, attracting over a million users within a week of its launch, and has been used by over 100 million people globally. It has also amassed over 9 billion all-time page visits. The World Customs Organization (WCO) is interested in exploring the use of ChatGPT and other Large Language Models in the Customs sector, provided that data security and privacy can be ensured.

OpenAI has stated on their website that ChatGPT users now have the option to disable chat history, which limits the data retained for training the models. When this feature is active, conversations are kept for only 30 days and are reviewed solely for abuse monitoring before permanent deletion. However, there are reports suggesting that ChatGPT still collects a substantial amount of user data, including conversations, location details, and IP addresses, raising concerns over data security and privacy.

Given these concerns, sharing confidential information from the WCO or Customs Organizations with ChatGPT poses a notable risk. Does this mean that Customs agencies should avoid using such emerging technologies for domain-specific tasks? Not necessarily. The rise of ChatGPT has also catalyzed the development of many open-source alternatives. These open-source Large Language Models and chatbots offer the option to keep user data private and allow organizations to integrate pre-trained models with their internal resources, thereby balancing innovation with security.

Before delving into the methodologies, let's first clarify some of the technical terminology commonly encountered in the realm of chatbots.

**Large Language Models (LLMs):** LLMs are sophisticated machine learning algorithms capable of understanding and generating text that closely mimics human language. Trained on extensive datasets often sourced from the internet, these models analyze linguistic patterns

and structures to carry out tasks like summarization, translation, text prediction, and content generation.

**GPT4All:** GPT4All is an open-source platform designed to facilitate the integration of Large Language Models into various applications without the need for costly subscriptions to specialized hardware or platforms. It enables individuals and organizations to train, customize, and deploy highly capable language models on conventional computing hardware. Its ultimate aim is to offer an instruction-optimized, assistant-style language model that is accessible, extendable, and free for both individual users and enterprises.

**LangChain:** LangChain is a Python-based library that serves as a comprehensive toolkit for building applications that leverage language models. Whether you aim to develop a chatbot, translation service, or automated content generator, LangChain provides the essential components to simplify the development process.

**API:** An Application Programming Interface (API) provides a list of things that a software program can do, like sending an email, looking up a location, or saving a file. When you, or another software, make a specific request from that list, the system carries out the task and returns the result. This allows different software programs to work together easily.

In the following, two methods with distinct advantages and drawbacks on open-source tools application for customized safe to use Chatbots will be introduced.

## 1. Using Python libraries to manually use open Source LLM

The first approach involves using open-source libraries, which provide convenient access to specialized Python libraries.

Python libraries are collections of pre-written code that streamline programming tasks. They enable users to perform complex operations without having to write the code from the ground up, thereby increasing efficiency.

Key libraries employed in this scenario include LangChain, GPT4All, PyPDF, and ChromaDB.

LangChain allows for easy integration with an open-source Large Language Model (LLM) named GPT4All. This model can be downloaded and run locally on your own computer, eliminating the need to connect to external servers via an API, as is required with paid models like GPT-4. Given the data security concerns mentioned earlier, running the model locally guarantees that data will not be exposed to third parties.

For those interested in the technical details, the source code is available for review. However, for the purpose of this discussion, we will focus solely on the outcomes of the implementation.

In our case, we used a publicly available PDF of the Data Analytics Capacity Building Framework to train the LLM. The model was trained to understand the contents of the document. After the training phase, we tasked the model with summarizing the document.

The model produced an accurate summary, although it is worth noting that the execution time was relatively slow.

The reason for this can be attributed to multiple factors, here are some examples:

- **Computation Power Limitations:** The code is executed on the free tier of Google Colaboratory, in alignment with the goal of making the project accessible to a wide range of users. It's important to note that this environment provides a limited 13GB of RAM. Large Language Models typically demand more computational power, which may contribute to delays in response time.
- **Document Size Limitations:** The size of the PDF used for training can also impact processing time. In this case, the Data Analytics Capacity Building Framework document spans 66 pages. Parsing and training on such a document using the PyPDF library may result in extended delays.
- **Library Efficiency:** The performance can also be influenced by the efficiency of the LangChain and GPT4All libraries themselves. These open-source solutions may not be as finely tuned for specific tasks as their commercial counterparts, possibly affecting overall performance.

Alongside the performance issue we can find that if we want to increase performance by utilizing Open AI API, we would be facing the security issues.

Still, there are advantages to be had by using this model, if proper optimization can be had:

As we are writing our own source code, there is the ability to customize the code to own needs, and flexibility to fine-tune parameters to improve precision and performance. Also we can find lots of documentation online on the libraries used and alternate methods of implementation and optimization.

In addition to performance limitations, shifting to OpenAI's API to improve speed would introduce data security concerns.

However, there are distinct advantages to using this approach if proper optimizations can be implemented:
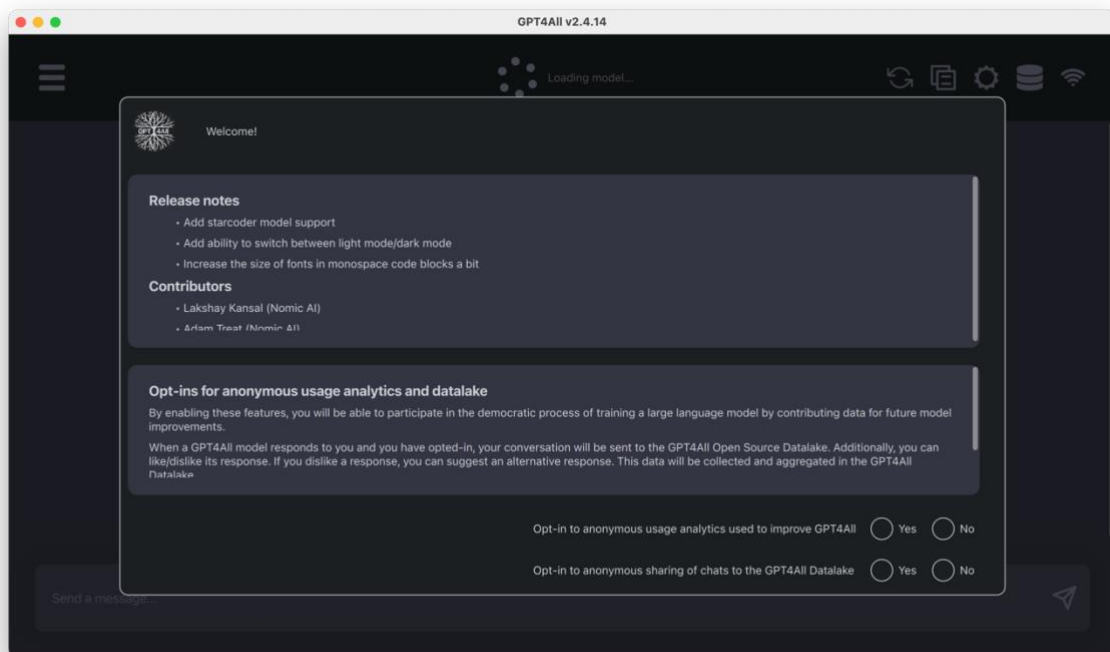
By writing our own source code, we gain the flexibility to tailor the code to meet specific requirements and fine-tune various parameters to enhance both precision and performance. Furthermore, an abundance of online documentation is available for the libraries in use, offering alternative strategies for implementation and optimization.
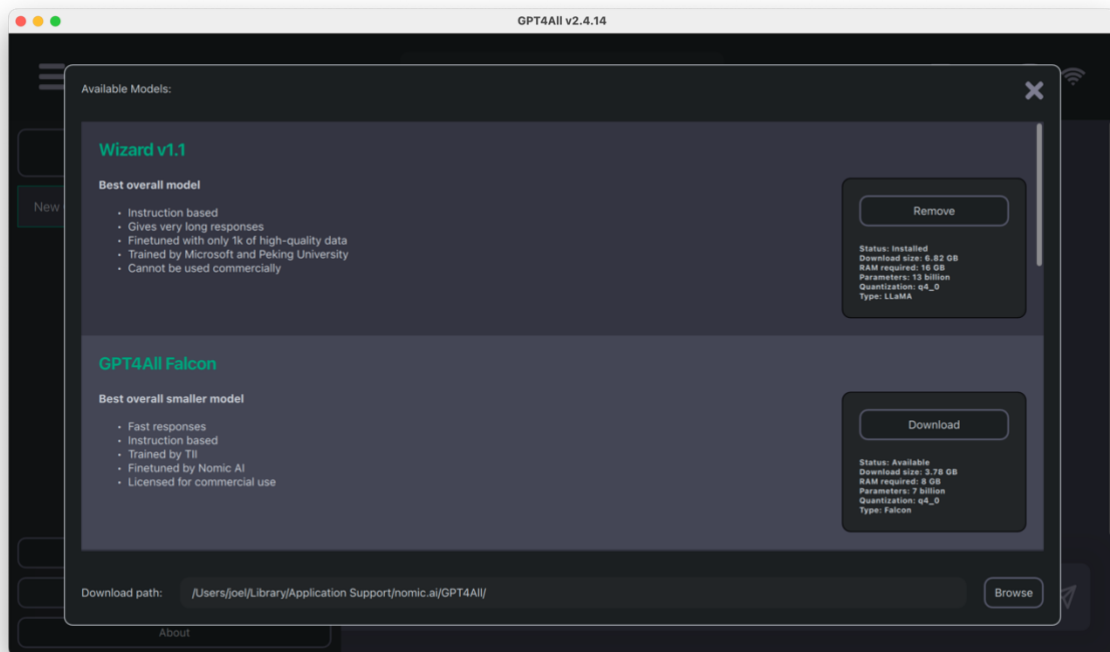
## 2. Using GPT4All Application

We already introduced the GPT4All LLM, but there is also a tool with the same name, that is a free-to-use, locally running, privacy-aware chatbot. No GPU or internet required. It can be installed on windows, MAC and linux. Much like the previous approach this can be downloaded and used locally to bypass any security and privacy concerns. Another Advantage is that it is a No-Code solution, meaning everyone can just follow some simple steps to make use of this tool:

Here is a short tutorial on how to use this tool to train models on your data and ask questions about your files:
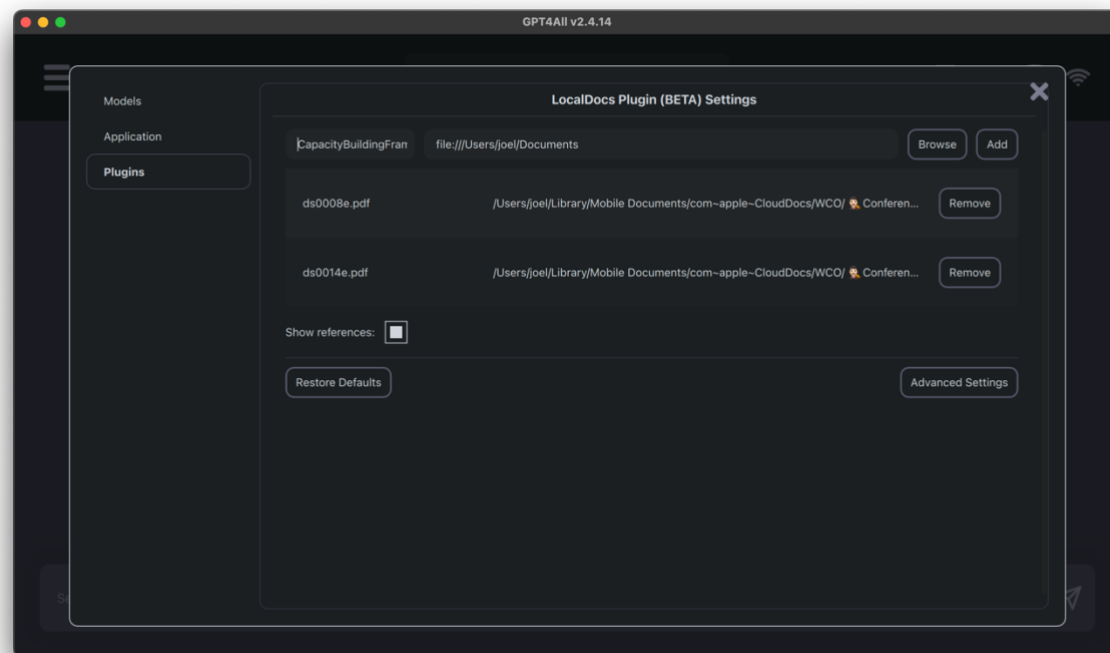
1. Install GPT4All from their website for your Operating System.
2. If you prefer not to share any data, make sure to select 'No' for both opt-in options that appear on the pop-up screen.
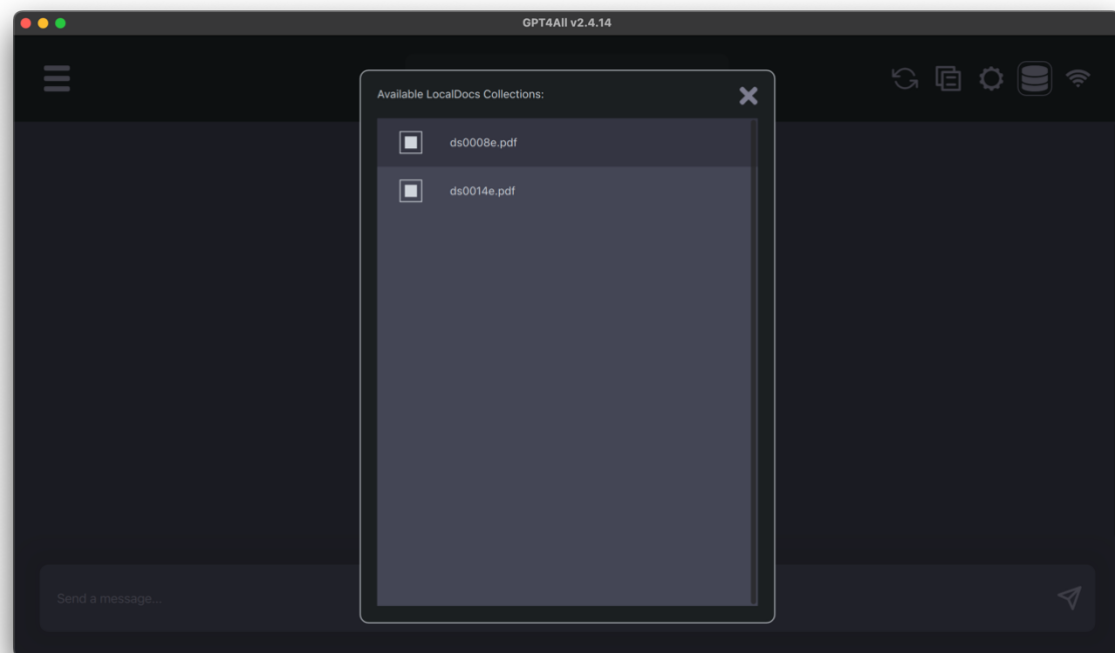


3. Next, you'll need to download a Large Language Model (LLM) to use. Free options include the Wizard or the GPT4All Falcon Model. You can also connect to OpenAI's GPT Model, though this would require data sharing via the API. Each model comes with detailed descriptions, so take a moment to read through these to find the model that best suits your needs. If you have ample storage space on your computer, I recommend downloading Wizard v1.1 by clicking the 'Download' button.

4. Next, click on the settings icon and navigate to the 'Plug-ins' section. Here, you can browse and add various types of files, including PDFs, Word documents, and text files, among others. The process is slightly intricate: first, click 'Browse' to navigate to the folder where your desired document is stored. Then, enter the name of the document in the 'Collection name' field.
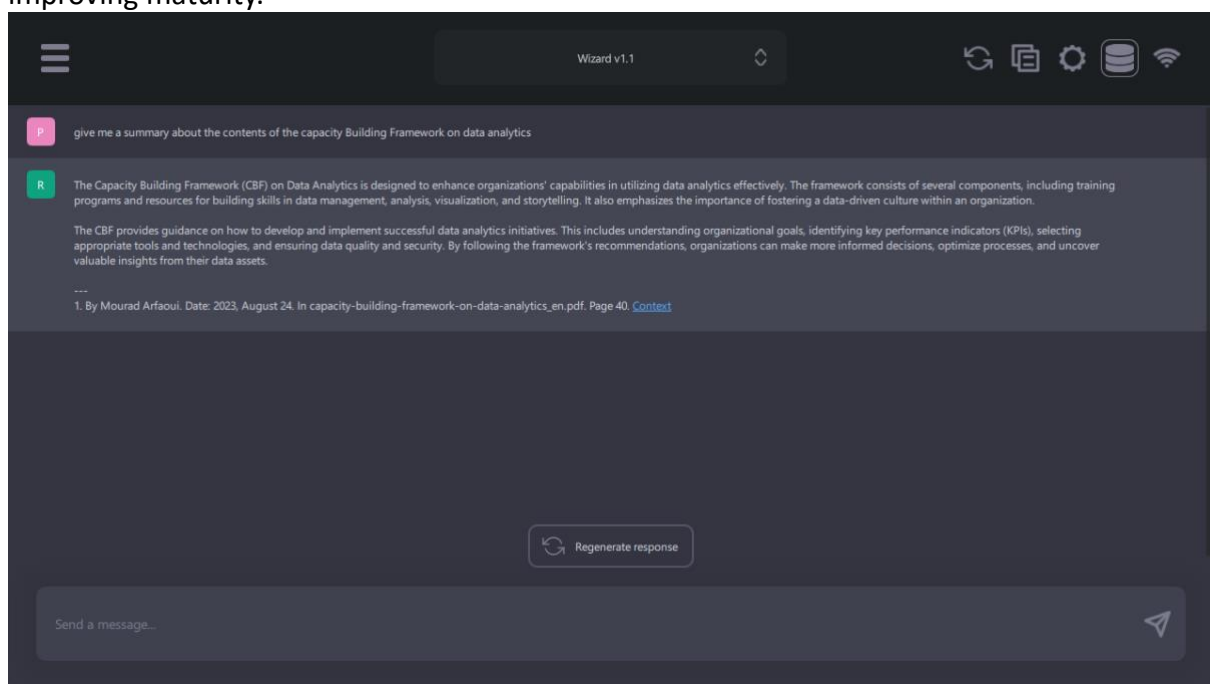
5. After completing these steps, you can close the pop-up window and click on the database icon. This will display the documents you've added, allowing you to select the specific files you wish to use.
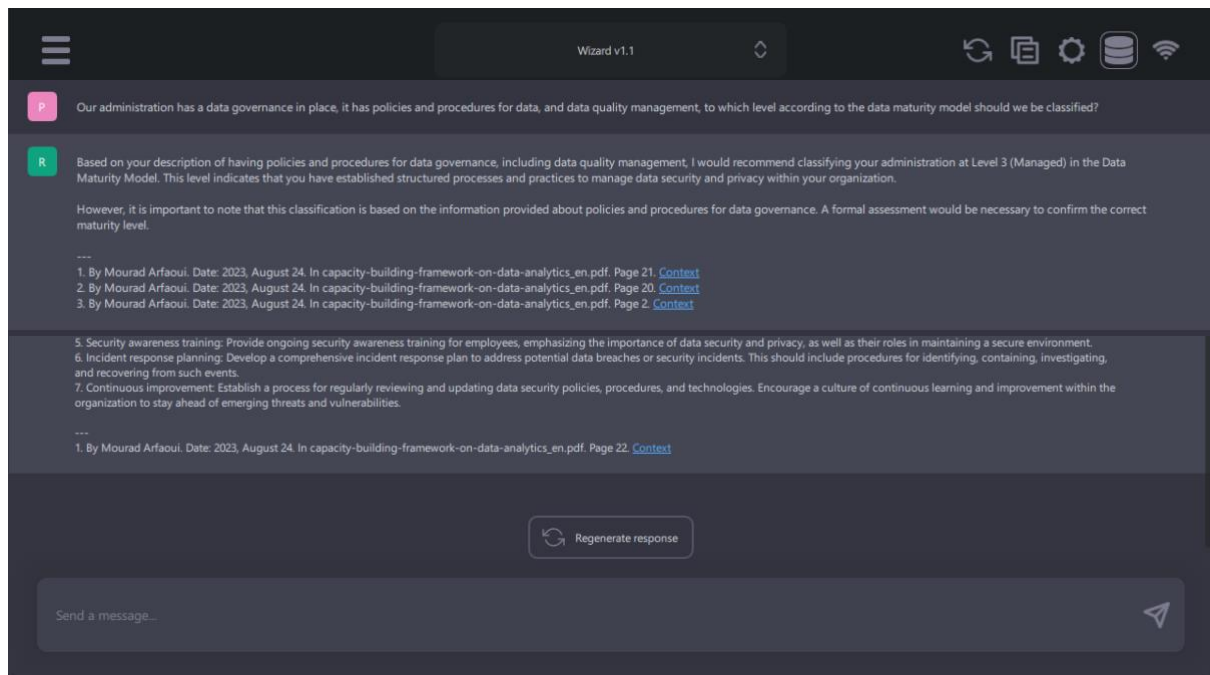


Now you're all set! To start asking questions, simply close the pop-up window and type your query into the message box. You're ready to go!

This tool you will find to be very helpful, here are some examples of how the tool helped me with assessing the maturity of my organization as well as giving me recommendations for improving maturity.
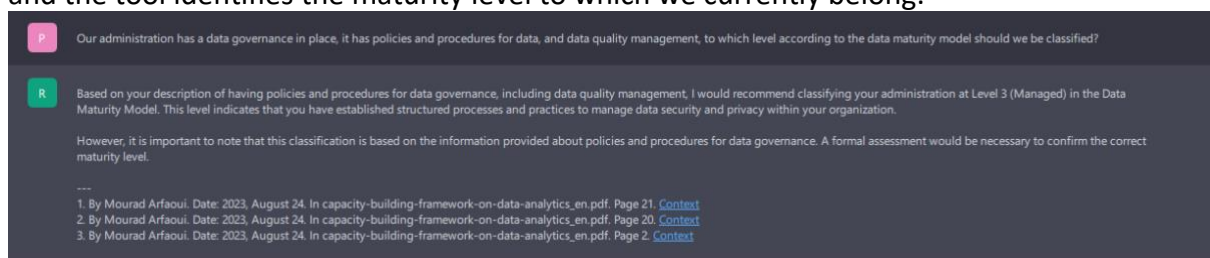
Here it gives me a summary on the Capacity Building Framework Document. Here it only takes a few seconds to produce a response.



In this section, it provides specific recommendations for advancing from the second to the third maturity level.

In this query, we outline the scope of our organization's efforts in Data Analytics Strategy and the tool identifies the maturity level to which we currently belong.



While this tool has its limitations—namely, it doesn't allow users to fine-tune its implementation and restricts them to its built-in features—it's worth noting that this No-Code solution offers the most user-friendly and effective way for Customs Officers to experiment and gain practical benefits for their day-to-day tasks.

In conclusion, it's crucial to highlight that the landscape of Large Language Models and chatbots is continuously evolving. During the course of our study, new models and methods were being developed on a weekly basis. Given this rapid pace of innovation, the insights presented here could become outdated in just a few months. Therefore, it's advisable for those interested to stay updated on the latest trends in the field of LLMs and chatbots.