

# **SISTEMAS COMPUTACIONAIS AVANÇADOS (SISTCA)**

## **ADVANCED COMPUTING SYSTEMS**

LICENCIATURA EM ENGENHARIA ELECTROTÉCNICA E DE  
COMPUTADORES

INSTITUTO SUPERIOR DE ENGENHARIA DO PORTO  
POLITÉCNICO DO PORTO

**Lab Classes Script:**

**Local Chat Agent**

## Table of Contents

1.Introduction .....	2
1.1 Scientific/Technological Context .....	2
1.2 Motivation for this topic/lab script .....	2
1.3 Objectives (of this lab script) .....	3
1.4 Structure (of this lab script).....	3
2 Theoretical (scientific/technological background) .....	3
2.1 Main theoretical concepts/terminology .....	3
2.2 The State of the Art.....	4
3. Outlook of the technology .....	5
3.1 Brief Overview of the TinyLlama-1.1B Model .....	5
3.2 Preparing the Environment .....	5
3.3 Hugging Face Access Token .....	6
3.4 Download and cache the model in drive D .....	8
3.5 Verify the installation .....	9
4.Ways to Interact with the Model .....	9
4.1 Exercise 1: Interaction via command line interface .....	10
4.2 Exercise 2: Interaction via Telegram Bot .....	12
4.3 Exercise 3: Interaction via Web Interface .....	15
4.4 Expansion and Integration Potential.....	17
5. LLaMA 3B v2 Model and Tiny LLaMA 1.1B Model .....	18
5.1 Model Comparison .....	18
6.References .....	19

**List of figures and tables**

Table 1- Comparison of various LLMs currently in use ..... 4

Table 2- Comparison between LLaMA 3B v2 and Tiny LLaMA 1.1B models .....18

Figure 1 - Flowchart of project stages ..... 5

Figure 2 - Telegram Interface.....12

Figure 3 - Token Access .....13

Figure 4 - Example of Interaction .....15

Figure 5 - Example of Web interaction .....17

# **1.Introduction**

## **1.1 Scientific/Technological Context**

Currently, artificial intelligence, more specifically machine learning, has been advancing at an exponential pace in recent years and has had a significant impact on different areas, such as industry, education and our society. The advancement of different language models that are based on deep neural networks, particularly those that use the Transformer architecture, is one of the most significant advances, as the models can understand and create text with a high level of cohesion and context, which has completely changed natural language processing.

Transformers work using different types of mechanisms that recognize connections between words in different sentences, regardless of their separation. This ability is useful for different types of tasks, such as natural language production, text summarization, automatic translation and answering different questions. Furthermore, thanks to the increase in computing power, we have been able to train increasingly complex and powerful models.

Thanks to improvements in model optimization and update approaches, new artificial intelligence systems can now operate on different devices with lower energy consumption, such as smartphones or personal servers. Thanks to this new technology, various types of tools have become very useful in assisting programming, such as interactive chatbots and virtual assistants.

On a technical and ethical level, the advancement of artificial intelligence has also brought about different new questions and difficulties, which has led the scientific and technological communities to pay increasing attention to issues such as data vision, information privacy and environmental impact. Due to these new difficulties and challenges, constant efforts are being made to improve the accessibility, security and transparency of these systems, so that everyone can truly benefit from their use.

## **1.2 Motivation for this topic/lab script**

We chose the Cloud AI Agents theme because of the great increase in the importance of artificial intelligence in our daily lives, especially in the area of natural language processing. Nowadays, the ability of computers to understand and create text in a cohesive way is transforming the way we interact with this technology in different areas, from education to health and even entertainment. With recent technological advances, it has become possible to study, experiment and test with different language models in an accessible way even with limited resources.

With this project we have the opportunity to deepen the knowledge about how these models work in practice, investigating different concepts, such as neural networks, model training and language interpretation. Our main motivation is, therefore, to learn in a practical and critical way how artificial intelligence is shaping the future of our communication and technology, while developing different types of technical skills.

### **1.3 Objectives (of this lab script)**

The main objective of this project is to analyze the installation, configuration and practical use of different open sources language models, with a special focus on the TinyLlama 1.1B and LLaMA 3B v2 models. With this project, we intend to explore how these models can be executed locally, even when machines have limited resources, and demonstrate how we can interact with them in different ways, such as through command lines and graphical interfaces.

In addition, the report aims to demonstrate how we can make the most of the capabilities of these models in natural language processing tasks, showing their potential in real-world contexts, such as virtual assistants, educational tools and prototypes of conversational applications.

### **1.4 Structure (of this lab script)**

The report is structured in several sections. It begins with a brief introduction to the topic, explaining the importance of artificial intelligence in natural language processing. A theoretical component is then presented, covering different types of concepts, such as the Transformer architecture, LLaMA models and the tools used.

The process of installing and configuring the TinyLlama 1.1B and LLaMA 3B v2 models for local execution is then described. It then presents the implementation of different interfaces for interacting with the models, via the command line, through a Telegram bot and through a web interface. Finally, practical exercises are presented and the main conclusions of the project are summarized.

## **2 Theoretical (scientific/technological background)**

### **2.1 Main theoretical concepts/terminology**

To understand the project proposal we have to explain different theoretical concepts that are necessary to understand the development of language models. As a first topic, we have to know what language models are, which are systems based on neural networks with the function of understanding and generating language in a natural way. Language models will be able to perform their tasks such as translation and text generation by learning large-volume statistical patterns.

Next we have the concept of transformer, which is an architecture that is based on attention mechanisms, this architecture will allow us to have parallel processing of tokens, which will guarantee us greater performance and efficiency. The parameters of a model are adjustable units and are related to the size and complexity of a model.

In a language model, this can be either open source or proprietary depending on its license, depending on the license it defines the rights of use and modification

of the model. When the model is open source there is transparency and free customization, but when it is proprietary the license restricts access.

Finally, a very important concept is the model interaction interface, which are the ways the user can communicate with the language model.

An interaction interface is the command line that will allow the user to interact with the model through the terminal. Another interface is the Telegram bot that will allow the user to interact through the Telegram platform, allowing users to communicate in a familiar and accessible environment. Finally, the web interface is a graphical environment accessible through the browser.

## 2.2 The State of the Art

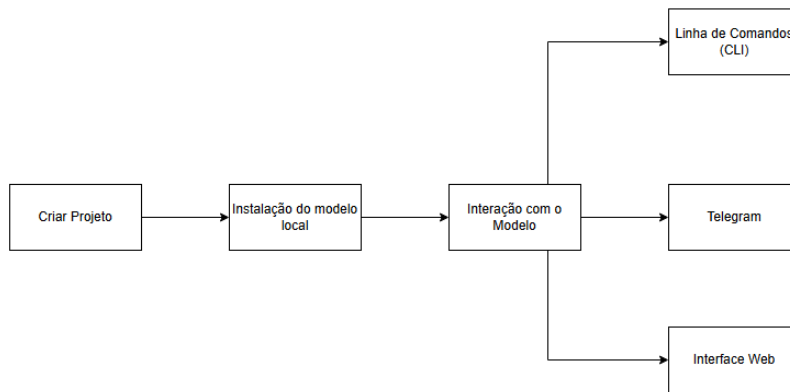
With the development of language models and the way they have revolutionized the field of natural language processing, there have been several advances in various tasks, such as text generation, automatic translation and even dialogue with users. Due to the importance of these factors, choosing the right model has become crucial.

To choose the right language model, we need to consider several factors, such as its computational capacity, its language, the complexity of the task and the type of license. Table 1 below summarizes and compares the different factors of several LLMs widely used today.

Modelo	Número de parâmetros	Idiomas	Requisitos Hardware	Licença	Treinado por
LLaMA 3B V2	3 Bilhões	Multilíngue	Médio/Alto	Open	Meta AI
TinyLlama 1.1B	1.1 Bilhões	Multilíngue	Baixo	Open	Comunidade
GPT-2	1.5 Bilhões	Inglês	Médio	Open	Open AI
GPT-3.5	6.7 Bilhões	Inglês/Multilíngue	Alto	Proprietária	Open AI
Mistral 7B	7 Bilhões	Multilíngue	Alto	Open	Mistral AI
BERT Base	110 Milhões	Inglês	Baixo	Open	Google
Falcon	40 Bilhões	Multilíngue	Alto	Open	Technology Innovation Institute
LLaMA 2	70 Bilhões	Multilíngue	Muito Alto	Open	Meta AI

**Table 1- Comparison of various LLMs currently in use**

As for the workflow that will exist in the project, it starts with the creation of the project and then its installation. After the installation of the model, the project will present different ways of interacting with the models. We can see the steps of the workflow in the figure represented below.



**Figure 1 - Flowchart of project stages**

## 3. Outlook of the technology

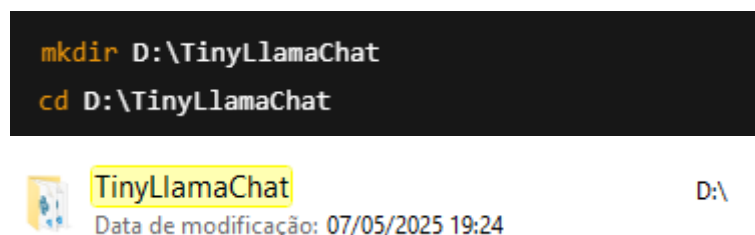
### 3.1 Brief Overview of the TinyLlama-1.1B Model

TinyLlama-1.1 B-Chat-v1.0 is a lightweight, open-source language model with approximately 1.1 billion parameters. It is optimized for chat-based tasks and balances performance and hardware efficiency well. Because of its compact size, it can run locally on personal machines. The model is available for free via the Hugging Face platform.

### 3.2 Preparing the Environment

- **Step 1: Create the working directory**

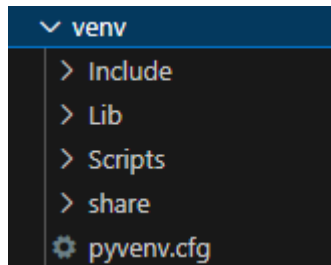
Start by creating the folder where everything related to the model will be set up.



- **Step 2: Create a virtual environment**

To isolate the project dependencies, create a virtual Python environment:

```
python -m venv venv
```



- **Step 3: Activate the virtual environment**

```
.\venv\Scripts\activate  
(venv) D:\TinyLlamaChat>
```

- **Step 4: Install Python dependencies**

Install the required Python packages:

```
pip install torch transformers
```

### 3.3 Hugging Face Access Token

While TinyLlama-1.1 B-Chat-v1.0 is public and does not require authentication, other large or gated models (such as Meta LLaMA or Mistral) require an access token and user permissions. Below are the steps to obtain and use such a token.

- **Step 1: If you don't already have an account, sign up for free:**

<https://huggingface.co/join>

- **Step 2: Generate an access token**

2.1 Go to your account settings:

<https://huggingface.co/settings/tokens>

2.2 Click on "**New token**".

2.3 Choose a name for the token (e.g., model\_access\_token).



2.4 Select the required permission level:

- read: to download models (recommended).
- write: to push or edit repositories.
- fine-grained: to allow access to specific resources.

2.5 Click "**Generate a token**" and copy the generated token (hf\_...).

<https://huggingface.co/docs/hub/en/security-tokens>

- **Step 3: Authenticate using the token**

3.1 Via Terminal:

```
huggingface-cli login
```

You'll be prompted to paste your token. Once authenticated, it will be stored locally.

3.2 Via Python

```
from huggingface_hub import login  
  
login(token="hf_...")
```

3.3 Using an environment variable

```
set HF_TOKEN=hf_...      # On Windows  
export HF_TOKEN=hf_...   # On Linux/macOS
```

- **Step 4: Request access for gated models**

Some models, like **Meta LLaMA**, require approval:

4.1 Visit the model page, for example:

<https://huggingface.co/meta-llama/Llama-2-7b-hf>

4.2 Click "**Request Access**" and accept the usage terms.

4.3 Wait for approval by the model administrators.

- **Step 5: Use the token to download the model**

Once access is granted, load the model with the token:

```
from transformers import AutoModelForCausalLM, AutoTokenizer

token = "hf_..."

tokenizer = AutoTokenizer.from_pretrained("meta-llama/Llama-2-7b-hf", token=token)
model = AutoModelForCausalLM.from_pretrained("meta-llama/Llama-2-7b-hf", token=token)
```

### 3.4 Download and cache the model in drive D

Script: download\_model.py

Now we will write a Python script to download and save the model in D:\TinyLlamaChat\hf\_cache.

```

1 # install_model.py
2 import os
3 from transformers import AutoTokenizer, AutoModelForCausalLM
4
5 # Desativa aviso de symlinks
6 os.environ["HF_HUB_DISABLE_SYMLINKS_WARNING"] = "1"
7
8 # ID do modelo no HuggingFace
9 model_id = "TinyLlama/TinyLlama-1.1B-Chat-v1.0"
10 cache_dir = "D:/TinyLlamaChat/hf_cache"
11
12 print("A descarregar o modelo...")
13
14 # Download do modelo e tokenizer
15 tokenizer = AutoTokenizer.from_pretrained(model_id, cache_dir=cache_dir)
16 model = AutoModelForCausalLM.from_pretrained(model_id, cache_dir=cache_dir)
17
18 print("Modelo instalado com sucesso em:", cache_dir)

```

```
(venv) D:\TinyLlamaChat>python download_model.py
A fazer download do modelo e tokenizer...
```

```
tokenizer.model: 100%|██████████| 500k/500k [00:00<00:00, 3.69MB/s]
tokenizer.json: 100%|██████████| 1.84M/1.84M [00:00<00:00, 2.76MB/s]
special_tokens_map.json: 100%|██████████| 551/551 [00:00<?, ?B/s]
config.json: 100%|██████████| 608/608 [00:00<?, ?B/s]
model.safetensors: 100%|██████████| 2.20G/2.20G [10:37<00:00, 3.45MB/s]
generation_config.json: 100%|██████████| 124/124 [00:00<?, ?B/s]
Download concluído.
```

### 3.5 Verify the installation

Script: verify\_model.py

Create a script called verify\_model.py to check whether the model has been properly installed.

```
verify_model.py > ...
1  # verify_model.py
2  import os
3  from transformers import AutoTokenizer, AutoModelForCausalLM
4
5  os.environ["HF_HUB_DISABLE_SYMLINKS_WARNING"] = "1"
6
7  model_id = "TinyLlama/TinyLlama-1.1B-Chat-v1.0"
8  cache_dir = "D:/TinyLlamaChat/hf_cache"
9
10 try:
11     tokenizer = AutoTokenizer.from_pretrained(model_id, cache_dir=cache_dir)
12     model = AutoModelForCausalLM.from_pretrained(model_id, cache_dir=cache_dir)
13     print("✅ Modelo carregado com sucesso a partir do disco D.")
14 except Exception as e:
15     print("❌ Erro ao carregar o modelo:", str(e))
```

```
(venv) D:\TinyLlamaChat>python verify_model.py
✅ Modelo carregado com sucesso a partir do disco D.
```

Next, a script (chat.py) will be created to test the model's functionality. However, this process is part of the model's practical usage and will be further detailed in **Section 5** of this report.

## 4. Ways to Interact with the Model

After the local installation of the model (described in Section 4), there are several ways to interact with it. These methods facilitate testing, application integration, and practical use, adapting to different user profiles, from technical users to end users. Each interaction method has its own characteristics and advantages, allowing the model to be used in various contexts, according to the needs of the project. In this project, although the step-by-step guide was done using the TinyLlama 1.1B model, the interactions described in this section were carried out with the Open LLaMA 3B-V2 model, which offers higher-quality responses and allows for the generation of more complete examples.

## 4.1 Exercise 1: Interaction via command line interface

- Description

Interaction via the command line interface (CLI) allows sending questions or prompts directly to the model and receiving the corresponding responses in the terminal. This is a quick and simple way to test the functioning of the model without the need for graphical interfaces or additional integrations. The CLI is particularly useful for technical users who want to quickly validate the model, experiment with prompts, or integrate the model into automated scripts.

- Solution

Prepare the Python file with the code to interact with the model via the command line, using the script that, in this case, was named `chat_CLI.py`.

- Setup and Code Used

The following shows the code used:

```
chat_CLI.py > ...
1  import torch
2  import os
3  from transformers import LlamaTokenizer, LlamaForCausalLM
4  from datetime import datetime
5
6  local_directory = "./open_llama_3b_v2_local"
7  device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
8  print(f"🖥️ Using device: {device}")
9
10 tokenizer = LlamaTokenizer.from_pretrained(local_directory, local_files_only=True, use_fast=False)
11 model = LlamaForCausalLM.from_pretrained(local_directory, local_files_only=True).to(device)
12
13 # Personalized initial message
14 if tokenizer.pad_token is None:
15     tokenizer.pad_token = tokenizer.eos_token
16
17 print("\n Conversation started with Open LLaMA 3B v2 (type '/exit' to exit)")
18
19 while True:
20     # Add current time in [HH:MM] format to user prompt
21     time_now = datetime.now().strftime("%H:%M")
22     user_input = input(f"\n[{time_now}] User: ")
23     command = user_input.strip().lower()
24
25     # Command to stop the interaction
26     if command == "/exit":
27         print("👋 Ending conversation. Goodbye!")
28         break
29
30     # Command to clear the terminal screen
```

```

chat_CLI.py > ...
30     # Command to clear the terminal screen
31     if command == "/clear":
32         os.system('cls' if os.name == 'nt' else 'clear')
33         continue
34
35     inputs = tokenizer(user_input, return_tensors="pt", padding=True).to(device)
36
37     output_ids = model.generate(
38         **inputs,
39         max_new_tokens=100,
40         pad_token_id=tokenizer.pad_token_id,
41         eos_token_id=tokenizer.eos_token_id,
42         do_sample=False,
43         repetition_penalty=1.2 # Penalize repetitions in the response
44     )
45
46     response = tokenizer.decode(output_ids[0], skip_special_tokens=True)
47
48     # Remove the question from the beginning of the response if echoed
49     if user_input.lower() in response.lower():
50         response = response.split(user_input, 1)[-1].strip()
51
52     # Remove simple repetitions (keep only the first line of the response)
53     response_lines = [line.strip() for line in response.split("\n") if line.strip()]
54     if response_lines:
55         response = response_lines[0]
56
57     # Add current time prefix to the model's response
58     time_now = datetime.now().strftime("%H:%M")
59     print(f"[{time_now}] LLama3Bv2: {response}")

```

The model can be customized in various ways, and in this case, some customizations were implemented to improve the user experience, such as adding the time to the messages, defining personalized names for the user and the model, the command to clear the terminal, the command to stop the interaction, and filtering responses to remove repetitions.

Then, run the script chat\_CLI.py in the terminal.

```

(venv) C:\modelo_open_llama_3b-v2>python chat_CLI.py
Using device: cpu
Loading checkpoint shards: 100% | 3/3 [00:00<00:00, 8.53it/s]

Conversation started with Open LLaMA 3B v2 (type '/exit' to exit)

[21:17] User:

```

- Example of Interaction

The following shows an example of interaction in the terminal, where it is possible to observe the customizations implemented in the code.

```

Conversation started with Open LLaMA 3B v2 (type '/exit' to exit)

[21:22] User: What is the capital of france?
[21:22] LLama3Bv2: Paris, France.

[21:22] User: /clear

```

```

[21:26] User: which is the biggest desert on earth?
[21:26] LLama3Bv2: The largest deserts in terms of area are: The Sahara Desert, located mostly within Africa and covering an estimated 3.5 million square miles (9.076 km2), or about one-third...

[21:26] User: /exit
👋 Ending conversation. Goodbye!

(venv) C:\modelo_open_llama_3b-v2>

```

- Final Remarks

These customizations make the interaction clearer, more organized, and easier to use, contributing to a satisfying and smooth user experience. Some of them were included simply to make the interaction more pleasant and intuitive, even though they are not essential for the functioning of the model.

## 4.2 Exercise 2: Interaction via Telegram Bot

- Description

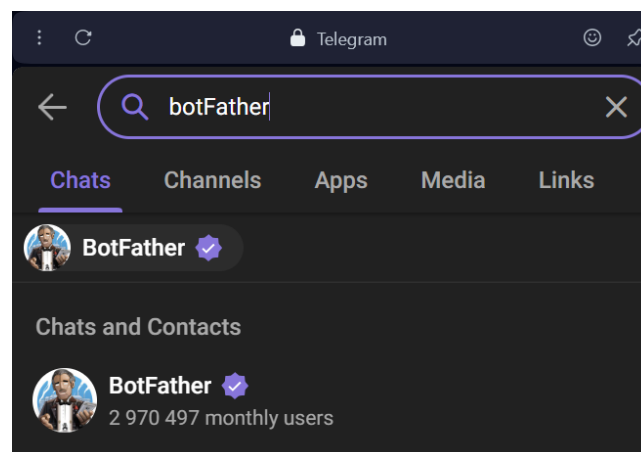
It is possible to connect the local model to a Telegram bot, allowing interaction directly through the app's chat, whether on a mobile phone, tablet, or PC. The user sends messages to the bot on Telegram, and the model, running on the computer, processes these messages and returns the responses through the same channel. For this interaction to work, it is important that both the computer running the model and the user's device interacting with Telegram are connected to the internet. Although the model's processing occurs locally, an internet connection is essential on both sides, as Telegram is a cloud-based application and depends on this connection to send and receive messages between the user and the bot.

- Solution

Create the Python file with the code to integrate and interact with the model through the Telegram bot, using the script that, in this case, was named `chat_telegram.py`.

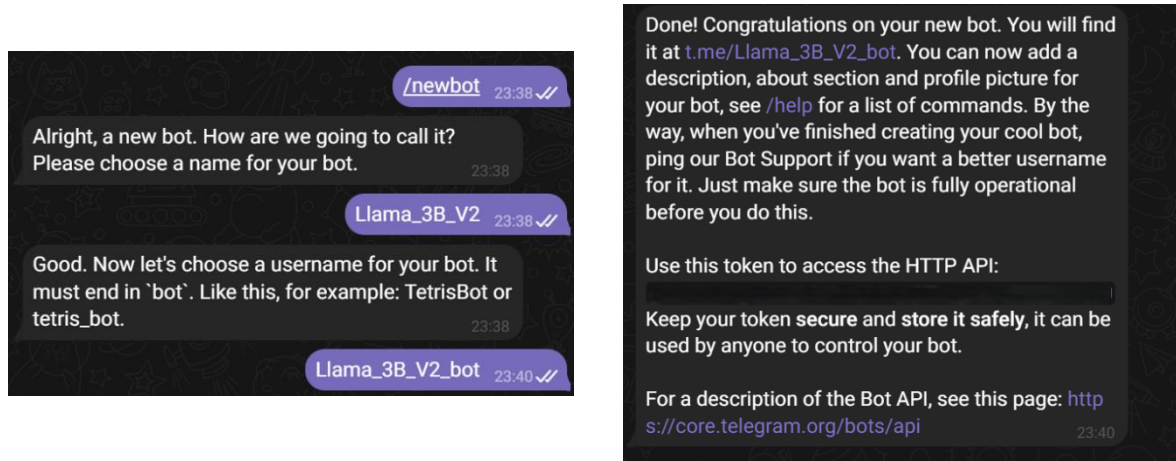
- Setup and Code Used

Open Telegram and use the search bar to look for "BotFather."



**Figure 2 - Telegram Interface**

In the chat with BotFather, the `/newbot` command was used to start the creation of a new bot, assigning it the name `Llama_3B_V2` and the username `Llama_3B_V2_bot` (with the required "bot" ending). After these steps, a token was generated that will be used in the Python script to integrate the model with Telegram.



**Figure 3 - Token Access**

It is necessary to install the required dependencies, such as the `python-telegram-bot` library, using the `pip` package manager.

```
pip install python-telegram-bot
```

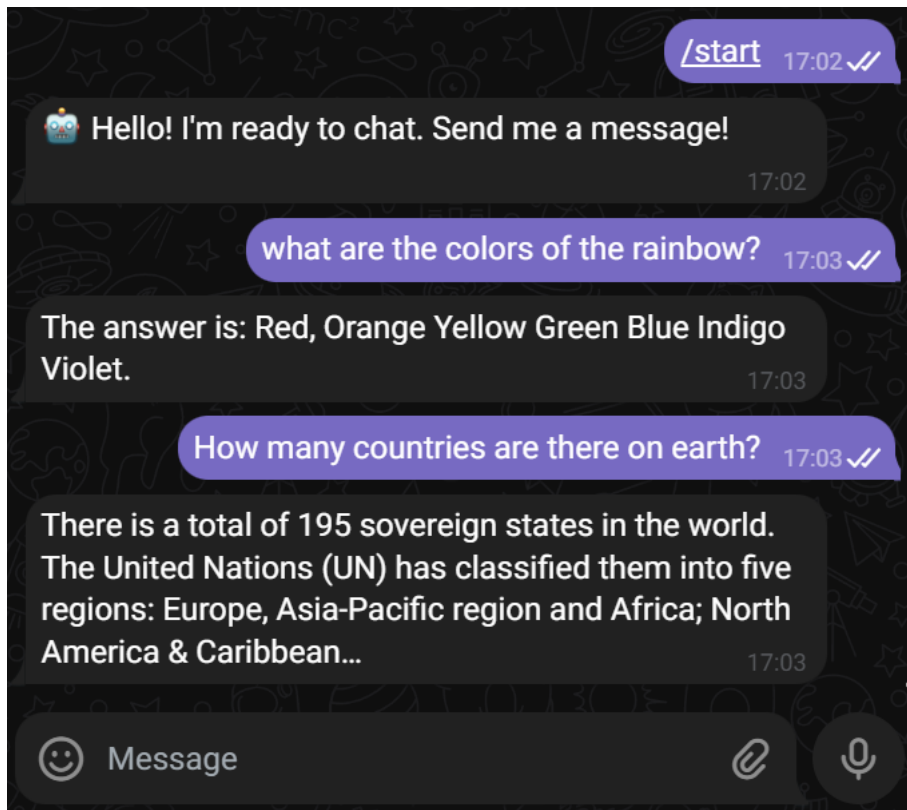
Then, create the Python file with the code that enables the integration and interaction of the model with the application's bot.

```
chat_telegram.py > handle_message
1  import torch
2  from transformers import LlamaTokenizer, LlamaForCausalLM
3  from telegram import Update
4  from telegram.ext import ApplicationBuilder, MessageHandler, CommandHandler, ContextTypes, filters
5  import logging
6
7  # Setup logging
8  logging.basicConfig(format='%(asctime)s - %(name)s - %(levelname)s - %(message)s', level=logging.INFO)
9
10 # Load model and tokenizer
11 local_directory = "./open_llama_3b_v2_local"
12 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
13 tokenizer = LlamaTokenizer.from_pretrained(local_directory, local_files_only=True, use_fast=False)
14
15 # Corrige pad_token_id se não existir
16 if tokenizer.pad_token is None:
17     tokenizer.pad_token = tokenizer.eos_token
18
19 model = LlamaForCausalLM.from_pretrained(local_directory, local_files_only=True).to(device)
20
21 # Define start command
22 async def start(update: Update, context: ContextTypes.DEFAULT_TYPE):
23     await update.message.reply_text("👋 Hello! I'm ready to chat. Send me a message!")
24
25 # Define message handler
```









**Figure 4 - Example of Interaction**

- Final Remarks

One of the main advantages of this integration is its portability and ease of use, allowing interaction with the model from different devices with access to Telegram, as long as they have permission to communicate with the bot. A very important point to pay attention to is security, ensuring that the bot token generated by BotFather is kept confidential and not shared publicly to prevent unauthorized access.

### **4.3 Exercise 3: Interaction via Web Interface**

- Description

Interaction through a web interface allows the model to be used directly in a browser, making the experience accessible from any device on the same network, such as a PC, tablet, or smartphone. This approach typically uses frameworks like Gradio, FastAPI, Flask, or Streamlit to create the backend and expose the application simply and intuitively for the end user.

- Solution

A Python file was developed using Gradio to build the web interface, using the script that, in this case, was named `interface_web.py`.

- Setup and Code Used

It is mandatory to install the required dependencies, such as the gradio library, using the pip package manager.

```
pip install gradio
```

To create the web interface, a Python file was prepared using Gradio. The code defines the input and output fields, the connection between the interface and the local model, and the clear, submit, and flag buttons.

```

1  import torch
2  from transformers import LlamaTokenizer, LlamaForCausalLM
3  import gradio as gr
4
5  local_directory = "./open_llama_3b_v2_local"
6  device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
7  tokenizer = LlamaTokenizer.from_pretrained(local_directory, local_files_only=True, use_fast=False)
8
9  if tokenizer.pad_token is None:
10     tokenizer.pad_token = tokenizer.eos_token
11
12  model = LlamaForCausalLM.from_pretrained(local_directory, local_files_only=True).to(device)
13
14  def responder_mensagem(user_input):
15     inputs = tokenizer(user_input, return_tensors="pt", padding=True).to(device)
16     output_ids = model.generate(
17         **inputs,
18         max_new_tokens=100,
19         pad_token_id=tokenizer.pad_token_id,
20         eos_token_id=tokenizer.eos_token_id,
21         do_sample=False,
22         repetition_penalty=1.2
23     )
24     response = tokenizer.decode(output_ids[0], skip_special_tokens=True)
25     if user_input.lower() in response.lower():
26         response = response.split(user_input, 1)[-1].strip()
27     response_lines = [line.strip() for line in response.split("\n") if line.strip()]
28     if response_lines:
29         response = response_lines[0]
30     return response
31
32  gr.Interface(fn=responder_mensagem, inputs="text", outputs="text", title="Open LLaMA 3B-V2 Chatbot").launch()

```

After the implementation of the code, the model becomes accessible through a browser, allowing the user to interact directly with the developed web interface. This solution makes it easy to send questions and view the responses generated by the model in a practical and intuitive way.

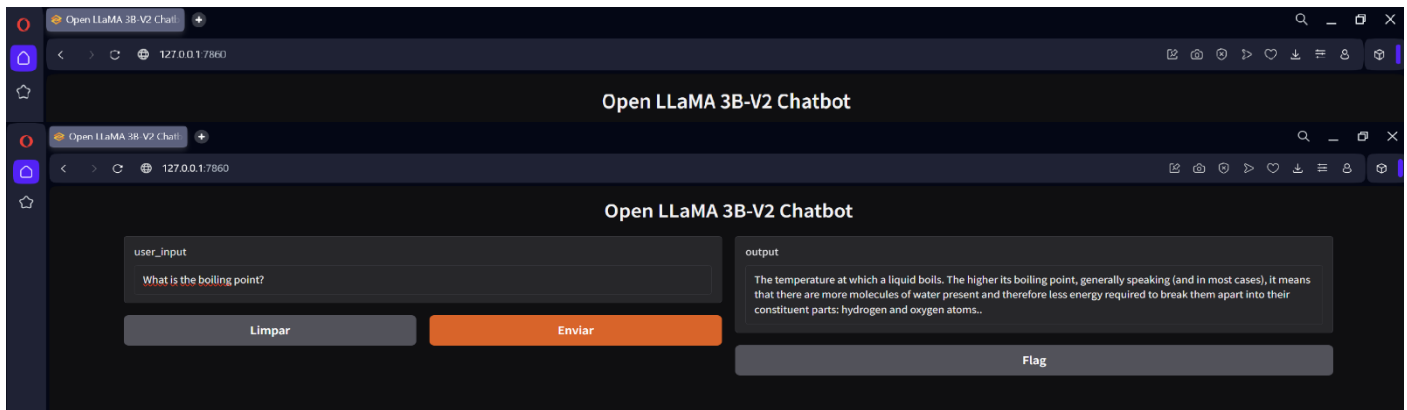
After that, run the command in the terminal to execute the script, in this case, `interface_web.py`.

```
(venv) C:\modelo_open_llama_3b-v2>python web_interface.py  
Loading checkpoint shards: 100%|██████████████████████████████████████████████████████████████████████████| 3/3 [00:00<00:00, 5.82it/s]  
* Running on local URL: http://127.0.0.1:7860  
* To create a public link, set `share=True` in `launch()`.
```

Next, copy the link provided in the terminal and paste it into the browser's address bar to access the web interface.

- Example of Interaction

Here is a screenshot of the browser showing the web interface, with the question field filled in and the response generated by the model.



**Figure 5 - Example of Web interaction**

- Final Remarks

The web interface facilitates access to the model, eliminating the need for interaction via terminal or other technical tools. In addition, it offers the possibility of being accessed from different devices on the local network. It also serves as a foundation that can be extended for online deployment, with adaptations to the layout and additional features such as chat history, conversation export, or integration with external APIs.

## 4.4 Expansion and Integration Potential

In addition to the interaction methods addressed in this project — CLI, Telegram, and web interface — there are still other possibilities that significantly expand the potential uses of the model. For example, automation with tools such as n8n, Node-RED, or Zapier allows the model to be integrated into automated workflows, generating responses, classifying information, or performing actions without manual intervention. Extensions or plugins for tools like VS Code, Obsidian, or Jupyter offer an integrated experience directly within development or writing environments, making the use of the model even more seamless. Finally, voice interaction, through systems like Google Assistant, Alexa, or custom solutions with speech recognition, paves the way for innovative and accessible applications.

## 5. LLaMA 3B v2 Model and Tiny LLaMA 1.1B Model

Meta developed the LLaMA 3B v2 model and is part of the second generation of the LLaMA family. It features 3 billion parameters and has been trained on large volumes of text, enabling it to perform tasks such as answering questions, generating text, summarizing, and translating, all with a high level of coherence. This version is robust enough for applications in chatbots, virtual assistants, and NLP tasks with limited resources. It is often used in projects where it is necessary to maintain efficiency and quality without relying on highly powerful hardware.

The Tiny LLaMA 1.1B model is one of the versions of the LLaMA architecture. It features 1.1 billion parameters and is both compact and highly efficient. This version was developed with a focus on environments with limited resources, such as mobile devices or edge applications, while maintaining strong capabilities in natural language understanding and generating. Although its performance is lower than that of larger models, it offers faster response times and reduced memory and processing usage, making it ideal for rapid testing, prototyping, or systems with infrastructure constraints.

### 5.1 Model Comparison

<i>Feature</i>	<i>LLaMA 3B v2</i>	<i>Tiny LLaMA 1.1B</i>
<b>Developer</b>	Meta	Community (based on LLaMA architecture)
<b>Architecture</b>	LLaMA v2	LLaMA-based (Tiny variant)
<b>Parameter Count</b>	3 billion	1.1 billion
<b>Performance</b>	High coherence, good general purpose NLP	Lower than larger models, but effective for lightweight tasks
<b>Resource Requiremen</b>	Moderate hardware	Very low – optimized for edge/mobile
<b>Ideal Use Cases</b>	Chatbots, assistants, general NLP	Prototyping, quick tests, mobile/embedded systems
<b>Advantages</b>	Balanced between quality and efficiency	Fast response, low memory/CPU usage
<b>Limitations</b>	Heavier than smaller models	Less accurate and powerful than larger models

**Table 2- Comparison between LLaMA 3B v2 and Tiny LLaMA 1.1B models**

The main difference between the models, as shown in Table 1, lies in performance and efficiency. The LLaMA 3B v2 model is more suitable for tasks that require more elaborate responses due to its higher capacity for text comprehension and generation, which demands more computational resources. On the other hand, the Tiny LLaMA 1.1B model is lighter and faster, making it ideal for cases where memory and processing consumption are limited.

## 6. References

Meta AI. (2024). *LLaMA 3.2 3B – Modelo de linguagem*. Disponível em: <https://huggingface.co/meta-llama/Llama-3.2-3B>

Meta AI. (2024). *LLaMA 3.2 Connect: Visão para dispositivos móveis e edge*. Meta AI Blog. Disponível em: <https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edge-mobile-devices/>

Zhang, L. et al. (2024). *LLaMA 3 Technical Report*. arXiv preprint. Disponível em: <https://arxiv.org/abs/2401.02385>

Asimov Academy. (2023). *Tutorial: Como acessar modelos restritos no Hugging Face*. Disponível em: <https://hub.asimov.academy/tutorial/como-acessar-modelos-restritos-no-hugging-face/>

Hugging Face. (2024). *Bem-vindo ao Hugging Face Hub*. Disponível em: <https://huggingface.co/welcome>

Hugging Face. (2024). *Verificação de segurança do utilizador*. Disponível em: <https://huggingface.co/security-checkup>

Hugging Face. (2024). *Guia rápido do Hugging Face Hub*. Disponível em: [https://huggingface.co/docs/huggingface\\_hub/en/quick-start](https://huggingface.co/docs/huggingface_hub/en/quick-start)

Meta AI. (2023). *LLaMA 2 7B – Modelo de linguagem*. Disponível em: <https://huggingface.co/meta-llama/Llama-2-7b-hf>

Hugging Face. (2024). *Guia sobre modelos com acesso restrito (Gated Models)*. Disponível em: <https://huggingface.co/docs/hub/en/models-gated>

Hugging Face. (2024). *Guia para uso de modelos privados com Transformers.js*. Disponível em: <https://huggingface.co/docs/transformers.js/en/guides/private>