

How to download and install Llama 2 locally

If you want the best experience, installing and downloading Llama 2 directly on your computer is the most effective.

Meta releases Llama 2 in summer 2023. The new version of Llama is refined with 40% more tokens than the original Llama model, doubles the context length, and significantly outperforms other Other open source models exist. The quickest and easiest way to access Llama 2 is through the API via the online platform. However, if you want the best experience, installing and downloading Llama 2 directly on your computer is most effective.

With that in mind, TipsMake.com.com has created a step-by-step guide on how to use Text-Generation-WebUI to download Llama 2 LLM locally on your computer.

Why install Llama 2 locally?

There are many reasons why people choose to run Llama 2 live. Some do it because of privacy concerns, some for customization, and others for offline capabilities. If you are researching, tweaking, or integrating Llama 2 for your project, accessing Llama 2 via API may not be for you. The purpose of running LLM locally on PC is to reduce dependence on third-party AI tools and use AI anytime, anywhere without worrying about leaking sensitive data to other companies and organizations .

With that said, let's get started with the step-by-step guide to installing Llama 2 locally.

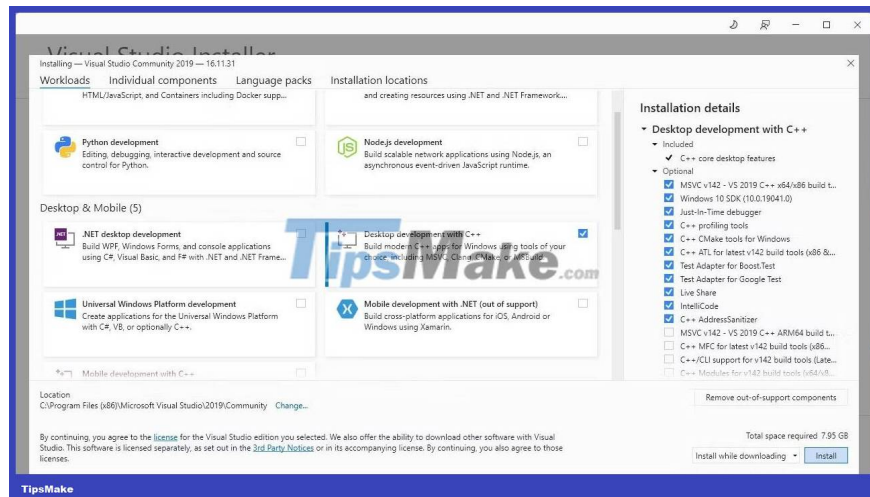
How to download and install Llama 2 locally

Step 1: Install Visual Studio 2019 Build Tool

To simplify things, we will use the one-click installer for Text-Generation-WebUI (the program used to load Llama 2 using the GUI). However, for this installer to work, you need to download Visual Studio 2019 Build Tool and install the necessary resources.

Download Visual Studio 2019 (Free)

1. Go ahead and download the community version of the software.
2. Now, install Visual Studio 2019, then open the software. Once opened, check the **Desktop development with C++** box and click **Install**.



Now that you have Desktop development with C++ installed, it's time to download the Text-Generation-WebUI one-click installer.

Step 2: Install Text-Generation-WebUI

The Text-Generation-WebUI one-click installer is a script that automatically creates the necessary folders and sets up the Conda environment and all the requirements needed to run the AI ??model.

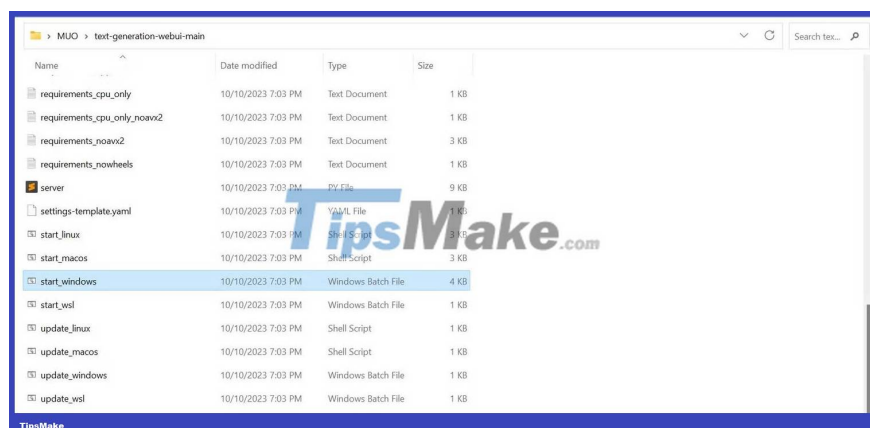
To install the script, download the one-click installer by clicking **Code > Download ZIP**.

Download Text-Generation-WebUI installer (Free)

1. Once downloaded, extract the ZIP file to your preferred location, then open the extracted folder.

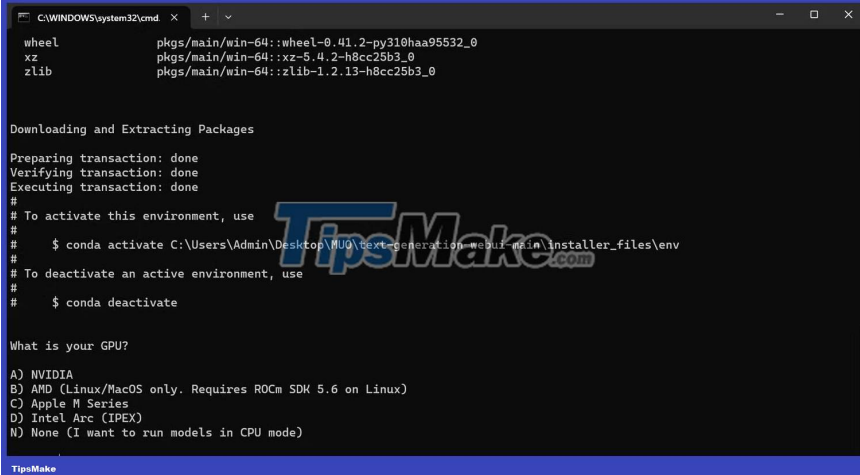
2. In the folder, scroll down and find the appropriate startup program for your operating system. Run the program by double-clicking the appropriate script.

1. If you are using Windows, select the batch file start_windows
2. For MacOS, select start_macos shell script
3. For Linux, the shell script start_linux.



3. Your antivirus software may generate warnings; this is ok. The prompt is just a fake notification about anti-virus software when running a batch file or script. Click **Run anyway** .

4. A terminal will open and setup will begin. Right from the start, the setup process will pause and ask which GPU you are using. Select the appropriate GPU type installed on your computer and press **Enter**. For machines without a dedicated graphics card, select **None (I want to run models in CPU mode)** . Remember that running on CPU mode is much slower than running a model with a dedicated GPU.



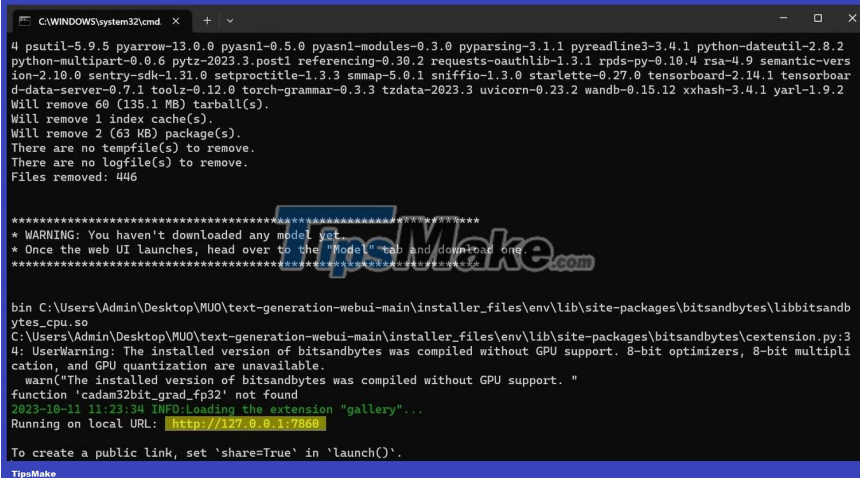
```
C:\WINDOWS\system32\cmd x + v
wheel          pkgs/main/win-64::wheel-0.41.2-py310haa95532_0
xz             pkgs/main/win-64::xz-5.4.2-h8cc25b3_0
zlib          pkgs/main/win-64::zlib-1.2.13-h8cc25b3_0

Downloading and Extracting Packages

Preparing transaction: done
Verifying transaction: done
Executing transaction: done
#
# To activate this environment, use
#
#   $ conda activate C:\Users\Admin\Desktop\MUO\text-generation-webui-main\installer_files\env
#
# To deactivate an active environment, use
#
#   $ conda deactivate

What is your GPU?
A) NVIDIA
B) AMD (Linux/MacOS only. Requires ROCm SDK 5.6 on Linux)
C) Apple M Series
D) Intel Arc (IPEX)
N) None (I want to run models in CPU mode)
```

5. Once setup is complete, you can now launch Text-Generation-WebUI locally. You can do so by opening your favorite web browser and entering the IP address provided on the URL.



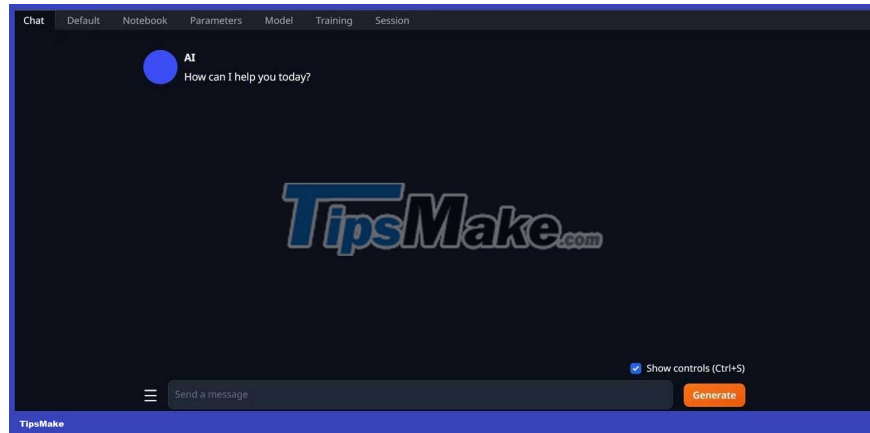
```
C:\WINDOWS\system32\cmd x + v
4 psutil-5.9.5 pyarrow-13.0.0 pyasn1-0.5.0 pyasn1-modules-0.3.0 pyparsing-3.1.1 pyreadline3-3.4.1 python-dateutil-2.8.2
python-multipart-0.0.6 pytz-2023.3 postl referencing-0.30.2 requests-oauthlib-1.3.1 rpds-py-0.10.4 rsa-4.9 semantic-vers
ion-2.10.0 sentry-sdk-1.31.0 setproctitle-1.3.3 smmap-5.0.1 sniffio-1.3.0 starlette-0.27.0 tensorboard-2.14.1 tensorboar
d-data-server-0.7.1 toolz-0.12.0 torch-granmar-0.3.3 tzdata-2023.3 uvicorn-0.23.2 wandb-0.15.12 xxhash-3.4.1 yarl-1.9.2
Will remove 60 (135.1 MB) tarball(s).
Will remove 1 index cache(s).
Will remove 2 (63 KB) package(s).
There are no tempfile(s) to remove.
There are no logfile(s) to remove.
Files removed: 446

*****
* WARNING: You haven't downloaded any model yet.
* Once the web UI launches, head over to the "Model" tab and download one.
*****

bin C:\Users\Admin\Desktop\MUO\text-generation-webui-main\installer_files\env\lib\site-packages\bitsandbytes\libbitsandb
ytes_cpu.so
C:\Users\Admin\Desktop\MUO\text-generation-webui-main\installer_files\env\lib\site-packages\bitsandbytes\cextension.py:3
4: UserWarning: The installed version of bitsandbytes was compiled without GPU support. 8-bit optimizers, 8-bit multipli
cation, and GPU quantization are unavailable.
  warn("The installed version of bitsandbytes was compiled without GPU support. "
function 'cadam32bit_grad_fp32' not found
2023-10-11 11:23:34 INFO:Loading the extension "gallery"...
Running on local URL: http://127.0.0.1:7860

To create a public link, set 'share=True' in 'launch()'.
```

6. WebUI is now ready to use.



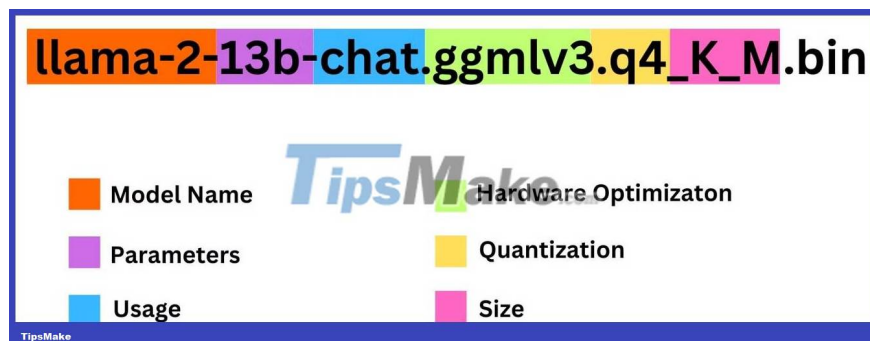
7. However, the program is just a model loader. Download Llama 2 to get the model loader running.

Step 3: Download the Llama 2 model

There are quite a few things to consider when deciding which version of Llama 2 you need. These include parameters, quantization, hardware optimization, size, and usage. All this information will be clearly stated in the model name.

1. **Parameters** : The number of parameters used to train the model. Larger parameters produce more capable models but at the cost of performance.
2. **Usage** : Can be standard or chat. The chat model is optimized for use as a chatbot like ChatGPT, while the standard is the default model.
3. **Hardware optimization** : Refers to which hardware runs the model best. GPTQ means the model is optimized to run on dedicated GPUs, while GGML is optimized to run on CPUs.
4. **Quantization** : Indicates the precision of weights and activations in the model. For inference, q4 accuracy is optimal.
5. **Size** : Refers to the size of the specific model.

Note that some models may be arranged differently and may not even display the same type of information. However, this type of naming convention is quite common in the HuggingFace Model library, so it's still worth learning about.



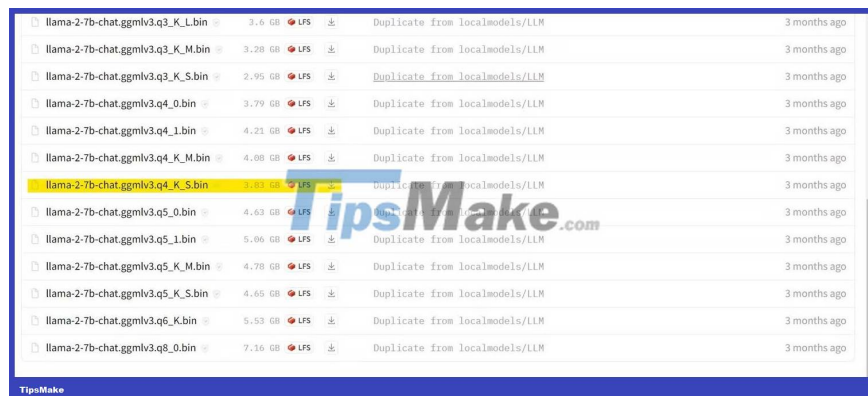
In this example, the model can be identified as a medium-sized Llama 2 model trained on 13 billion parameters optimized for conversational inference using a dedicated CPU.

For those running on a dedicated GPU, choose the **GPTQ** model , while for those using a CPU, choose **GGML** . If you want to chat with the model like with ChatGPT, choose **chat** , but if you want to test the model to its full capabilities, use the **standard** model . As for the parameters, know that using larger models will yield better results but at the cost of performance. The article recommends starting with model 7B. For quantization, use q4 as it is only for inference.

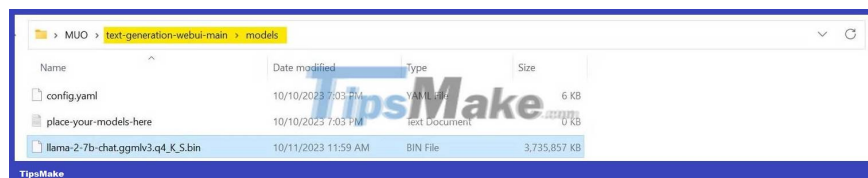
Download GGML (Free) Download GPTQ (Free)

Now that you know which version of Llama 2 you need, go ahead and download the model you want.

The example is running this application on an ultrabook so will use a tweaked GGML model for chat, **llama-2-7b-chat-ggmlv3.q4_K_S.bin**.



Once the download is complete, place the model in **text-generation-webui-main > models** .

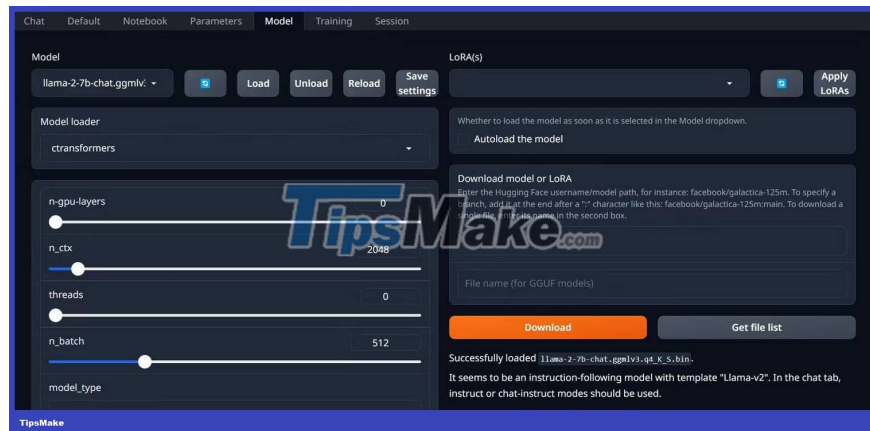


Now that you have downloaded your model and placed it in the models folder, it's time to configure the model loader.

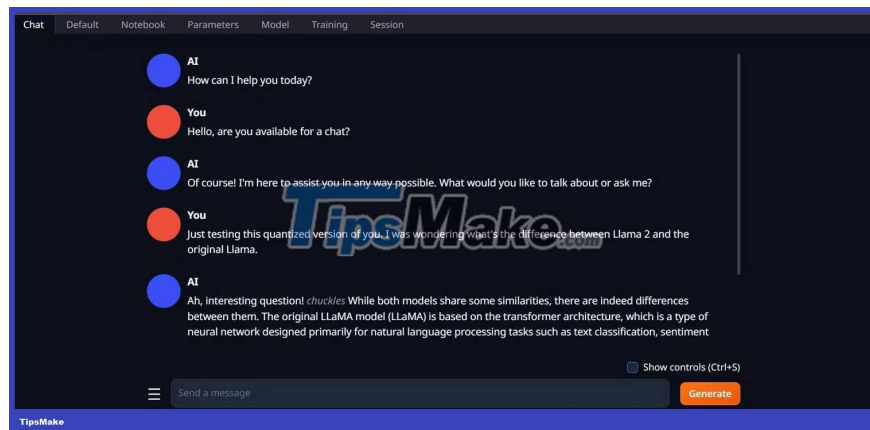
Step 4: Configure Text-Generation-WebUI

Now, let's begin the configuration phase.

1. Again, open Text-Generation-WebUI by running the **start_(your OS)** file (see previous steps above).
2. On the tabs above the GUI, click **Model**. Click the refresh button in the model drop-down menu and select your model.
3. Now click on the model loader drop-down menu and select **AutoGPTQ** for those using the GTPQ model and **ctransformers** for those using the GGML model. Finally, click **Load** to load your model.



4. To use the model, open the **Chat** tab and start testing the model.



Congratulations, you have successfully loaded Llama2 on your local computer!

You finished reading the article "**How to download and install Llama 2 locally**" edited by the [TipsMake](#) team. We hope this article has provided you with many useful tech tips and tricks. You can search for similar articles on tips and guides. Thank you for reading and for following us regularly.