

# ComfyUI-llama-cpp\_vllm 安装步骤

## 目录

- 1、NB管理器安装方法
- 2、TE管理器安装方法
- 3、CMD命令安装办法

-----感谢夜风大佬制作

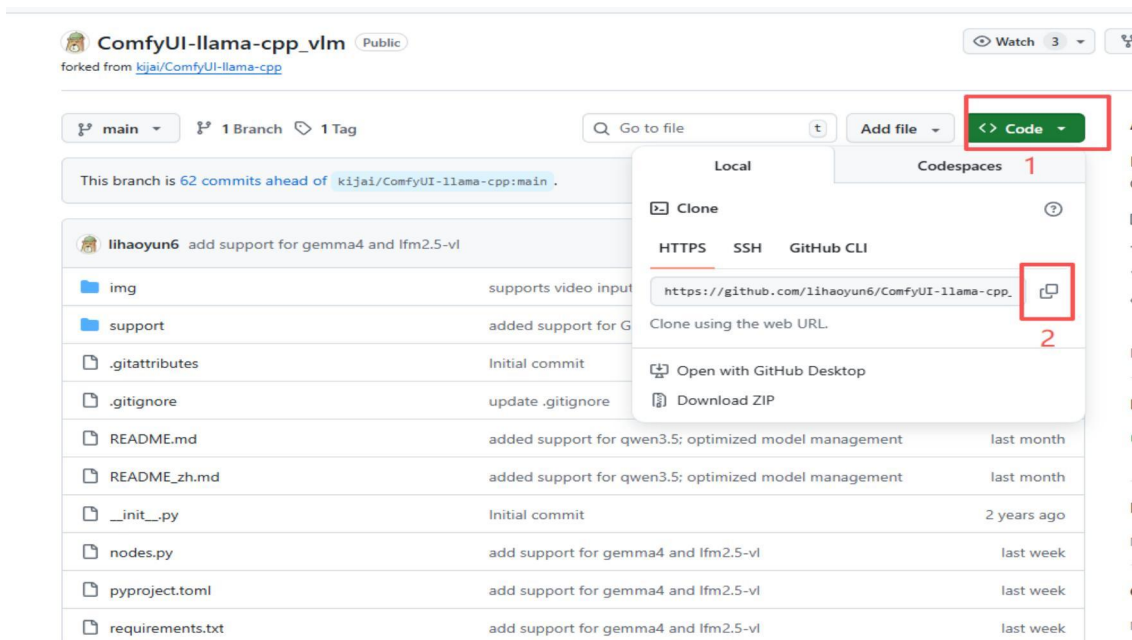
## 1 NB 管理器安装

安装前打开MF

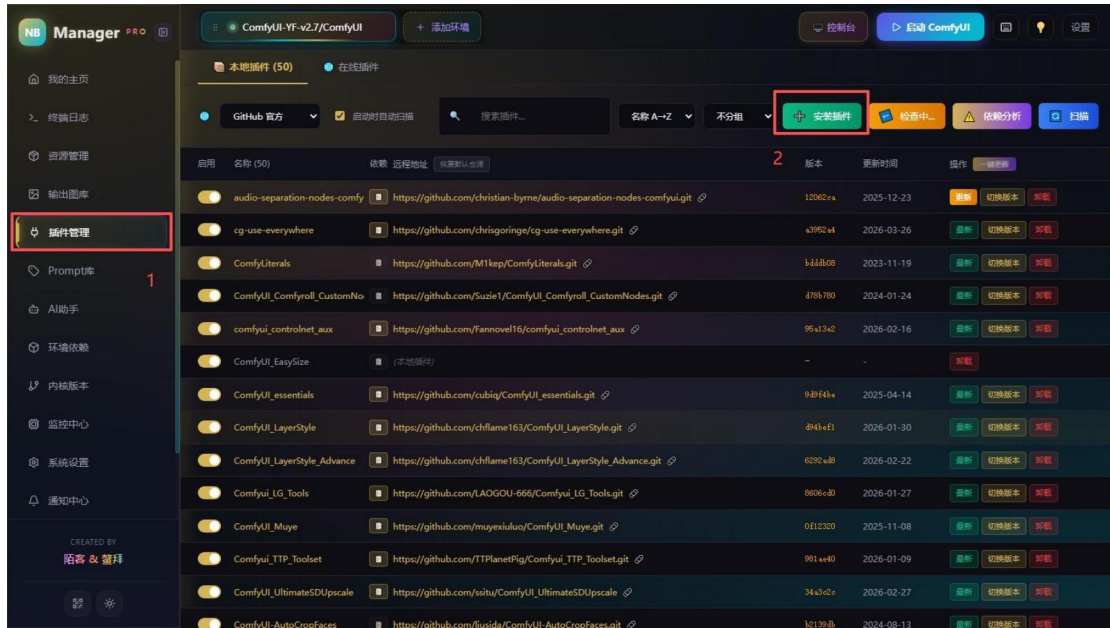
### 1.1 安装插件

1. 打开网址，复制 git 地址：

Github 网址：[https://github.com/lihaoyun6/ComfyUI-llama-cpp\\_vlm](https://github.com/lihaoyun6/ComfyUI-llama-cpp_vlm)

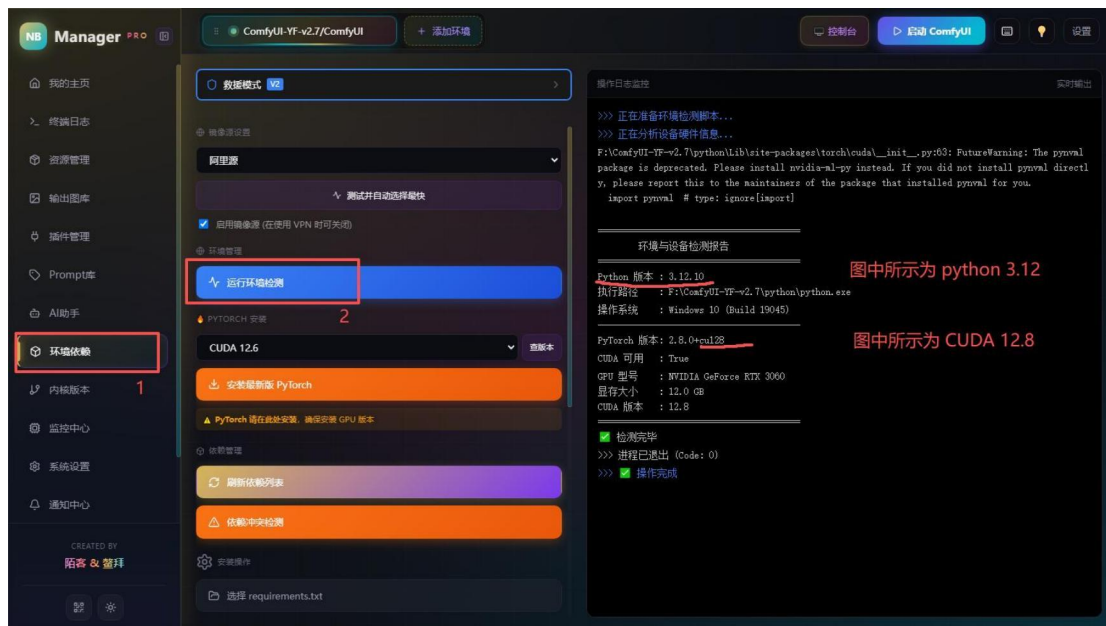


## 2. 打开 NB 启动器



## 1.2 安装轮子

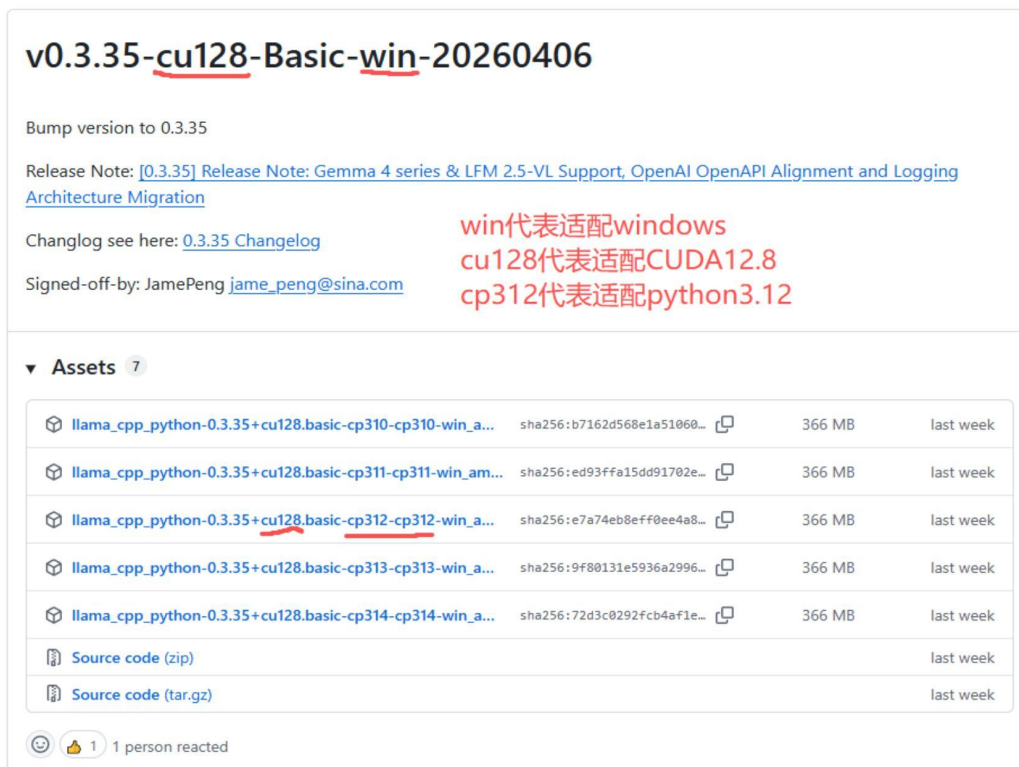
### 1. 查看 python 版本、CUDA 版本



### 2. 下载对应版本的轮子.whl 文件，此处以上图版本为例

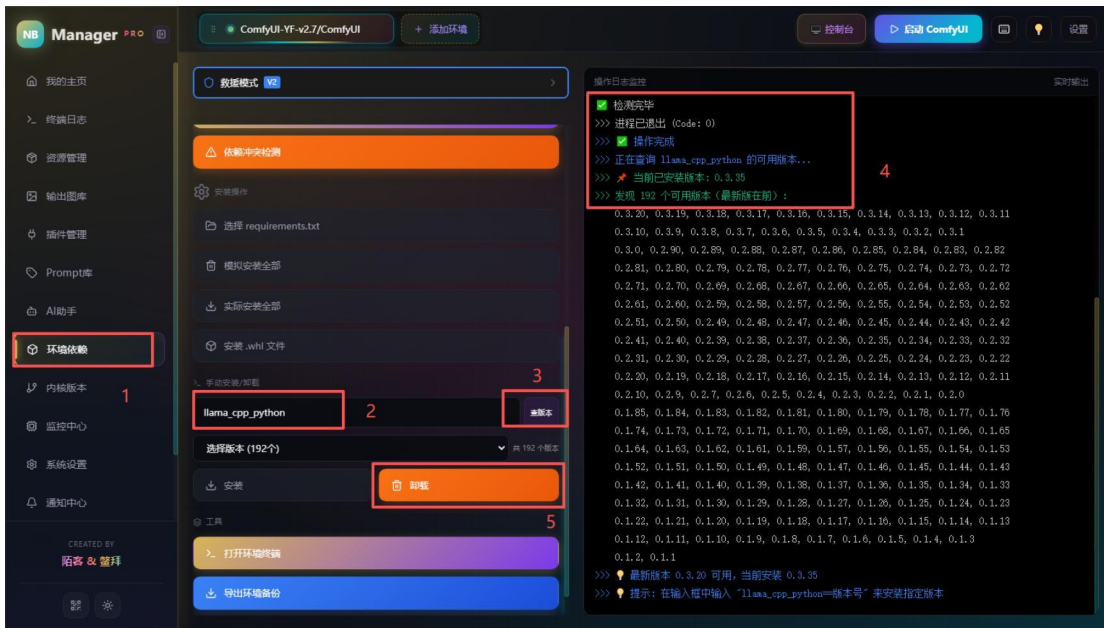
打开网址：<https://github.com/JamePeng/llama-cpp-python/releases>

### 3. 找到适配自己 python 版本+CUDA 版本的轮子，点击蓝字下载

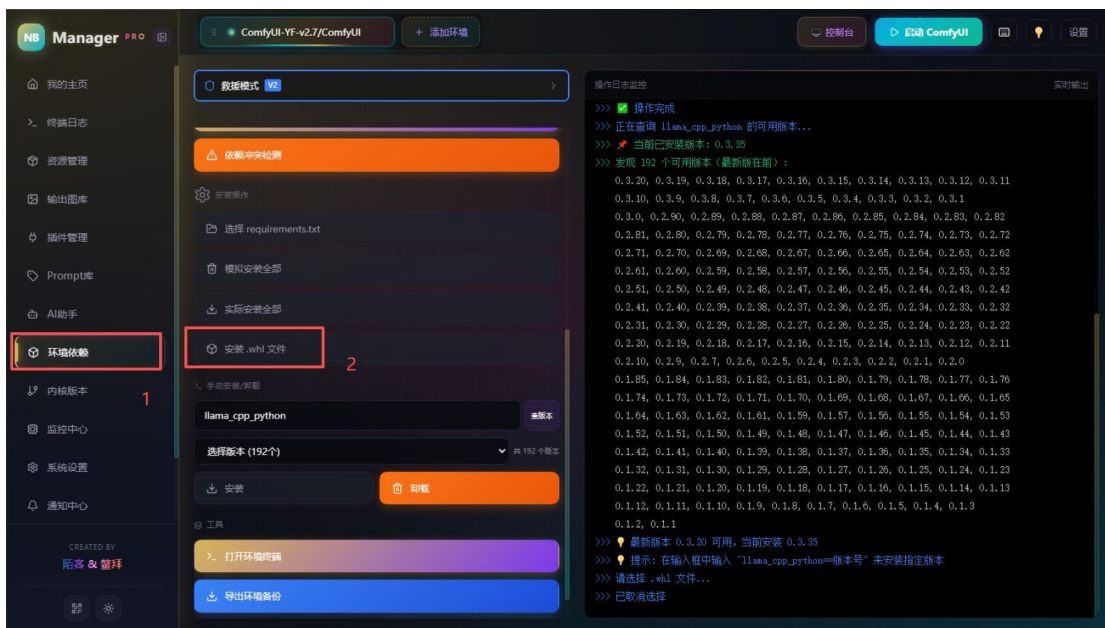


4. 卸载当前安装的轮子（由于不知道你之前是否安装对了轮子，所以先卸载）：

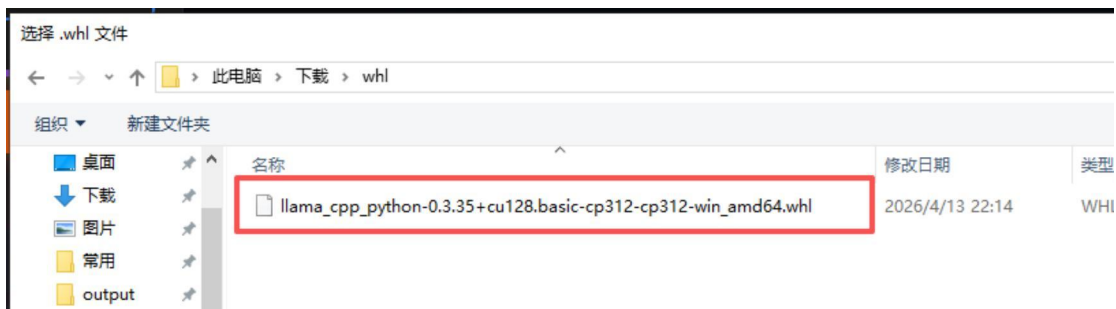
输入查找：llama\_cpp\_python



5. 安装正确版本的轮子：



6. 选择刚刚下载的轮子文件



等待安装结束即可。

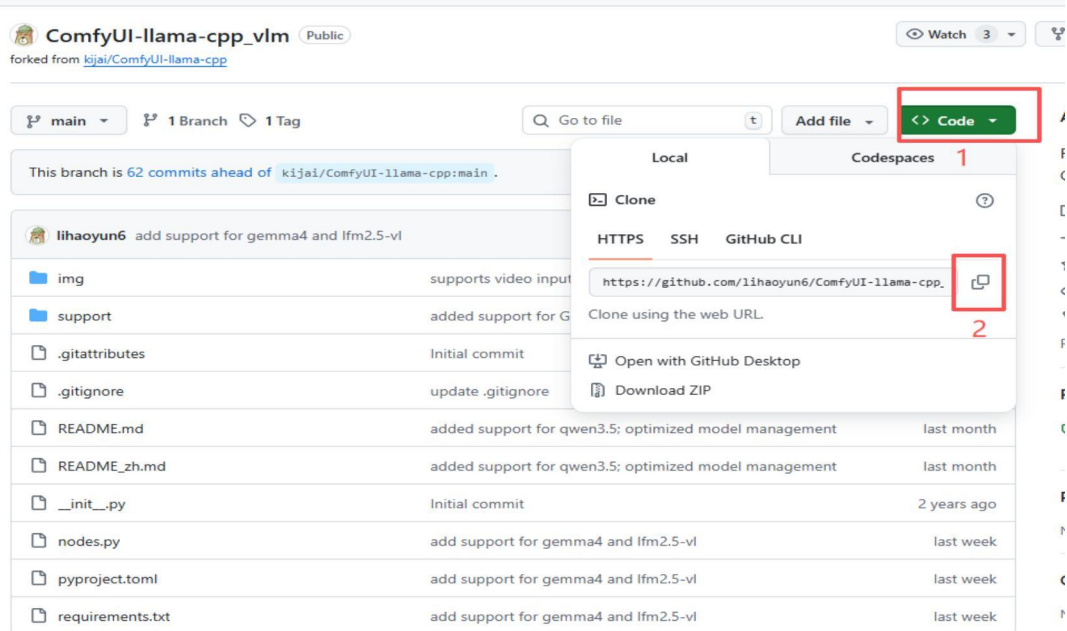
# 2 TE 启动器安装

安装前打开MF

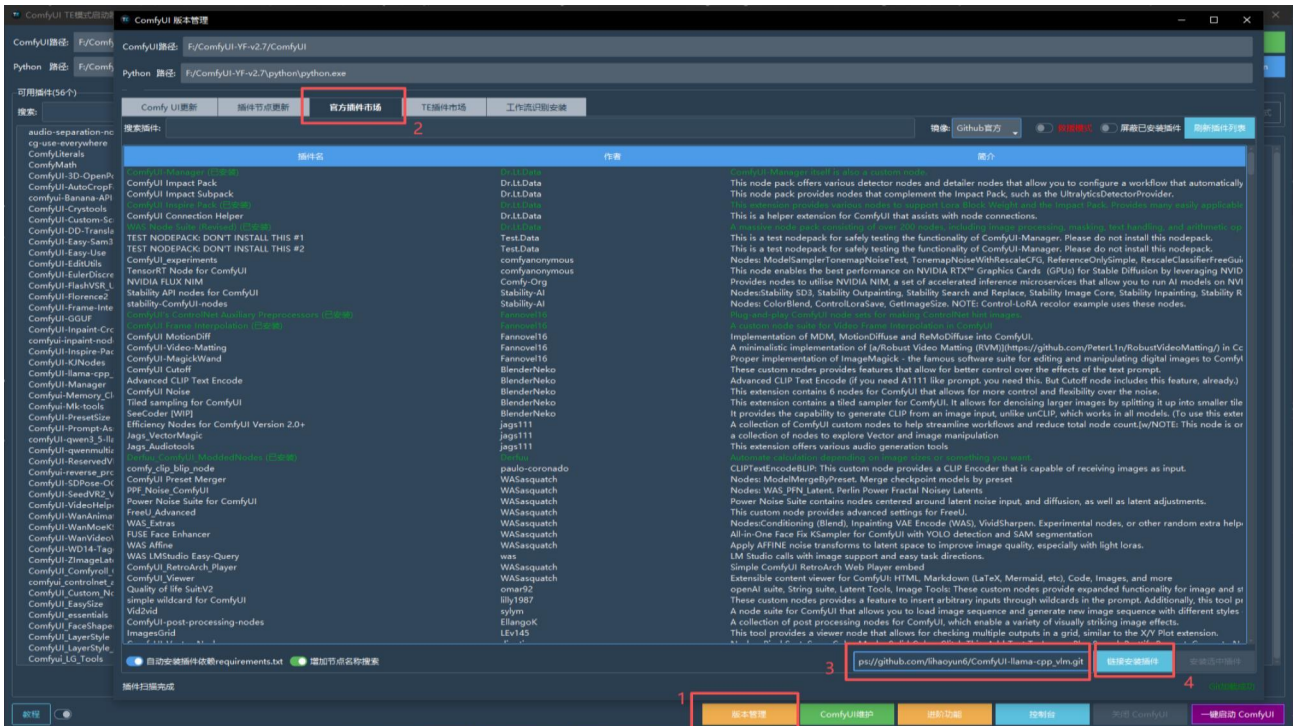
## 2.1 安装插件

1. 打开网址复制 git 地址:

Github 网址: [https://github.com/lihaoyun6/ComfyUI-llama-cpp\\_vlm](https://github.com/lihaoyun6/ComfyUI-llama-cpp_vlm)

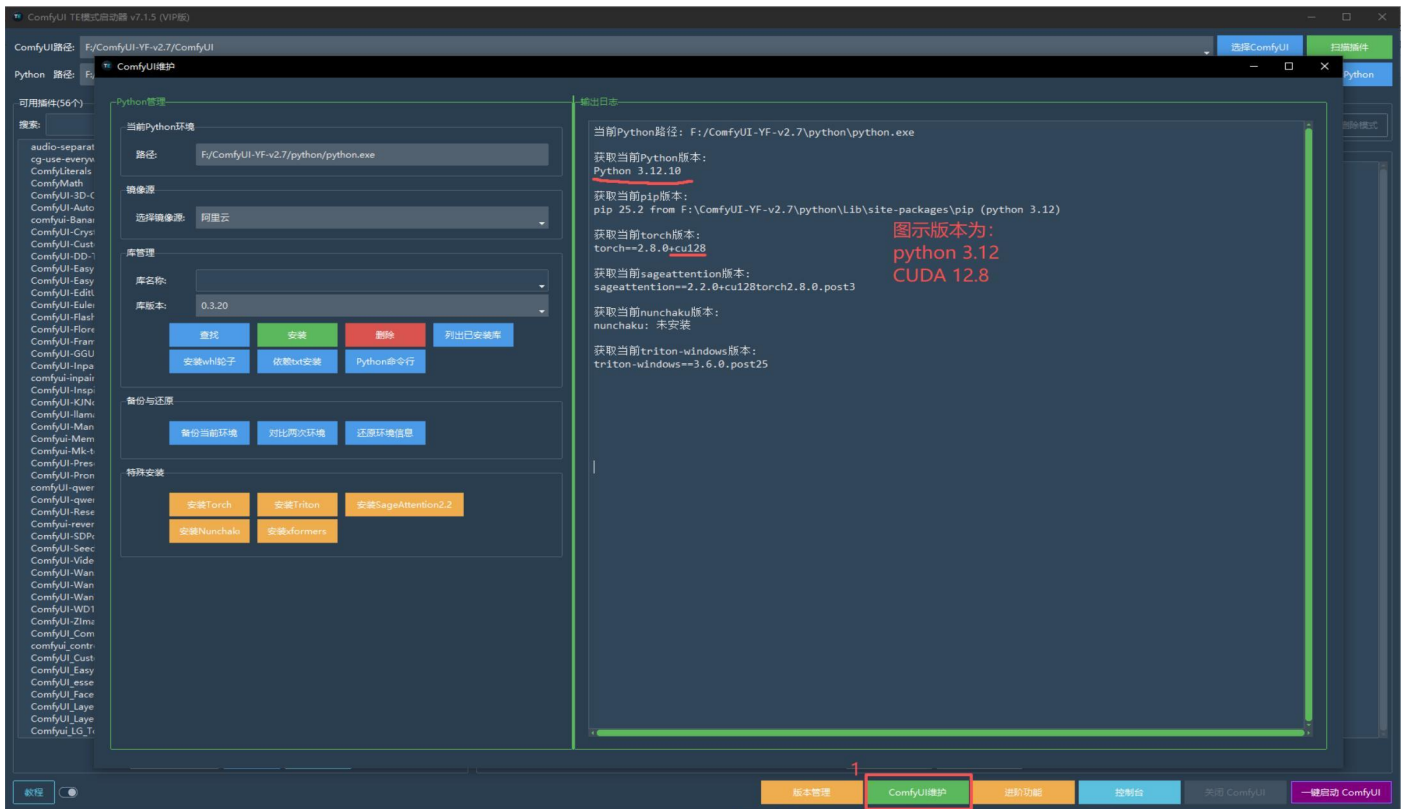


2. 打开 TE 启动器



## 2.2 安装轮子

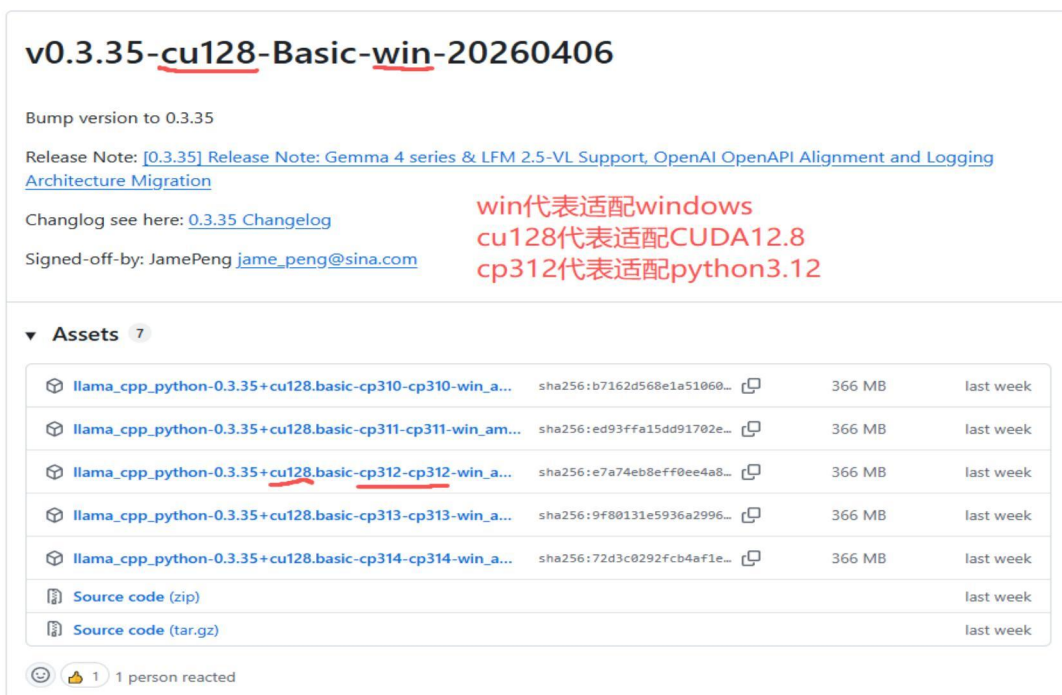
### 1. 查看 python 版本、CUDA 版本



### 2. 下载对应版本的轮子.whl 文件，此处以上图版本为例

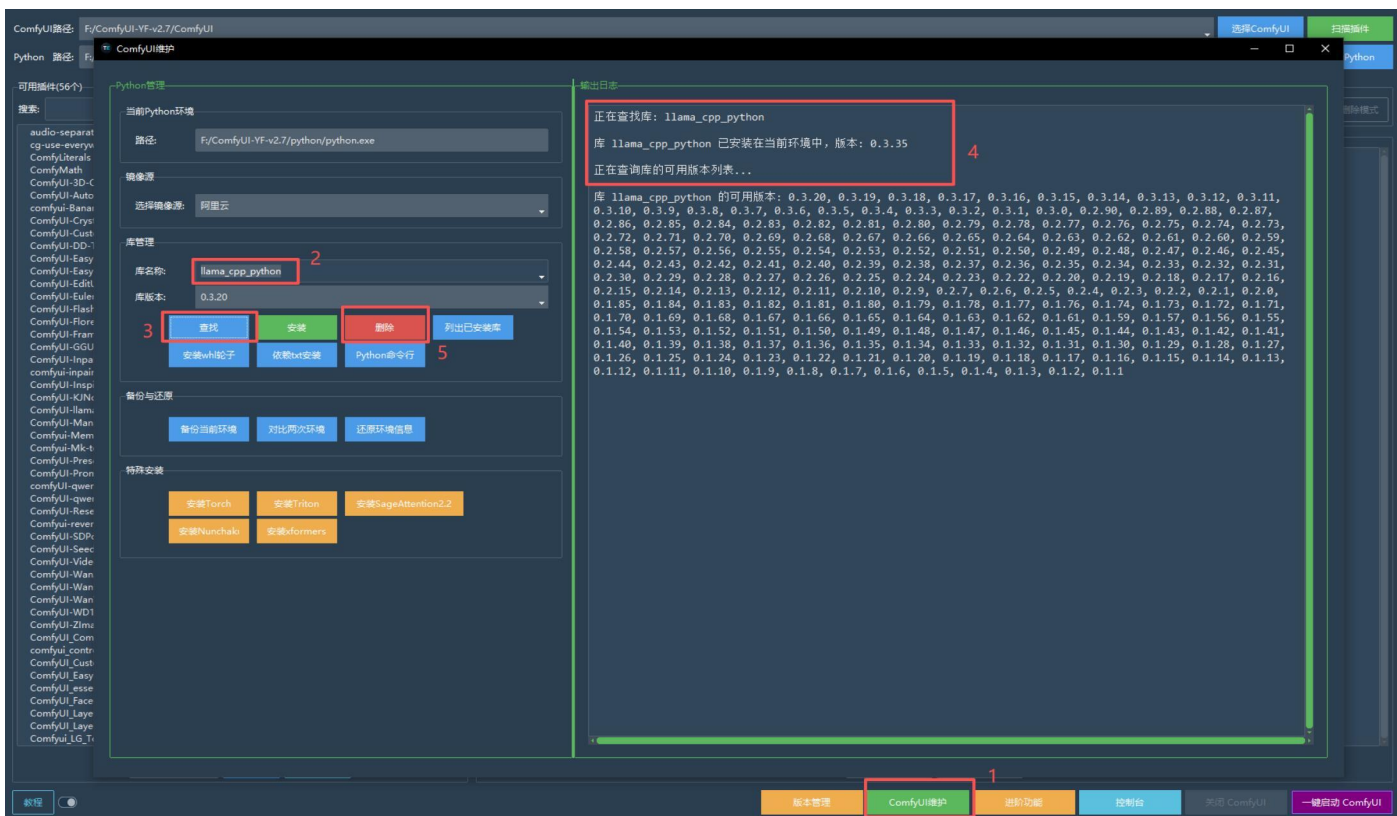
打开网址：<https://github.com/JamePeng/llama-cpp-python/releases>

### 3. 找到适配自己 python 版本+CUDA 版本的轮子，点击蓝字下载

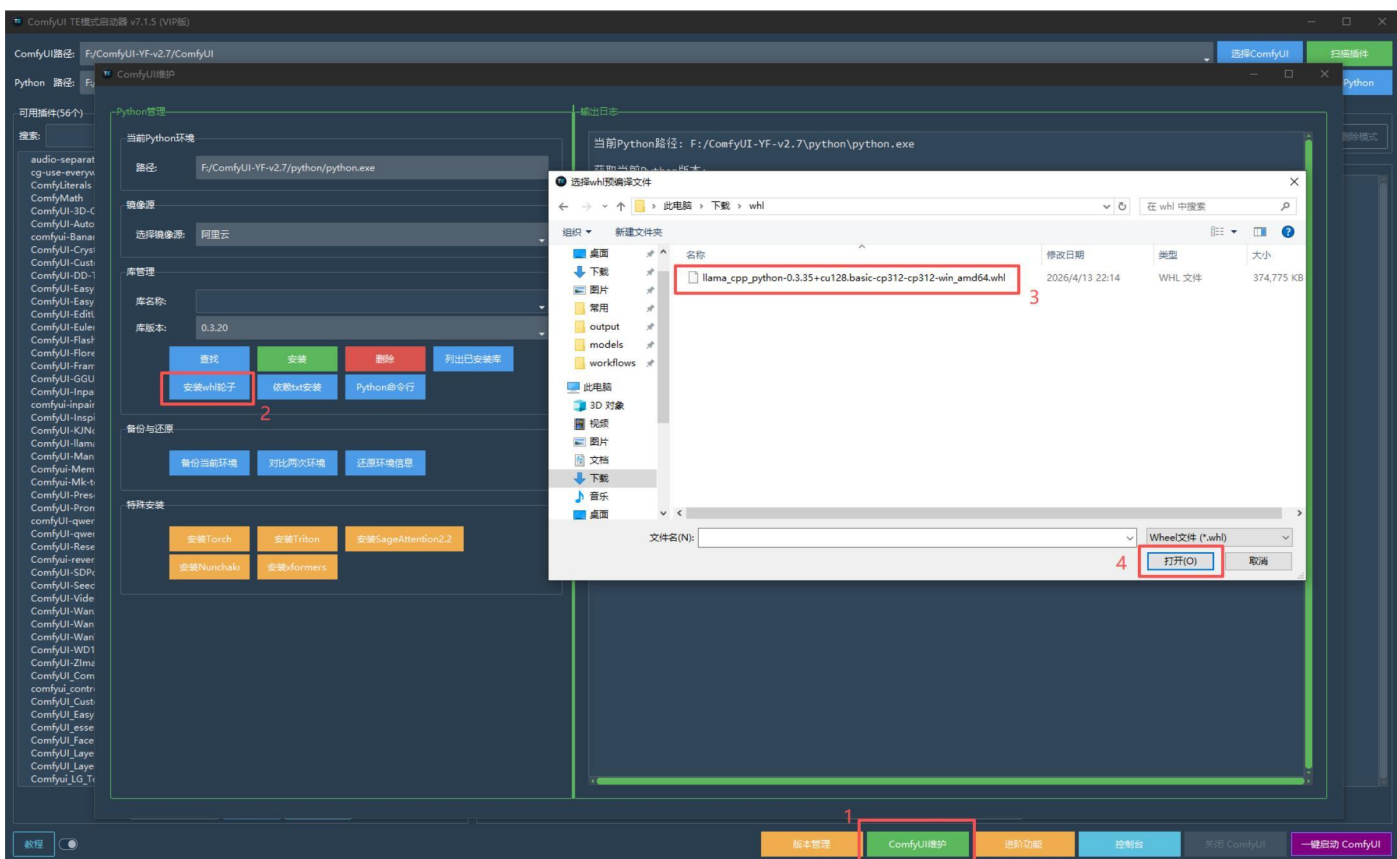


#### 4. 卸载当前安装的轮子（由于不知道你之前是否安装对了轮子，所以先卸载）：

输入查找：llama\_cpp\_python



#### 5. 安装正确版本的轮子：



等待安装结束即可。

# 3 CMD 命令安装

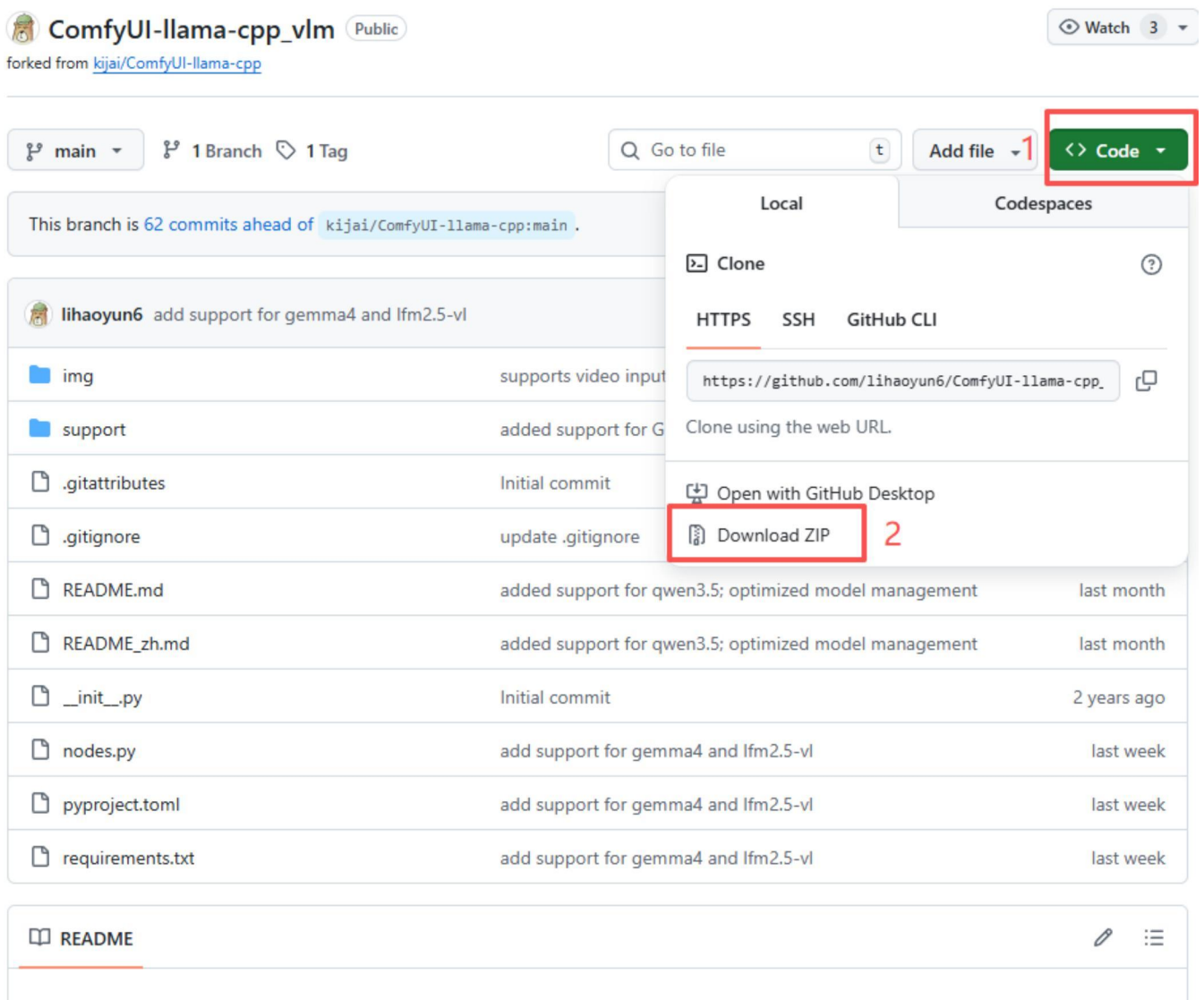
安装前打开魔法

## 3.1 安装插件

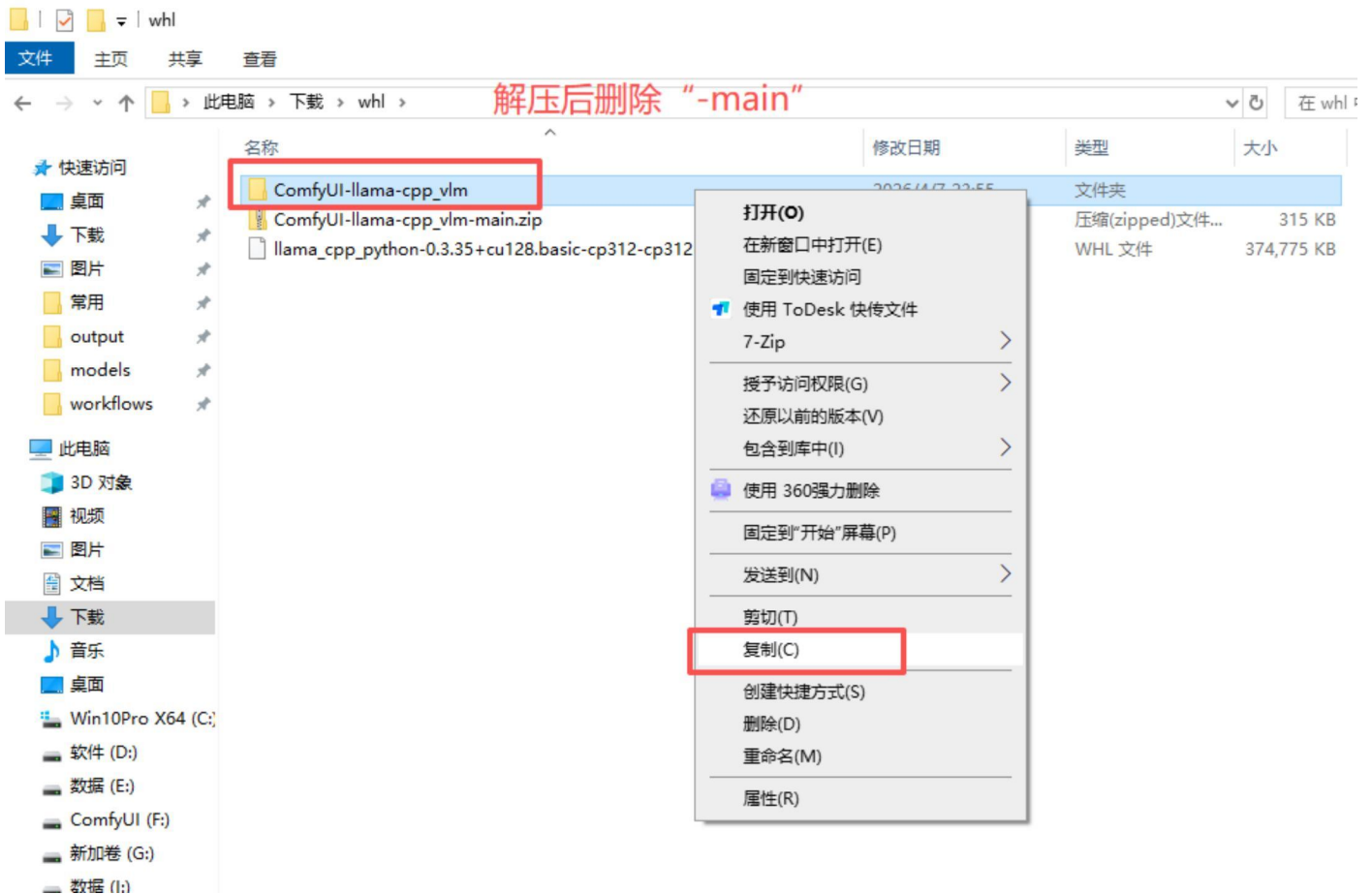
1. 打开网址复制 git 地址:

Github 网址: [https://github.com/lihaoyun6/ComfyUI-llama-cpp\\_vlm](https://github.com/lihaoyun6/ComfyUI-llama-cpp_vlm)

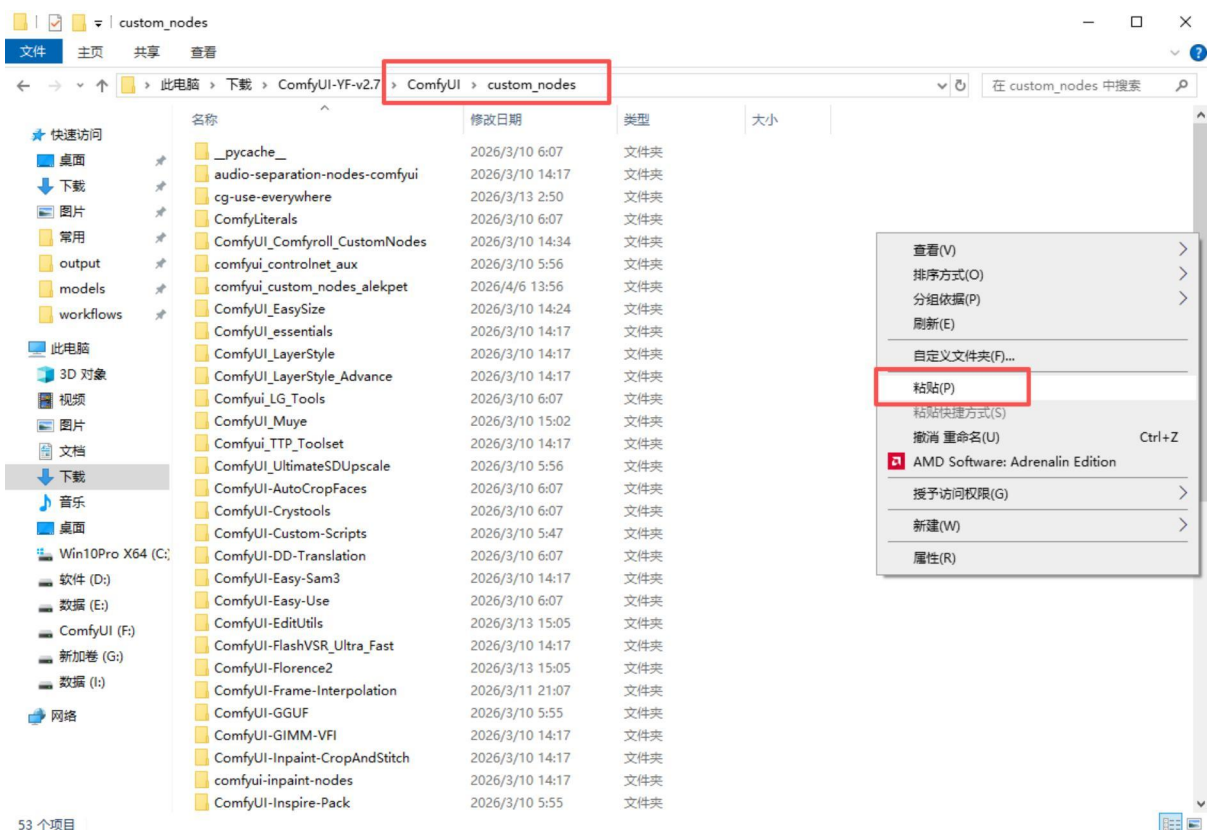
2. 下载 Zip 压缩包



3. 解压刚刚下载的 Zip 压缩包, 然后删除后缀 "-main"

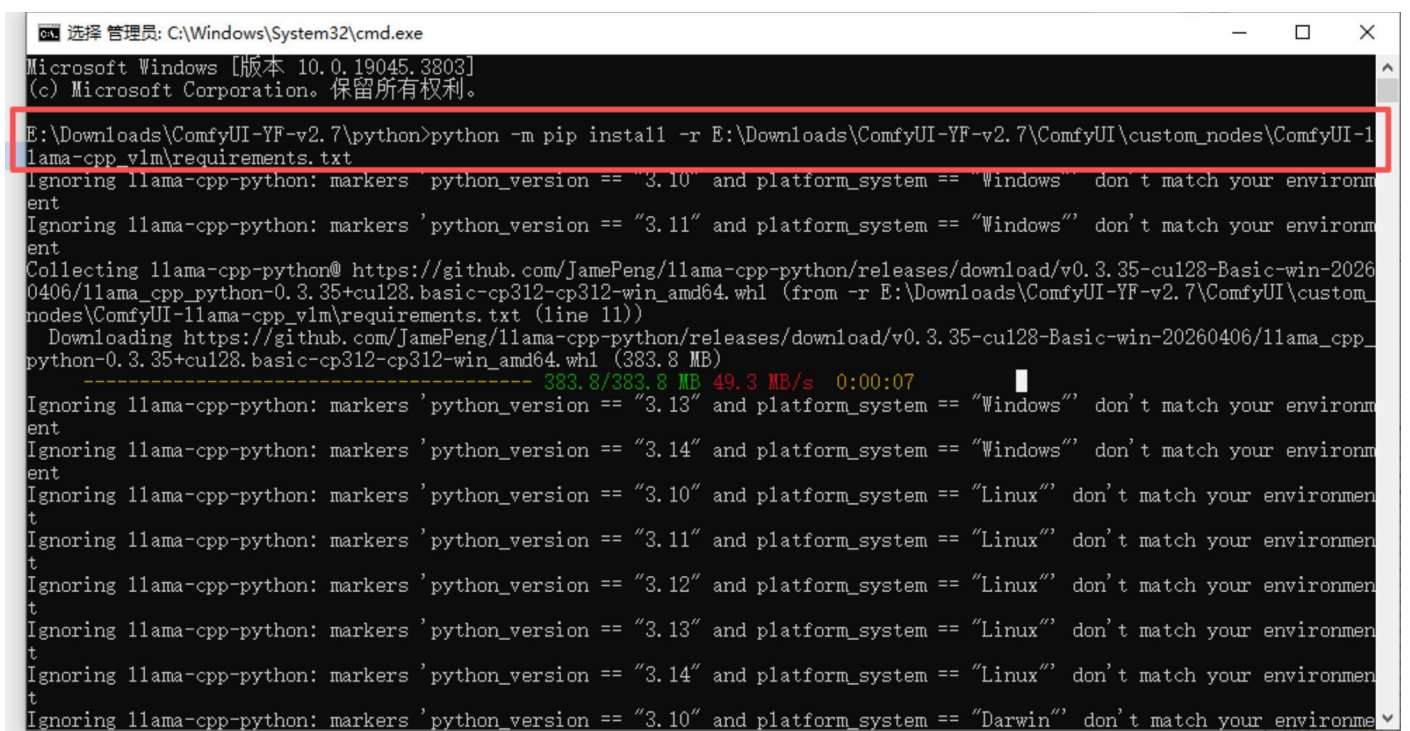
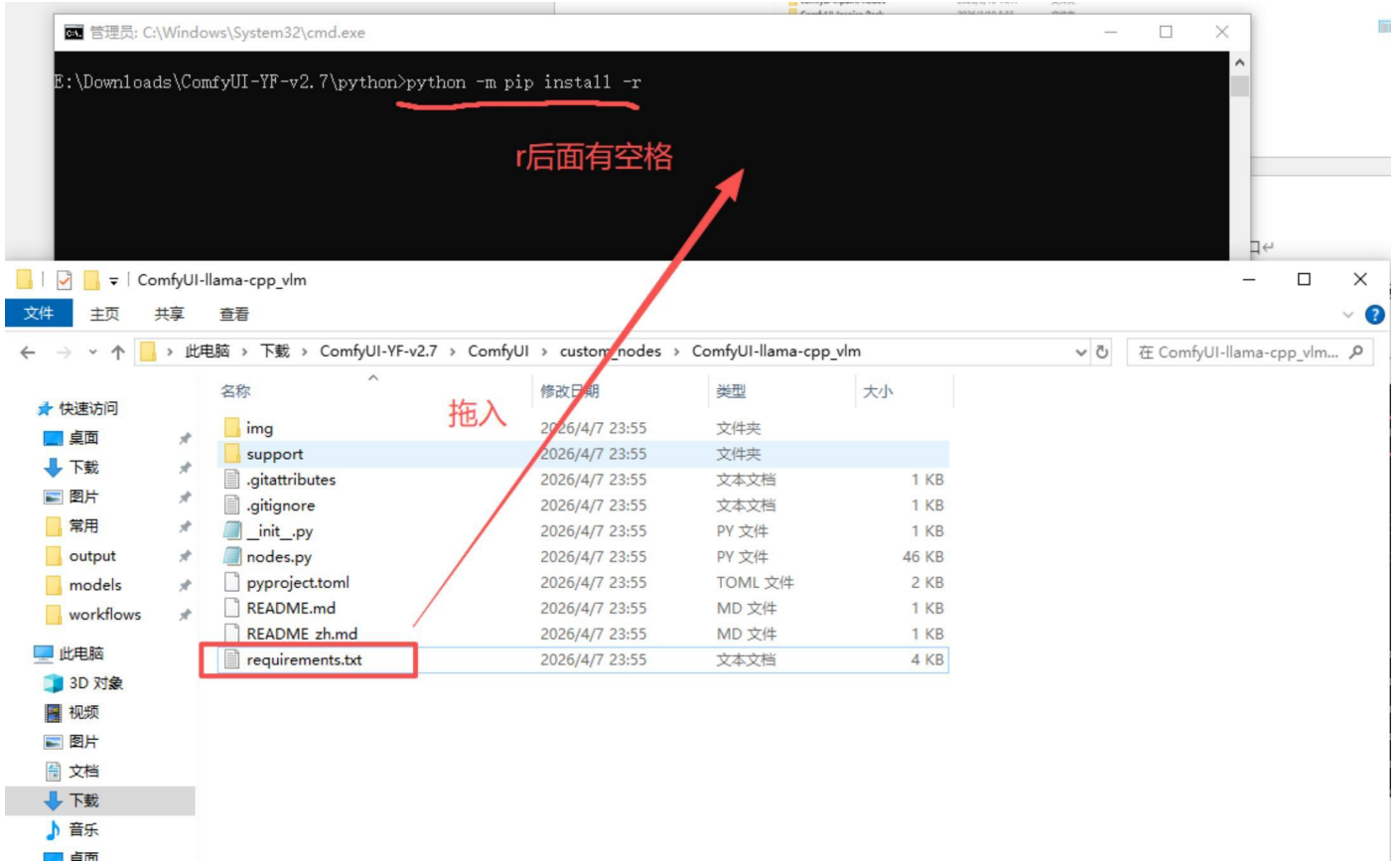


4. 打开 ComfyUI 插件文件夹 custom\_nodes 所在目录 X:\ ComfyUI-YF-v2.7\ComfyUI\custom\_nodes, 将刚刚解压改名的文件夹复制进去



## 3.2 安装插件依赖:

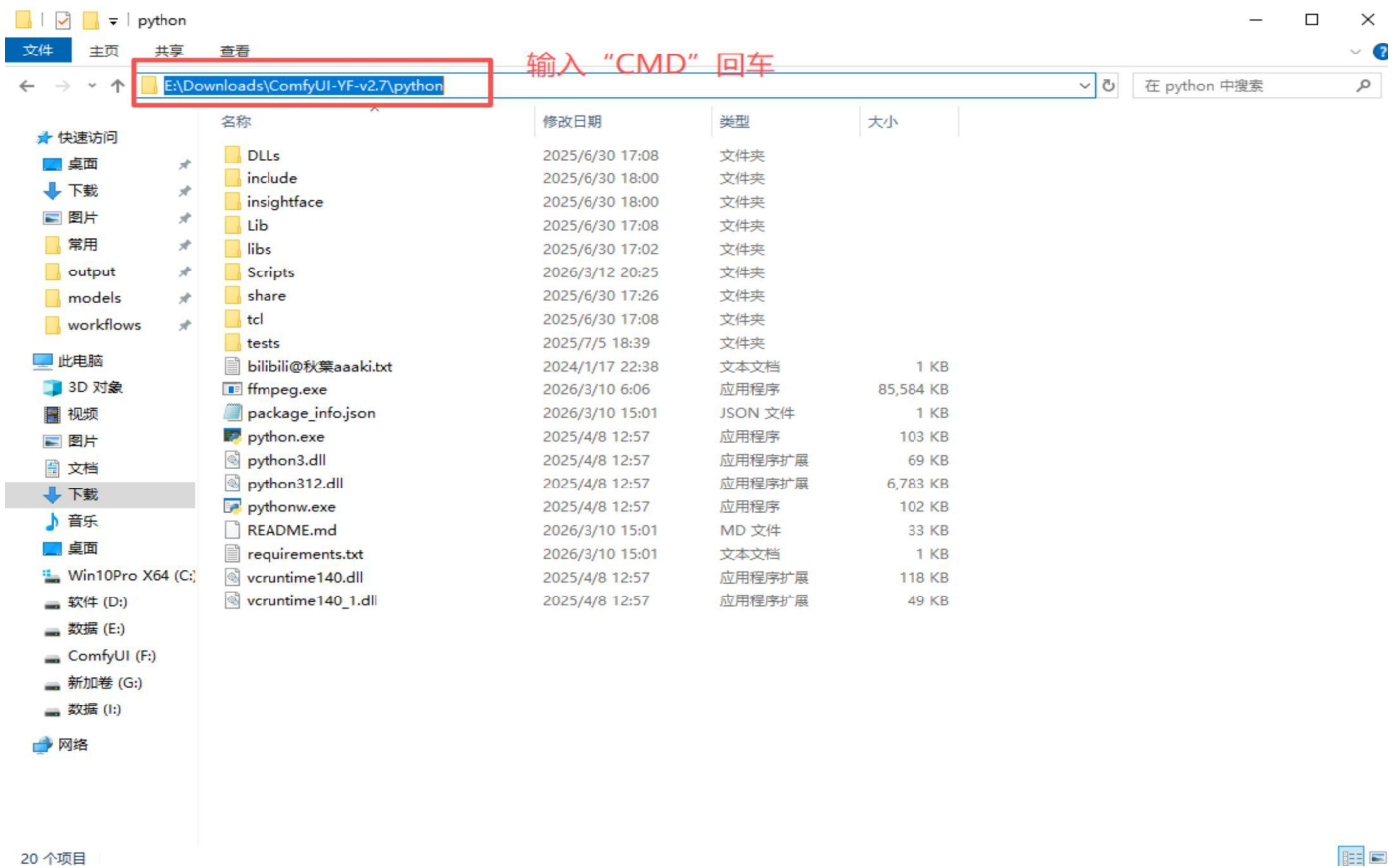
1. 进入 ComfyUI\python 文件夹，在地址栏输入 CMD 回车，打开命令提示符窗口
2. 输入 `python -m pip install -r` [然后将 omfyUI-llama-cpp\_vlm 文件夹中的 requirements.txt 文件拖动进 CMD 黑框中，回车]



### 3.3 安装轮子

1. 查看 python 版本、CUDA 版本

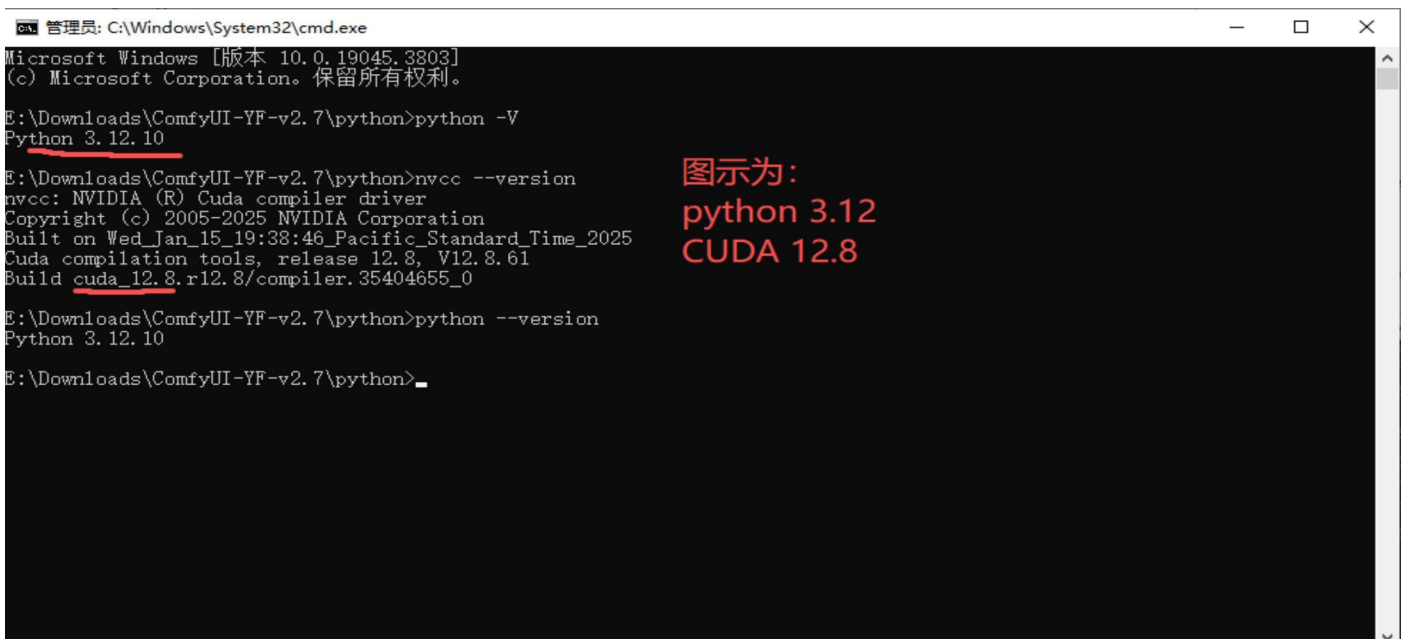
打开 ComfyUI/python 文件夹所在目录，在地址栏输入“CMD”回车



2. 分别输入以下命令查看 python 版本、CUDA 版本

```
python -V
```

```
nvcc --version
```



3. 下载对应版本的轮子.whl 文件，此处以上图版本为例

打开网址：<https://github.com/JamePeng/llama-cpp-python/releases>

4. 找到适配自己 python 版本+CUDA 版本的轮子，点击蓝字下载

**v0.3.35-cu128-Basic-win-20260406**

Bump version to 0.3.35

Release Note: [\[0.3.35\] Release Note: Gemma 4 series & LFM 2.5-VL Support, OpenAI OpenAPI Alignment and Logging Architecture Migration](#)

Changelog see here: [0.3.35 Changelog](#)

Signed-off-by: JamePeng [jame\\_peng@sina.com](mailto:jame_peng@sina.com)

win代表适配windows  
cu128代表适配CUDA12.8  
cp312代表适配python3.12

▼ Assets 7

llama_cpp_python-0.3.35+cu128.basic-cp310-cp310-win_a...	sha256:b7162d568e1a51060...		366 MB	last week
llama_cpp_python-0.3.35+cu128.basic-cp311-cp311-win_am...	sha256:ed93ffa15dd91702e...		366 MB	last week
llama_cpp_python-0.3.35+cu128.basic-cp312-cp312-win_a...	sha256:e7a74eb8eff0ee4a8...		366 MB	last week
llama_cpp_python-0.3.35+cu128.basic-cp313-cp313-win_a...	sha256:9f80131e5936a2996...		366 MB	last week
llama_cpp_python-0.3.35+cu128.basic-cp314-cp314-win_a...	sha256:72d3c0292fcb4af1e...		366 MB	last week
Source code (zip)				last week
Source code (tar.gz)				last week

👍 1 person reacted

5. 卸载当前安装的轮子（由于不知道你之前是否安装对了轮子，所以先卸载）：

输入卸载命令：`python -m pip uninstall llama_cpp_python`

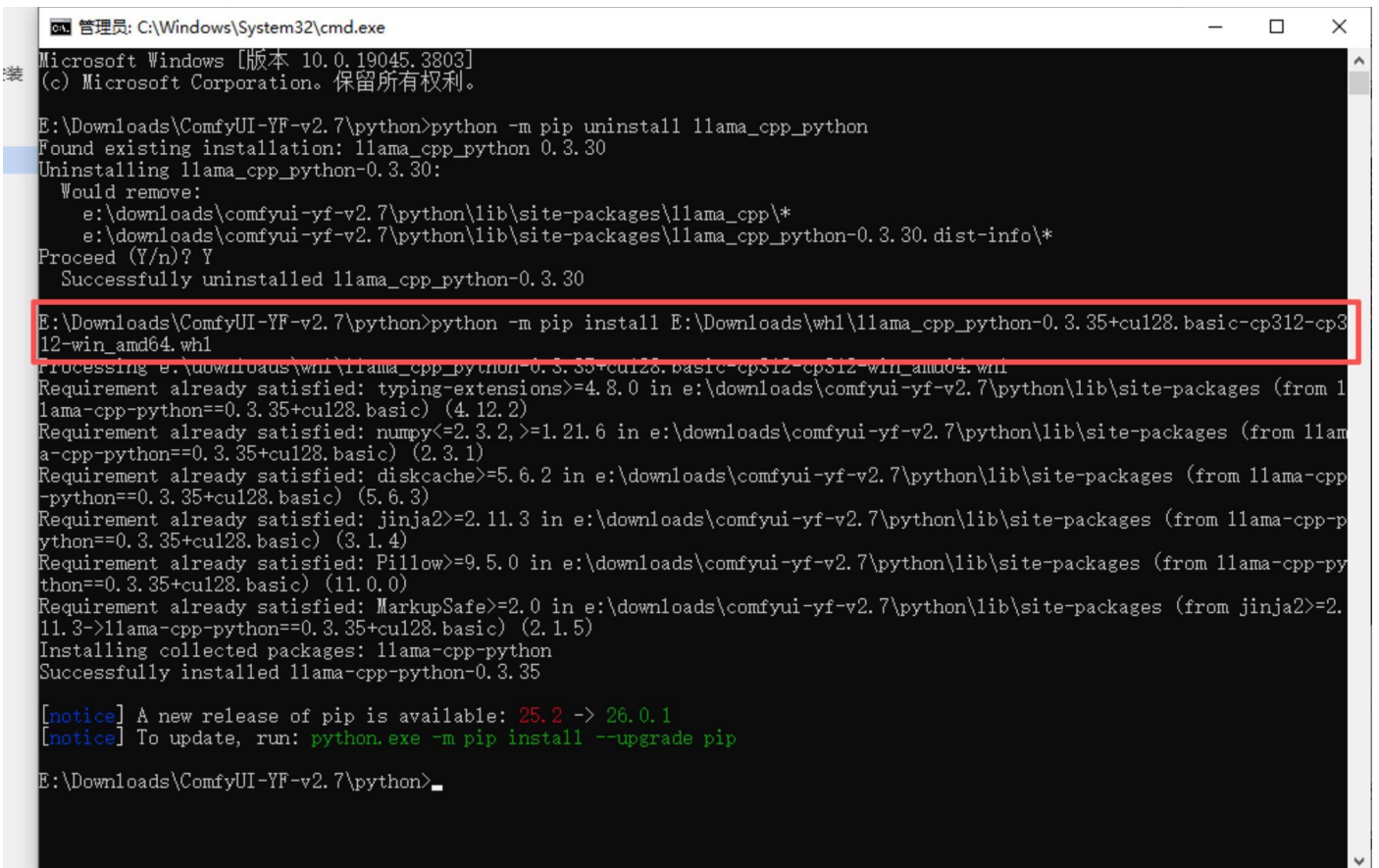
```
管理员: C:\Windows\System32\cmd.exe
Microsoft Windows [版本 10.0.19045.3803]
(c) Microsoft Corporation. 保留所有权利。

E:\Downloads\ComfyUI-YF-v2.7\python>python -m pip uninstall llama_cpp_python
Found existing installation: llama_cpp_python 0.3.30
Uninstalling llama_cpp_python-0.3.30:
  Would remove:
    e:\downloads\comfyui-yf-v2.7\python\lib\site-packages\llama_cpp\*
    e:\downloads\comfyui-yf-v2.7\python\lib\site-packages\llama_cpp_python-0.3.30.dist-info\*
Proceed (Y/n)? Y
Successfully uninstalled llama_cpp_python-0.3.30

E:\Downloads\ComfyUI-YF-v2.7\python>
```

## 6. 安装正确版本的轮子:

输入安装命令: `python -m pip install` [然后将下载好的轮子文件.whl 拖动进 CMD 黑框中, 回车]



等待安装结束即可。