# 金水Task2笔记
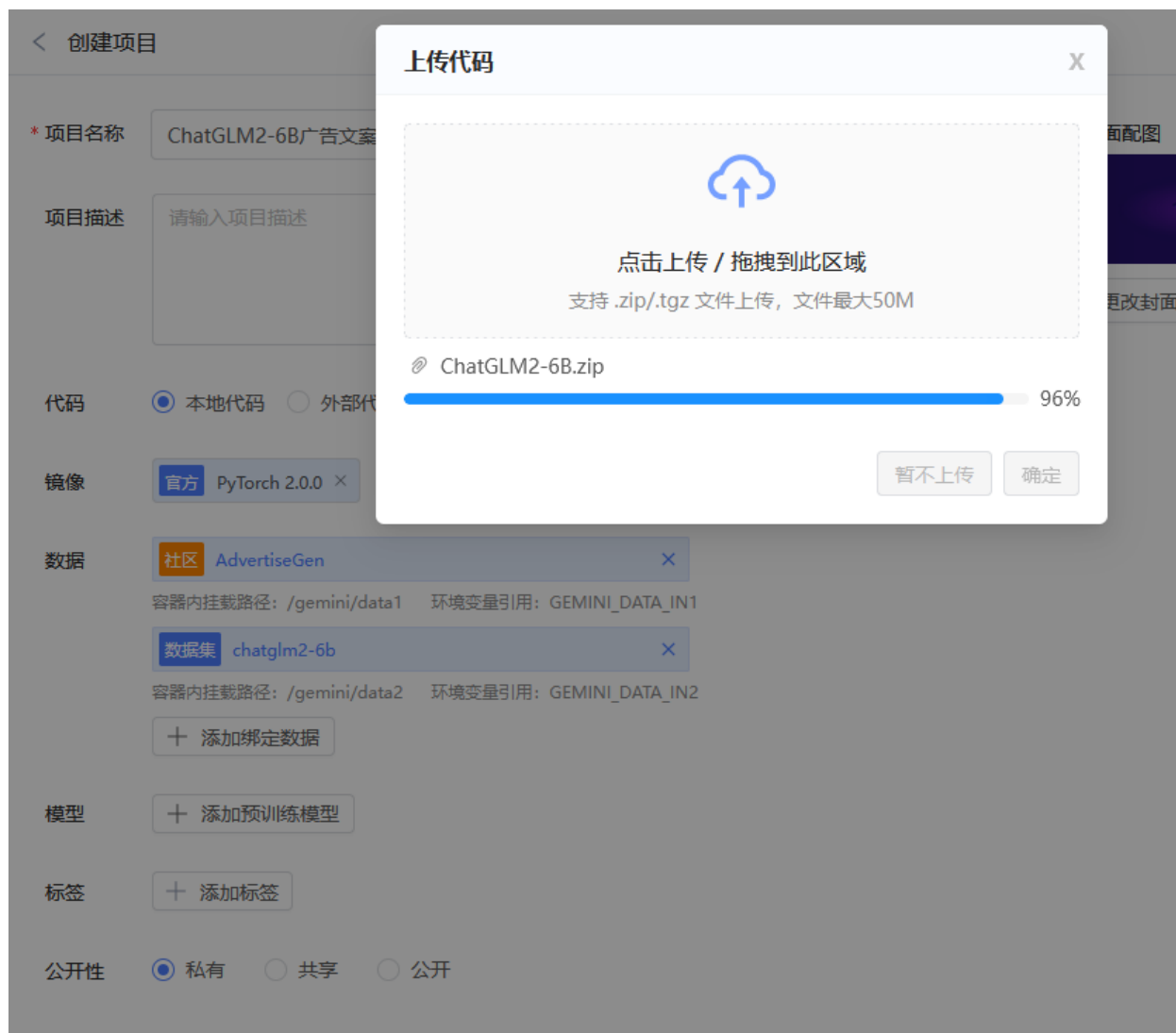
1.首先和Task1流程类似，先配置好环境，数据，和要运行的代码



2.镜像换为chatglm2-6b,并进行端口77的开放

# ID: 360747710078259200

**开发环境实例**

修改实例化规格　修改挂载数据　修改镜像　修改 SSH 配置　修改最长运行时间

**实例配置**

| 实例 | GPU | 挂载数据 | 代码版本 | 镜像 |
|---|---|---|---|---|
| B1.large | 1 | 2 数据集 | latest | ChatGLM2-... |

**运行记录**　历史记录　事件

○ 创建　2023-09-23 21:22

○ 等待中

　✓ 等待配额　　　　　　　　　　　　　2023-09-23 21:22

　✓ 任务调度　　　　　　　　　　　　　2023-09-23 21:22

　　进入任务池 ✓　　　　　　　　　　　2023-09-23 21:22

　　等待调度 ✓　　　　　　　　　　　　2023-09-23 21:22

　　▼ 任务调度 ✓　　　　　　　　　　　2023-09-23 21:22

　　代码准备 ✓　　　　　　　　　　　　2023-09-23 21:22

状态　　　　　　　　　　　停止运行

○ 等待中 ⓘ　　　　　　　　　设为紧急

开放端口 ⓘ　　　　　　　　　添加

外部访问地址需在开发环境启动，且程序
正常运行后使用。 查看文档

内部端口：77 ⓘ ｜ TCP 协议

外部访问：需在开发环境启动后分配

使用用途：chatglm2-6b 网页形式访问

开发者工具

🖥 JupyterLab　　　　　　　　 ＞

🖥 网页终端　　　　　　　　　 ＞

开始时间

2023-09-23 21:25

自动停止

已设置最大运行时长为8 h

## 3,改下requirements.txt



```
requirements.txt                    ×    +

1   protobuf
2   transformers==4.30.2
3   cpm_kernels
4   torch>=2.0
5   gradio
6   mdtex2html
7   sentencepiece
8   accelerate
9   sse-starlette
10  streamlit>=1.24.0
11  rouge_chinese
12  nltk
13  jieba
14  datasets
```

## 4.在终端运行

```
pip install -r ChatGLM2-6B/requirements.txt -i https://pypi.virtaicloud.com/repository/pypi/sim
```
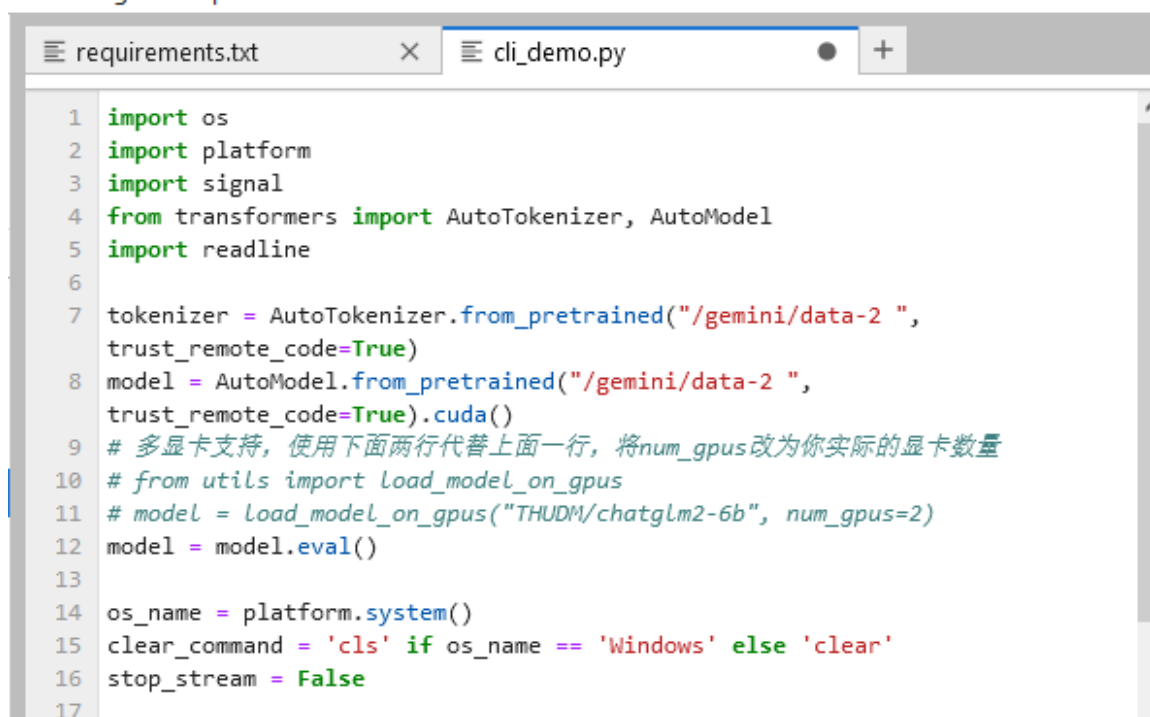
## 安装成功

```
t (line 5)) (0.19.3)
Collecting mdurl~=0.1 (from markdown-it-py>=2.2.0->rich<14,>=10.14.0->streamlit>=1.24.0->-r ChatGLM2
-6B/requirements.txt (line 10))
  Downloading https://pypi.virtaicloud.com/repository/pypi/packages/mdurl/0.1.2/mdurl-0.1.2-py3-none
-any.whl (10.0 kB)
Installing collected packages: watchdog, validators, tzlocal, toml, tenacity, smmap, safetensors, md
url, blinker, pydeck, markdown-it-py, gitdb, transformers, sse-starlette, rich, gitpython, streamlit
  Attempting uninstall: transformers
    Found existing installation: transformers 4.27.1
    Uninstalling transformers-4.27.1:
      Successfully uninstalled transformers-4.27.1
Successfully installed blinker-1.6.2 gitdb-4.0.10 gitpython-3.1.37 markdown-it-py-3.0.0 mdurl-0.1.2
pydeck-0.8.1b0 rich-13.5.3 safetensors-0.3.3.post1 smmap-5.0.1 sse-starlette-1.6.5 streamlit-1.27.0
tenacity-8.2.3 toml-0.10.2 transformers-4.30.2 tzlocal-5.0.1 validators-0.22.0 watchdog-3.0.0
WARNING: Running pip as the 'root' user can result in broken permissions and conflicting behaviour w
ith the system package manager. It is recommended to use a virtual environment instead: https://pip.
pypa.io/warnings/venv

[notice] A new release of pip is available: 23.1.1 -> 23.2.1
[notice] To update, run: python -m pip install --upgrade pip
root@360748546191790080-taskrole1-0:/gemini/code#
```

5.修改 cli_demo.py 中模型的地址，将 `THUDM/chatglm2-6b` 替换为 `/gemini/data-2` （即模型实际挂载到环境的地址）

```
≣ requirements.txt          ×   ≣ cli_demo.py          ●   +

 1  import os
 2  import platform
 3  import signal
 4  from transformers import AutoTokenizer, AutoModel
 5  import readline
 6
 7  tokenizer = AutoTokenizer.from_pretrained("/gemini/data-2 ",
    trust_remote_code=True)
 8  model = AutoModel.from_pretrained("/gemini/data-2 ",
    trust_remote_code=True).cuda()
 9  # 多显卡支持，使用下面两行代替上面一行，将num_gpus改为你实际的显卡数量
10  # from utils import load_model_on_gpus
11  # model = load_model_on_gpus("THUDM/chatglm2-6b", num_gpus=2)
12  model = model.eval()
13
14  os_name = platform.system()
15  clear_command = 'cls' if os_name == 'Windows' else 'clear'
16  stop_stream = False
17
```

6.1. 切换至 网页终端。执行如下命令唤醒交互式对话。等待最终 `Loading checkpoint shards: 100%` 且返回 `用户：` 字样。。

```
python ChatGLM2-6B/cli_demo.py
```

成功