

在AutoDL中快速部署ChatGLM3-6B模型

本节内容介绍的是在AutoDL平台上租赁算力资源，并完成ChatGLM3-6B的部署调用的流程。根据官方介绍，ChatGLM3-6B 目前支持GPU运行（需要英伟达显卡）、CPU运行以及Apple M系列芯片运行。其中GPU运行需要至少6GB以上显存（4Bit精度运行模式下），而CPU运行则需要至少32G的内存。而由于Apple M系列芯片需要最少13G内存运行。在正式安装之前，需要先确保拥有足够的算力资源，以下推荐的是一种轻量化的部署方式，非常适合入门级的测试开发：AutoDL租赁服务器进行快速部署。

1. 物理机 or 云服务

- 完全小白，对大模型技术没有了解，建议用新人账号白嫖各大云服务平台的免费算力，再考虑购买或者租赁。
- 如果经常做微调实验，或实验室学生系统学习，有自己的物理机将更加方便，按照学习实践部分内容采购即可。
- 为用户提供相关的推理服务，首选云服务，有更大参数量，更好性能的模型选择，随用随停，按量计费。
- “独角兽”公司AI应用/大模型AI技术创新公司……，需要大规模大批量的微调训练或者对内/对外提供大量推理服务，按需配备高性能GPU服务器。

物理机部分大家可以按照前序了解的自行购买，但是这里再次强调，**购买需谨慎**，尤其在二手平台购买二手显卡需要更仔细专业的判断。目前国内市场也会有A、H系列显卡流通，可能是存货、二手、……渠道，但是这类高性能显卡要更专业细致的判断，谨防被骗。

云服务厂商

国内主流

- 阿里云：<https://www.aliyun.com/product/ecs/gpu>
- 腾讯云：<https://cloud.tencent.com/product/gpu>
- 火山引擎：<https://www.volcengine.com/product/gpu>

国外主流

- AWS: <https://aws.amazon.com/cn/campaigns/aws-gpu/>
- Vultr: <https://www.vultr.com/products/cloud-gpu/>
- TPU: https://cloud.google.com/gpu/?hl=zh_cn

算力平台

主要适用于学习和训练，不适用于企业级部署提供服务。

- ModelScope: 阿里出品，中国的“HuggingFace”，模型开源社区，绑定阿里云有（24GB显存+36小时）GPU环境。 <https://www.modelscope.cn/home>
- Colab: 谷歌出品，升级服务仅需 9 美金。 <https://colab.research.google.com/>
- Kaggle: 免费，每周 30 小时 T4，P100 可用。 <https://www.kaggle.com/>
- AutoDL: 价格亲民，支持 Jupyter Notebook 及 ssh，国内可用。
<https://www.autodl.com/home>

2. 算力准备

在正式安装之前，需要先确保拥有足够的算力资源，以下推荐的是一种轻量化的部署方式，非常适合入门级的测试开发：首先可以从以下链接中进入AutoDL的官方网址，在右上角的选项里可以注册/登录。

<https://www.autodl.com/home>

AutoDL AI算力云

弹性、好用、省钱

立即注册 了解详情

注册礼包
注册立送30天会员

GPU选型
如何选择合适的GPU

开具发票
简单快速开具发票

新手入门
简单几步，创建实例

炼丹会员及租用价格

AutoDL坚持为您提供服务稳定、价格公平的GPU租用服务。更为学生提供免费升级会员通道，享极具性价比的会员价格。 [如何升级会员？](#)



登录

注册

+86 请输入手机号

请输入验证码

请输入8~16位包含数字和字母组合

☒ 我已阅读并同意 [《AutoDL服务协议》](#) 和 [《隐私协议》](#)

已有账号, [点击登录](#)

进入界面之后点击右上角的用户信息可以查看余额和进行充值，其金额可以自定义。注意：只有账户有余额才能在后续算力市场租赁主机。

炼丹师0909



炼丹会员

认证学生升级炼丹师
等级与会员福利

成长值



距离升级还需44

[进入成长值主页](#)

炼丹师0909

未实名

ID: 2bd2ceda-18bb-4d5c-a8c7-679b33394f82



炼丹会员

可用余额: ¥ 42.00

冻结余额: ¥ 0.00

代金券: ¥ 0.00

容器实例: 1

充值

退出登录

费用信息

充值

账户余额: ¥42.00

充值金额:

¥50

成长值+50

¥100

成长值+100

¥500

成长值+500

¥1000

成长值+1000

¥2000

成长值+2000

其他金额

[了解会员成长值>>](#)

请输入充值额度

¥

支付方式:

微信支付

支付宝

对公汇款

[如何开发票?](#) [如何对公汇款?](#)

充值

领优惠券

消息

评论

分享

确认用户余额充裕后点击左上角的算力市场，租赁合适的主机，推荐的配置为：计费方式选择按量计费、地区任选、GPU型号选择RTX3090/24GB、GPU数量选择为1。

选择RTX3090/24GB卡的理由是ChatGLM3-6B的GPU运行需要至少6GB以上显存（4Bit精度运行模式下），而CPU运行则需要至少32G的内存。其中CPU运行模式下内存占用过大且运行效率较低，GPU模式部署才能有效的进行大模型的学习实践。基于性能和性价比进行考量，我们建议选择以上参数进行部署。

① 严禁使用WebUI等算法生成违禁图片、严禁挖矿，一经发现立即封号！

计费方式:

按量计费

包日

包周

包月

选择地区:

重庆A区

西北B区

北京A区

北京B区

佛山区

内蒙A区

3090专区

L20专区

V100专区

A800专区

GPU型号:

☐ 全部☒ RTX 3090 (279/1072)☐ V100-SXM2-32GB (1/32)☐ L40 (9/24)☐ RTX 2080 Ti (219/616)☐ RTX 3080 (74/400)☐ RTX A4000 (5/16)

GPU数量:

1

2

3

4

5

6

7

8

10

12

北京A区 / 607机 可租用至: 2024-10-01

RTX 3090 / 24 GB

空闲/总量 1 / 8

每GPU分配

CPU: 15 核, Xeon(R) Platinum 8358P

内存: 80 GB

硬盘

系统盘: 30 GB

数据盘: 50 GB, 可扩容 560 GB

其它

GPU驱动: 525.105.17

CUDA版本: ≤ 12.0 ?

¥1.58/时

¥4.66/时

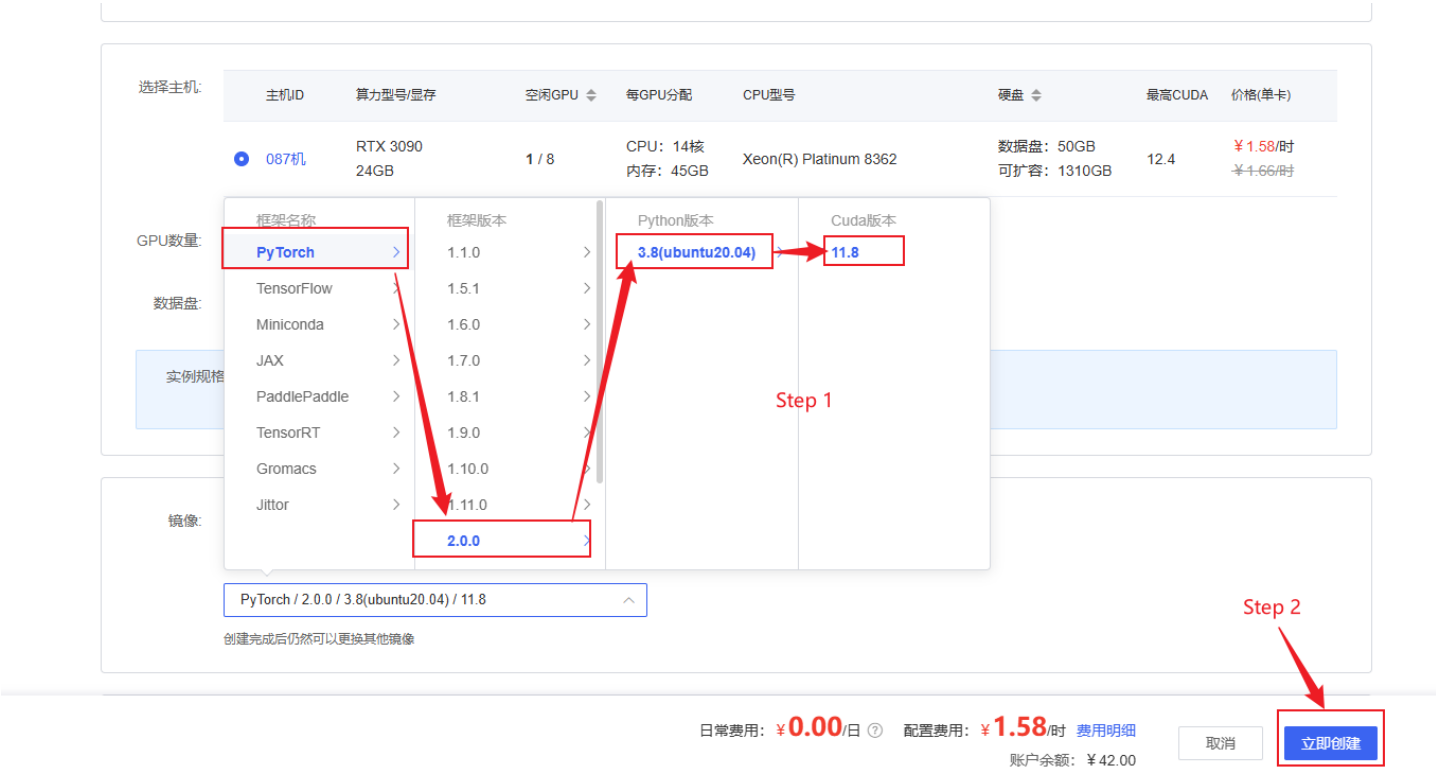
9.5折

会员最低享9.5折 ¥1.58/时

1卡可租

选择合适的主机后需要在下方的镜像栏中选择适合的框架——框架名称：PyTorch，框架版本：2.0.0，Python版本：3.8（Ubuntu20.04），Cuda版本11.8.选择好之后点击右下角的立即创建便可完成配置。

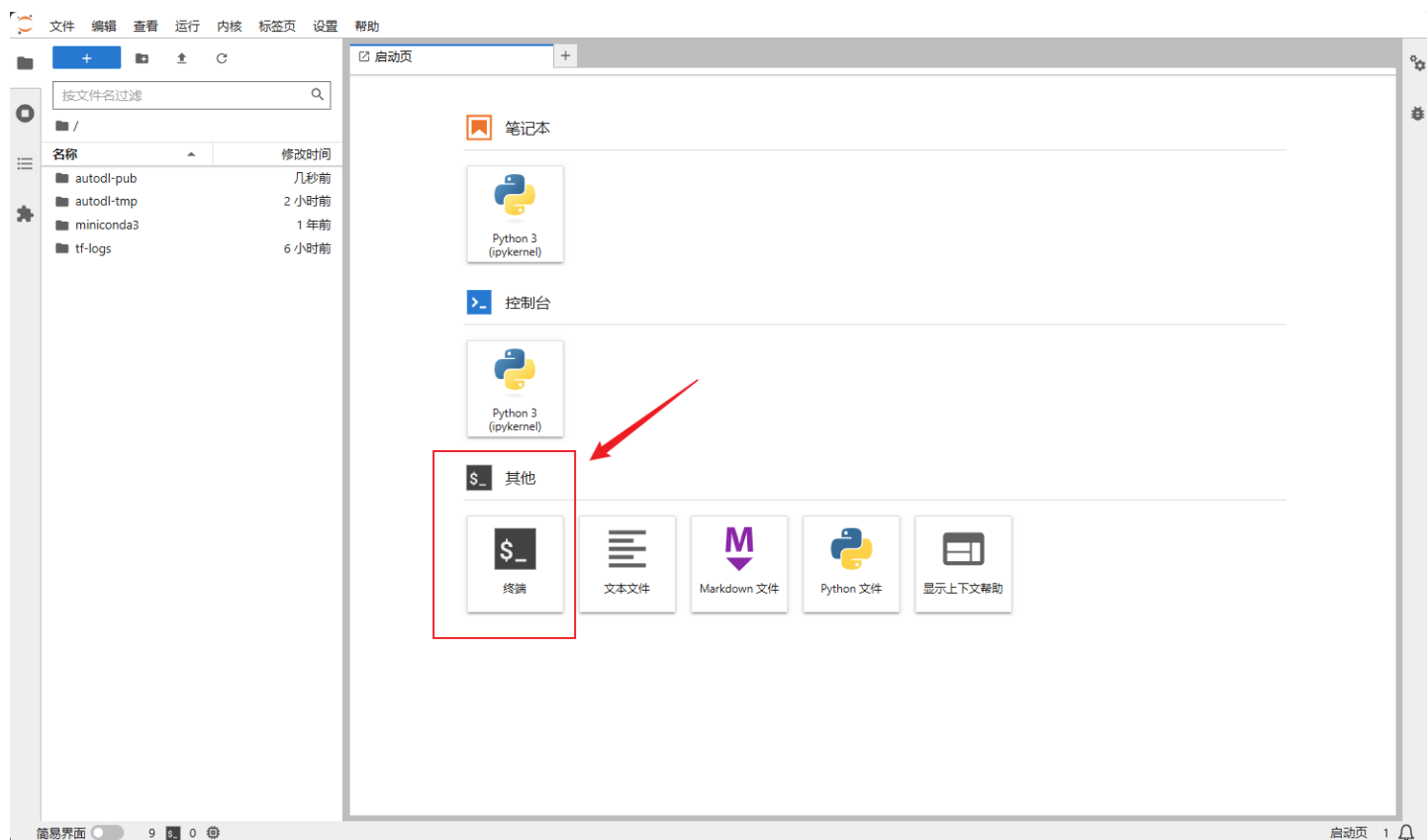
其中 PyTorch 是一个流行的深度学习框架，支持大规模模型的训练和推理。Python 3.8 是一个稳定且常用的版本，兼容大多数机器学习库和工具。选择 Ubuntu 20.04 作为操作系统版本是因为其长期支持和广泛使用，特别适合在生产环境中部署。Cuda 是 NVIDIA 提供的并行计算平台和编程模型，支持 GPU 加速。选择 Cuda 11.8 版本是因为它与 PyTorch 2.0.0 兼容。



创建完成后，点击左边栏的容器实例便可随时找到配置好的实例，在快捷工具栏中点击Jupyter lab开始模型的安装部署。



3. 换源和安装依赖包



进入Jupyter lab打开终端开始环境配置，首先要进行的是 pip 换源和安装依赖包。点击启动终端，在其中逐行输入以下代码以实现功能。

在终端通过命令升级 pip，确保使用的是最新版本的 pip，这样可以避免在安装库时出现兼容性问题。更换 pip 的默认源为清华大学的镜像源，以加速 Python 库的下载和安装。

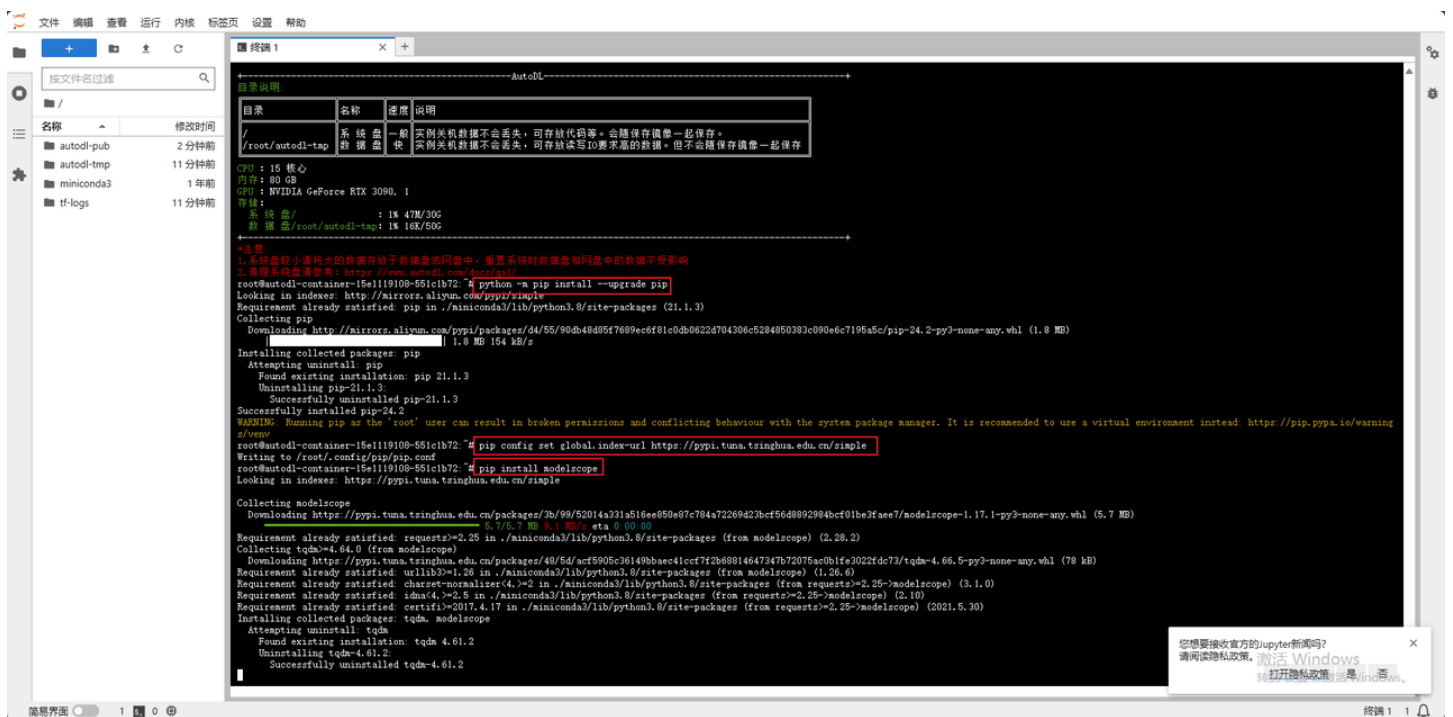
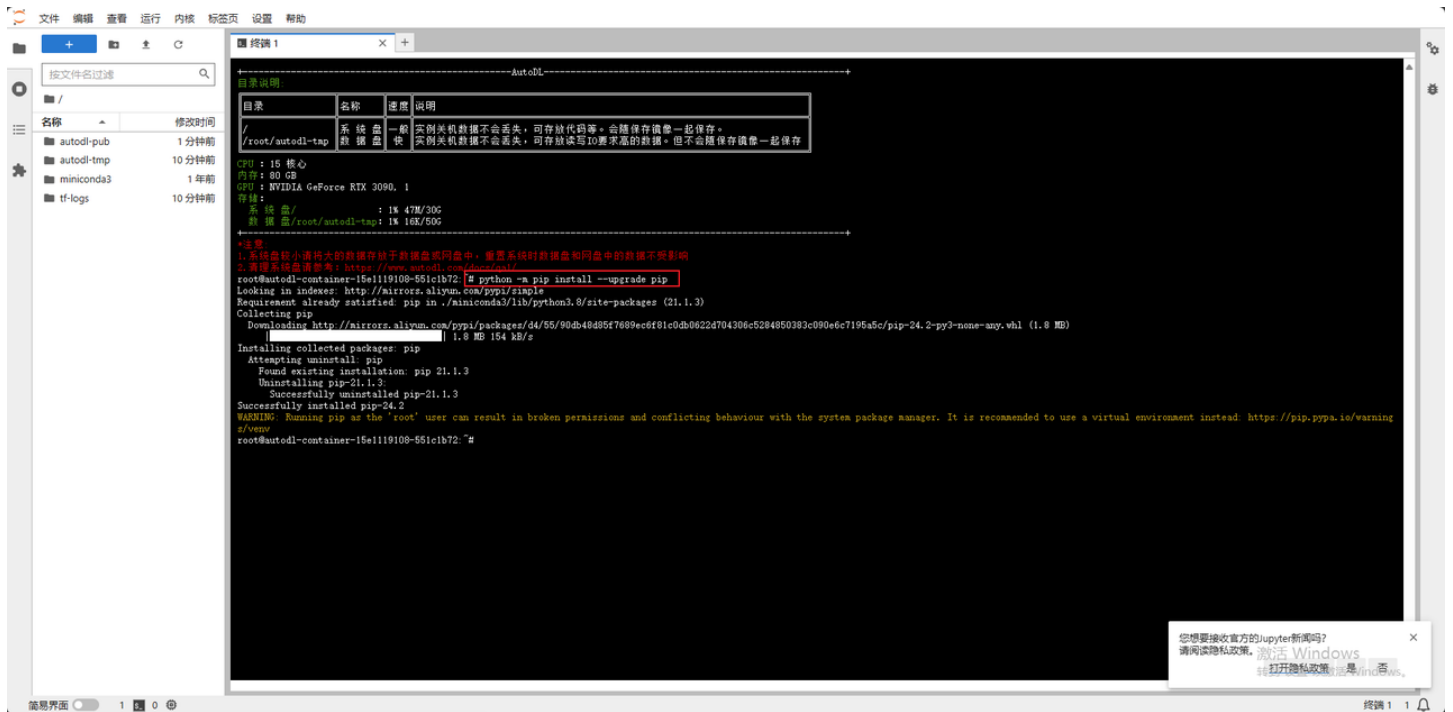
以下是安装的库的介绍：

modelscope: 用于模型推理和部署的库，支持多种机器学习和深度学习模型。

transformers: 包含了大量预训练的 Transformer 模型，包括 BERT、GPT 等等。

sentencepiece: 一个用于处理文本的库，特别是对子词单元进行分词操作，常用于自然语言处理任务。

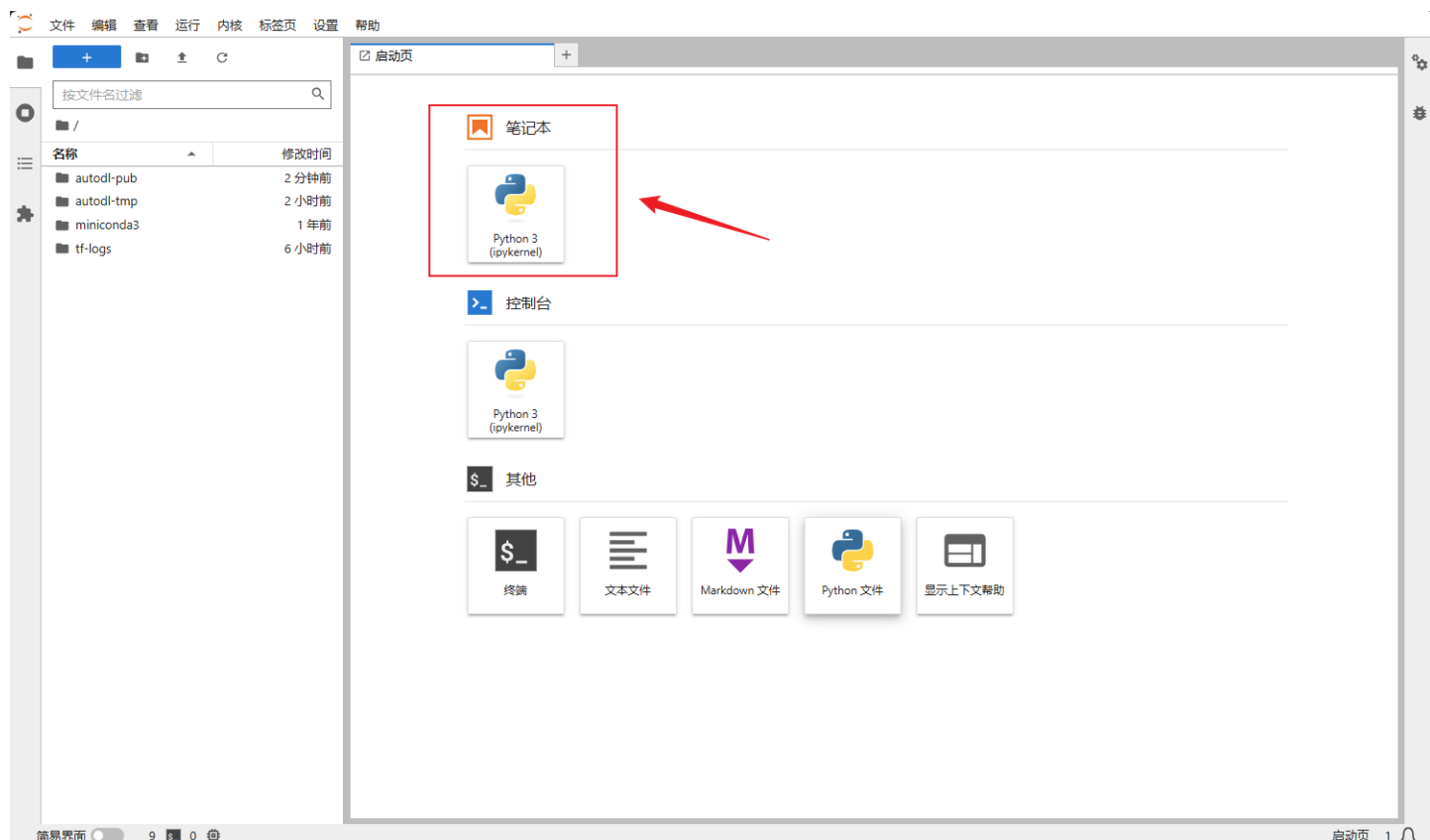
```
1 # 升级pip
2 python -m pip install --upgrade pip
3 # 更换 pypi 源加速库的安装
4 pip config set global.index-url https://pypi.tuna.tsinghua.edu.cn/simple
5
6 pip install modelscope
7 pip install transformers
8 pip install sentencepiece
```





4. 模型下载

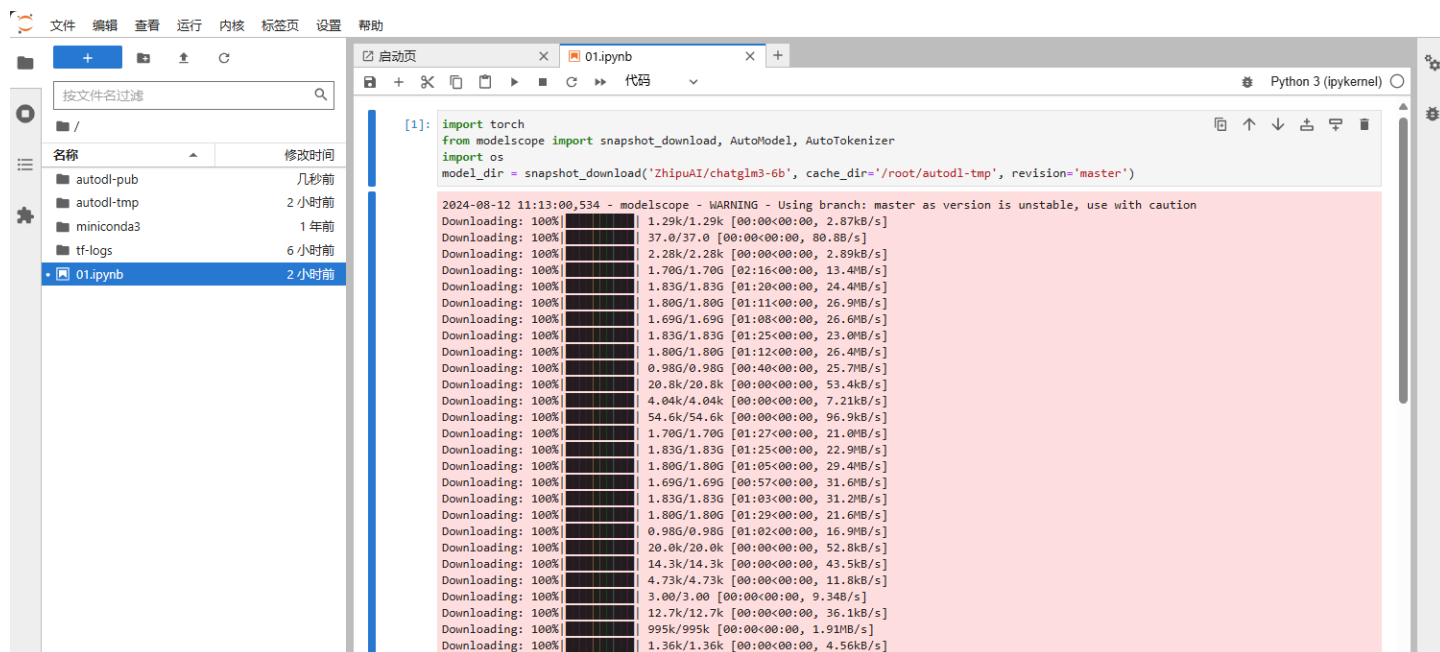
在启动页打开新的Jupyter notebook进行模型的下载。



这里选择的是使用 modelscope 中的 snapshot_download 函数下载模型，这个函数中的第一个参数为模型名称，第二个参数 cache_dir 为模型的下载路径。

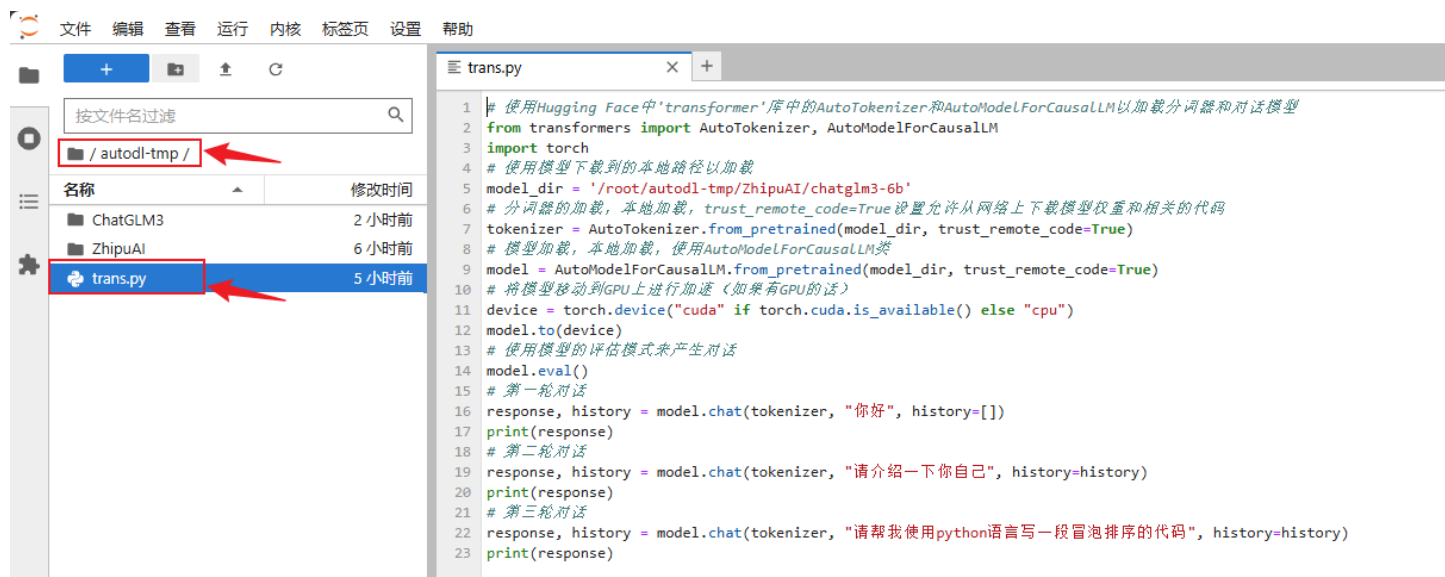
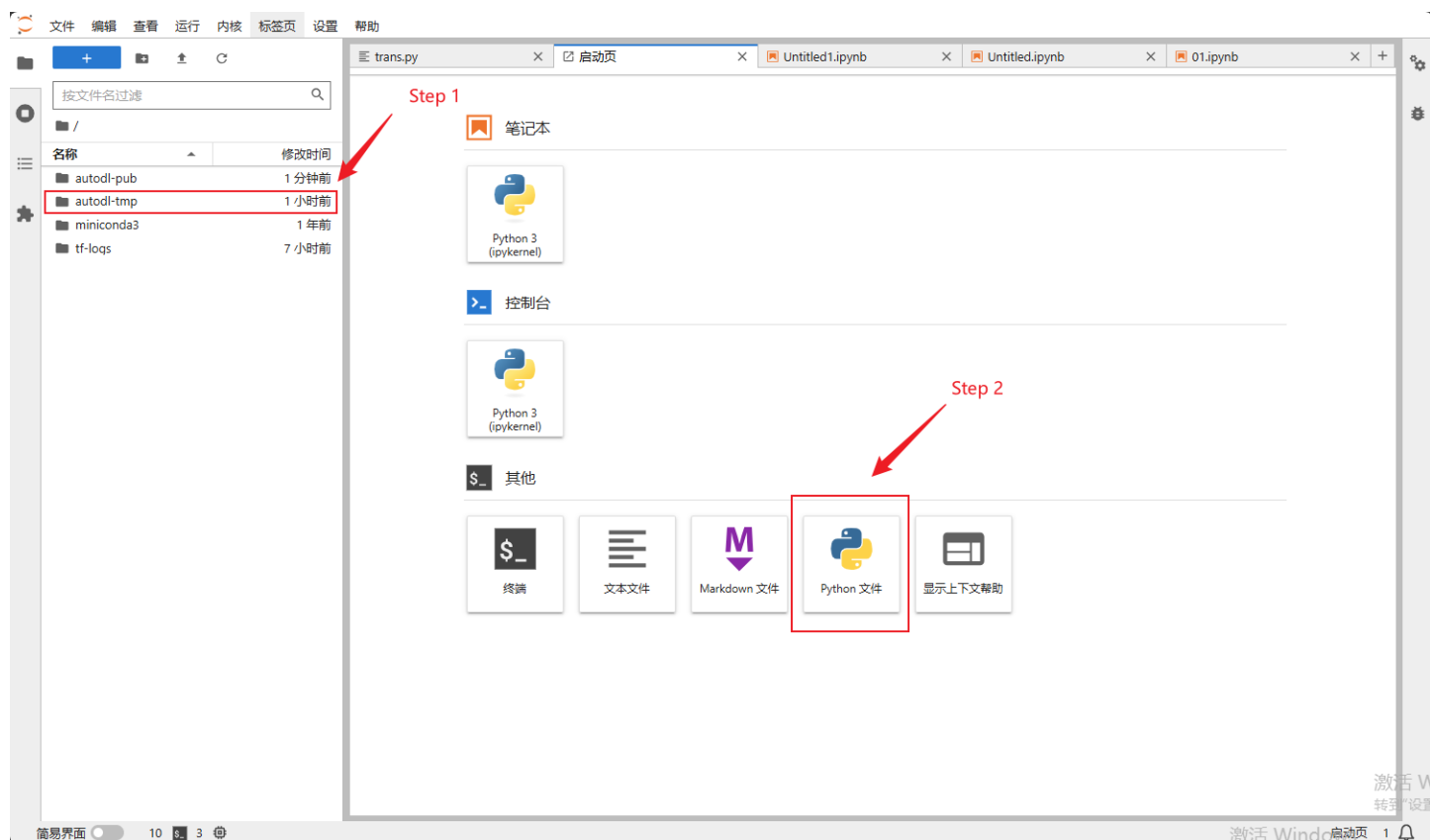
打开Jupyter Lab文件执行以下代码进行下载，ChatGLM3-6B模型大小为 14 GB，下载模型大概需要 20~25 分钟。

```
1 import torch
2 from modelscope import snapshot_download, AutoModel, AutoTokenizer
3 import os
4 model_dir = snapshot_download('ZhipuAI/chatglm3-6b', cache_dir='/root/autodl-tmp', revision='master')
```



5. 启动模型的代码

在/root/autodl-tmp路径下新建trans.py文件并在其中输入以下内容



```
1 from transformers import AutoTokenizer, AutoModelForCausalLM # 使用Hugging Face
  中'transformer'库中的AutoTokenizer和AutoModelForCausalLM以加载分词器和对话模型
2 import torch
3 model_dir = '/root/autodl-tmp/ZhipuAI/chatglm3-6b' # 使用模型下载到的本地路径以加
  载
4 tokenizer = AutoTokenizer.from_pretrained(model_dir, trust_remote_code=True)
  # 分词器的加载，本地加载，trust_remote_code=True设置允许从网络上下载模型权重和相关的代
  码
5 model = AutoModelForCausalLM.from_pretrained(model_dir,
  trust_remote_code=True) # 模型加载，本地加载，使用AutoModelForCausalLM类
6 device = torch.device("cuda" if torch.cuda.is_available() else "cpu") # 将模型移
  动到GPU上进行加速 (如果有GPU的话)
```

```

7 model.to(device)
8 model.eval() # 使用模型的评估模式来产生对话
9 # 第一轮对话
10 response, history = model.chat(tokenizer, "你好", history=[])
11 print(response)
12 # 第二轮对话
13 response, history = model.chat(tokenizer, "请介绍一下你自己", history=history)
14 print(response)
15 # 第三轮对话
16 response, history = model.chat(tokenizer, "请帮我使用python语言写一段冒泡排序的代码", history=history)
17 print(response)

```

6. 部署运行

需要注意的是，如果transformers的版本不匹配会导致报错，因此我们需要先降其版本。回到启动页打开终端分别输入以下指令将版本确定至4.37.2：

```

1 pip uninstall transformers #卸载当前版本
2 pip install --upgrade transformers==4.37.2 #安装指定版本

```

```

root@autodl-container-15e119108-551c1b72: /autodl-tmp# pip uninstall transformers
Found existing installation: transformers 4.43.0
Uninstalling transformers-4.43.0:
Would remove:
  /root/.miniconda3/bin/transformers-cli
  /root/.miniconda3/lib/python3.8/site-packages/transformers-4.43.0.dist-info/*
  /root/.miniconda3/lib/python3.8/site-packages/transformers/*
Proceed (Y/n)? y
Successfully uninstalled transformers-4.43.0
WARNING: Running pip as the 'root' user can result in broken permissions and conflicting behaviour with the system package manager, possibly rendering your system unusable. It is recommended to use a virtual environment instead: https://pip.pypa.io/warnings/venv. Use the --root-user-action option if you know what you are doing and want to suppress this warning.
root@autodl-container-15e119108-551c1b72: /autodl-tmp# pip install --upgrade transformers==4.41.2
Looking in indexes: https://pypi.tuna.tsinghua.edu.cn/simple
Collecting transformers==4.41.2
  Downloading https://pypi.tuna.tsinghua.edu.cn/packages/d9/b7/98f921d70102e2d38483bbb7013a689d2d646daa4495377bc910374ad727/transformers-4.41.2-py3-none-any.whl (9.1 MB)
    9.1 MB 31.8 MB/s eta 0:00:00

```

随后在终端输入以下指令，可以发现在平台上部署成功。可以看到，终端返回了前面trans.py文件提出的三个问题。

```

1 cd /root/autodl-tmp #将路径导向指定位置
2 Python trans.py #执行对应文件

```

```
root@autodl-container-15ell19108-55:~/b72 /autodl-tmp# python trans.py
Setting eos_token is not supported. use the default one.
Setting pad_token is not supported. use the default one.
Setting unk_token is not supported. use the default one.
Using cuda_device: -1, cuda_device: 100%|
你好! 我是人工智能助手 ChatGLM3-6B, 很高兴见到你, 欢迎向我任何问题。
我是一个人工智能助手 ChatGLM3-6B, 是清华大学 KEG 实验室和智谱 AI 公司于 2023 年共同训练的语言模型。我的任务是针对用户的问题和要求提供适当的答复和支持。我是一段程序, 并没有物理意义上的实体, 所以不能像真正的中国人那样和你进行交流。
当然可以, 以下是使用 Python 编写的冒泡排序代码:

'''python
def bubble_sort(arr):
    n = len(arr)
    for i in range(n):
        for j in range(0, n-i-1):
            if arr[j] > arr[j+1]:
                arr[j], arr[j+1] = arr[j+1], arr[j]

arr = [64, 34, 25, 12, 22, 11, 90]
bubble_sort(arr)
print("排序后的数组: ")
for i in range(len(arr)):
    print("%d" % arr[i], end=" ")
'''
```