1. 前期准备

本文档目的是帮助在本地linux服务器端,部署并成功运行Meta开源LLM llama2。基于本文档已经成功在A100和8卡4090服务器部署llama2-7b-chat-hf。

开始本文档,需要提前准备:

- 1. 支撑llama2运行的GPU服务器资源,使用基于Linux内核的服务器操作系统。
- 2. 确保可以通过远程连接,连接到服务器资源。(比如使用vscode remote-ssh)
- 3. 确保有足够的存储空间,最好有50G的富裕空间。
- 4. 服务器已经预装python3, CUDA等基础软件,由于本教程使用docker部署,所以不需要安装miniconda。

2. 下载llama-2

llama2可以在Meta官网和Huggingface网站上获取,为了安装的便捷性。我们最好遵从以下工作流完成:

1. Meta官网,使用美国节点+Gmail注册,并申请使用llama2的权限。

- 2. Huggingface官网,使用1中的gmail注册,并在Huggingface界面再次申请llama2权限。
- 3. 两个权限都申请通过以后, 就可以开始准备下载环节。

注意: 这里的两次注册最好使用同一个邮箱, 且确保自己在美国节点的代理条件之下, 预计两次审核时间不会超过一个小时。

Research Blog Resources About Q

Meta官网在完成如下表格的申请后,将会发送反馈邮件:

∞	Meta	

Request access to the next version of Llama

First Name	Last Name
Email	
Country	~
Organization / Affiliation	
Select the models you would like access to	x
Select the models you would like access to Llama 2 & Llama Chat Code Llama	x
Select the models you would like access to Llama 2 & Llama Chat Code Llama	x.
Select the models you would like access to Llama 2 & Llama Chat Code Llama Lama 2 Version Release Date: July 18, 2023	x

6	🔒 meta-Ilama/Llama-2-7b-chat-hf 🗙 🕂	- 0							
\leftarrow (https://huggingface.co/meta-llama/Llama-2-7b-chat-hf	A G 🗘 C C I D 🕼 🐨 📽 🚇 …							
	Hugging Face Q Search models, datasets, users	A C C C C C C C C C C C C C C C C C C C							
	🕫 meta-llama/Llama-2-7b-chat-hf 🖆 🛛 🖓 🗽 💷								
	😳 Text Generation 🔒 Transformers 🌔 PyTorch 😂 Safetensors 🚯 English llama facebook	meta llama-2 🕈 text-generation-inference 🗅 arxiv:2307.09288							
	Model card → Files and versions	: 《 Train · 《 Deploy · 《 Use in Transformers							
	Sated model You have been granted access to this model	Downloads last month 950,482							
	Llama 2								
	Llama 2 is a collection of pretrained and fine-tuned generative text models ranging in	Safetensors Model size 6.74B params Tensor type F16							
	scale from 7 billion to 70 billion parameters. This is the repository for the 7B fine-tuned	E. Test Connection							
	model, optimized for dialogue use cases and converted for the Hugging Face	V real Generation							
	Transformers format. Links to other models can be found in the index at the bottom.	Inference API has been turned off for this model.							
	Model Details	B Spaces using meta-llama/Llama-2-7b-chat-hf 483							
		IuggingFaceH4/open_llm_leaderboard ≤ qingxu98/gpt-academic							

Note: Use of this model is governed by the Meta license. In order to download the model

在申请完成之后,我们就可以下载Huggingface的llama2版本。

Read our paper, learn more about the model, or get started with code on Gitmub.

⊘ Llan	na Model In	dex		
Model	Llama2	Llama2-hf	Llama2-chat	Llama2-chat-hf
7B	<u>Link</u>	<u>Link</u>	Link	Link
13B	<u>Link</u>	Link	Link	Link
70B	<u>Link</u>	<u>Link</u>	<u>Link</u>	<u>Link</u>

llama2有许多版本,如上图,7B13B是参数规模,不确定计算资源能力的情况下,先选择7B版本, 实测3卡4090可以轻松运行7B版本。

后缀主要是功能和版本的不同,如本文档选择的llama2-7b-chat-hf便是7B参数规模的对话机器人, 自带使用Gradio构建的简易交互界面。

考虑到网络问题,我们需要使用的是Huggingface的镜像站和Hugg官方下载工具 `huggingfacecli

镜像站在: hf-mirror.com - Huggingface 镜像站

按照以下方式执行,可以得到下载在当前目录下的具体文件。

```
1 pip install -U huggingface_hub # 安装
2
3 export HF_ENDPOINT=https://hf-mirror.com #切换镜像源,每次使用huggingface-cli镜像
源下载都需要提前执行这一条命令
4
5 huggingface-cli download --resume-download meta-llama/Llama-2-7b-chat-hf --
local-dir Llama-2-7b-chat-hf --local-dir-use-symlinks False --token
hf_**Your token** # 替换自己的Huggingface token后,即可下载到本地,注意参数--local-
dir-use-symlinks 必须添加,不然大于5M的文件会被自动生成软链接,后续llama2无法启动
6
7
```

这样我们就下载完成了llama2-7b-chat-hf。

3. Docker运行

本文档运行环境为docker,优点是环境隔离足够强,且由于docker对端口需要映射,可以解决后续一 个很麻烦的端口问题。

(目前开源模型, llama2和stable diffusion等等, 开源仓库都使用7860端口作为交互端口, 如果服务器7860端口被占用, 就需要腾出端口或者修改工程文件来实现。但是修改工程文件需要花费更多时间寻找对应工程文件, docker 的端口映射可以帮助我们很简单地解决这个问题。)

要在docker环境下运行,自然需要安装docker,具体可以参考参考链接一节。自己配置docker file和 环境依赖很耗时间,本着不重复造轮子的原则,使用前辈已经配置好的docker 库文件来完成。

库文件地址: <u>soulteary/docker-llama2-chat: Play LLaMA2 (official / 中文版 / INT4 / llama2.cpp)</u> <u>Together! ONLY 3 STEPS! (non GPU / 5GB vRAM / 8~14GB vRAM) (github.com)</u>

该库已经有安装说明,可在参考链接中查阅。由于本文档定位为解决从0到1的问题,所以对该库文件的文件不作修改,力求最快速度跑通。

下载github 仓库,需要使用镜像站,可以使用kkgithub.com来下载。由于镜像站稳定性差,建议搜索需要下载时可以使用的镜像站,并参考对应命令来下载。也可以下载到本地后,使用vscode 配置 SFTP上传文件。

下面我们需要完成一系列调整工作:

1. 调整目录结构

```
    # 创建一个新的目录,用于存放我们的模型
    mkdir meta-llama
    # 将下载好的模型移动到目录中
    mv Llama-2-7b-chat-hf meta-llama/
    mv Llama-2-13b-chat-hf meta-llama/
    mv Llama-2-70b-chat-hf meta-llama/
```

2. 完整的目录结构类似下面这样,所有的模型都在我们创建的 meta-11 ama 目录的下一级中

```
1 # tree -L 2 meta-llama
2 meta-llama
3 ├── Llama-2-13b-chat-hf
4 │ ├── added_tokens.json
```

5	
6	
7	LICENSE.txt
8	
9	
10	
11	— model.safetensors.index.json
12	pytorch_model-00001-of-00003.bin
13	
14	
15	│
16	README.md
17	│
18	∣ ⊣ special_tokens_map.json
19	│
20	
21	USE_POLICY.md
22	└── Llama-2-7b-chat-hf
23	├── added_tokens.json
24	├── config.json
25	├── generation_config.json
26	LICENSE.txt
27	<pre>model-00001-of-00002.safetensors</pre>
28	├── model-00002-of-00002.safetensors
29	— model.safetensors.index.json
30	— modelsmeta-llamaLlama-2-7b-chat-ht
31	- pytorch_model-00001-of-00003.bin
32	- pytorch_model-00002-of-00003.bin
33	
34	<pre>pytorcn_model.bin.index.json</pre>
35	F README.MO
20 27	- special_cokens_map.json
20	- tokenizer_coning.json
20	tokenizer model
10	
40	
4 L	

× doaltar llamaΩ shat		з аоске
	• •	4
> docker		
> llama.cpp		
> llama2-7b		11000000111111111111111111111111111111
> llama2-7b-cn		
> llama2-7b-cn-4bit		SSH_CLI
> llama2-13b		LC_TIME
\sim meta-llama/Llama-2-7b-chat-hf		PATH=/u
♦ .gitattributes	U	DBUS_SE
🕄 config.json	U	LC_NUME
generation_config.json	U	_=/usr/ Removing of
LICENSE.txt	U	99d7768e.lo
≡ model-00001-of-00002.safetensors	U	Spawned rem
≡ model-00002-of-00002.safetensors	U	Waiting for
model.safetensors.index.json	U	Waiting for
≡ pytorch model-00001-of-00002.bin	U	Waiting for
≡ pytorch model-00002-of-00002.bin	U	Waiting for
Pytorch model.bin.index.json	U	Waiting for
① README.md	U	Waiting for
special tokens map.ison	U	Waiting for
B tokenizer config.ison	U	9d9aae13e3a
tokenizer ison	U	SSH_AUTH_SO
E tokenizer model	U	listeningOn
		osReleaseId
		arch==x86_6
		vscodeArch= bitness==64
		tmpDir==/ru
		platform==1
		didLocalDow

开始构建并运行:

1	# 进入程序目录
2	cd docker-llama2-chat
3	# 构建 7B 镜像
4	bash scripts/make-7b.sh
5	# 或者,构建 13B 镜像
6	bash scripts/make-13b.sh
7	
8	#构建完成后
9	
10	# 运行 7B 镜像,应用程序
11	bash scripts/run-7b.sh
12	# 或者,运行 13B 镜像,应用程序
13	bash scripts/run-13b.sh
14	

```
15
    #命令执行后,如果一切顺利,你将看到类似下面的日志:
16
17
    _____
18
    == PyTorch ==
19
    _____
20
21
    NVIDIA Release 23.06 (build 63009835)
22
    PvTorch Version 2.1.0a0+4136153
23
24
    Container image Copyright (c) 2023, NVIDIA CORPORATION & AFFILIATES. All
    rights reserved.
25
26
    Copyright (c) 2014-2023 Facebook Inc.
    Copyright (c) 2011-2014 Idiap Research Institute (Ronan Collobert)
27
28
    Copyright (c) 2012-2014 Deepmind Technologies
                                                    (Koray Kavukcuoglu)
    Copyright (c) 2011-2012 NEC Laboratories America (Koray Kavukcuoglu)
29
30
    Copyright (c) 2011-2013 NYU
                                                     (Clement Farabet)
    Copyright (c) 2006-2010 NEC Laboratories America (Ronan Collobert, Leon
31
    Bottou, Iain Melvin, Jason Weston)
32
    Copyright (c) 2006
                           Idiap Research Institute (Samy Bengio)
    Copyright (c) 2001-2004 Idiap Research Institute (Ronan Collobert, Samy
33
    Bengio, Johnny Mariethoz)
34
    Copyright (c) 2015
                           Google Inc.
35
   Copyright (c) 2015
                           Yangging Jia
36
    Copyright (c) 2013-2016 The Caffe contributors
37
    All rights reserved.
38
    Various files include modifications (c) NVIDIA CORPORATION & AFFILIATES. All
39
    rights reserved.
40
41
    This container image and its contents are governed by the NVIDIA Deep
    Learning Container License.
    By pulling and using the container, you accept the terms and conditions of
42
    this license:
43
    https://developer.nvidia.com/ngc/nvidia-deep-learning-container-license
44
45
    WARNING: CUDA Minor Version Compatibility mode ENABLED.
      Using driver version 525.105.17 which has support for CUDA 12.0. This
46
    container
47
      was built with CUDA 12.1 and will be run in Minor Version Compatibility
    mode.
48
      CUDA Forward Compatibility is preferred over Minor Version Compatibility
    for use
      with this container but was unavailable:
49
50
      [[Forward compatibility was attempted on non supported HW
    (CUDA_ERROR_COMPAT_NOT_SUPPORTED_ON_DEVICE) cuInit()=804]]
51
      See https://docs.nvidia.com/deploy/cuda-compatibility/ for details.
52
    Loading checkpoint shards: 100%|
53
                                                             2/2
    [00:05<00:00, 2.52s/it]
   Caching examples at: '/app/gradio_cached_examples/20'
54
55
    Caching example 1/5
```

```
56 /usr/local/lib/python3.10/dist-
    packages/transformers/generation/utils.py:1270: UserWarning: You have
    modified the pretrained model configuration to control generation. This is a
    deprecated strategy to control generation and will be removed soon, in a
    future version. Please use a generation configuration file (see
    https://huggingface.co/docs/transformers/main_classes/text_generation )
57
      warnings.warn(
    Caching example 2/5
58
    Caching example 3/5
59
60
    Caching example 4/5
61
    Caching example 5/5
62
    Caching complete
63
    /usr/local/lib/python3.10/dist-packages/gradio/utils.py:839: UserWarning:
64
    Expected 7 arguments for function <function generate at 0x7f3e096a1000>,
    received 6.
65
      warnings.warn(
    /usr/local/lib/python3.10/dist-packages/gradio/utils.py:843: UserWarning:
66
    Expected at least 7 arguments for function <function generate at
    0x7f3e096a1000>, received 6.
      warnings.warn(
67
    Running on local URL: http://0.0.0.0:7860
68
69
   To create a public link, set `share=True` in `launch()`.
70
71
72
```

然后使用浏览器打开https://localhost:7860即可运行。

如果是使用vscode远程开发, 会先出现下图的提示:

docker-llama2-chat > scripts > \$ run-7b.sh 1 3 4 ik=-1ulimit stack=67108864rm -it -v `pwd`/meta-llama:/app/meta-llama -p 7865:7860 soulteary/lla 4 问题 输出 调试控制台 终端 端口 3 • (base) fuhongye@amax:~\$ ls	
2 3 g 4 (base) fuhongye@amax:~\$ 1s	
3 [:k=-1ulimit stack=67108864rm -it -v `pwd`/meta-llama:/app/meta-llama -p 7865:7860 soulteary/lla 4 问题 输出 调试控制台 终端 端口 3 • (base) fuhongye@amax:~\$ ls	
问题 输出 调试控制台 终端 端口 3 ● (base) fuhongye@amax:~\$ 1s	na2:7b
问题 输出 调试控制台 终端 端口 3 ● (base) fuhongye@amax:~\$ 1s	
问题 输出 调试控制台 终端 端口 3 ● (base) fuhongye@amax:~\$ 1s	
●(base) fuhongye@amax:~\$ ls	$+ \cdot \cdot$
	∑ install
docker-llama2-chat llama2-l/D-nt	>_ install
<pre>(base) fuhongye@amax:~/docker-llama2-chat\$ bash scripts/run-7b.sh</pre>	() Dasii oc
== PyTorch ==	
NVIDIA Release 23.06 (build 63009835) PyTorch Version 2.1.0a0+4136153	
Container image Copyright (c) 2023, NVIDIA CORPORATION & AFFILIATES. All rights reserved.	
Copyright (c) 2014-2023 Facebook Inc. Copyright (c) 2011-2014 Idiap Research Institute (Ronan Collobert) Copyright (c) 2012-2014 Deepmind Technologies (Koray Kavukcuoglu) Copyright (c) 2011-2014 NEC Laboratories (Margica (Koray Kavukcuoglu)	

如果端口7860不可用,则可以修改为7865或者其他端口,如上图就修改为外部7865端口映射docker 运行的7860端口。

然后vscode会提示是否需要转发端口上的内容,可以添加7865端口进入转发设置,然后打开任意浏览器使用llama2。

i localhost:7860		A ^N Z	ን ር	Ġ		£≦	@ %3	
Llama-2 7B Chat								
This Space demonstrates model <u>Llama-2-7b-chat</u> by Meta, a Llama 2 mode run your own service, you can also <u>deploy the model on Inference Endpoint</u>	el with 7B parameters fine-tuned for chat instructions. Feel f <u>ts</u> .	free to play with it, or duplic	ate to run ş	eneratio	ns withou	ut a queue	l If you wa	nt to
For more details about the Llama 2 family of models and how to use the	em with transformers, take a look <u>at our blog post</u> .							
K Looking for an even more powerful model? Check out the <u>13B version</u> or	r the large <u>70B model demo</u> .							
Duplicate Space for private use								
🖙 Chatbot								
 In control to the state of the control to the control to the table of the control to th								
	Zhejiang University							
Ah, my apologies for the confusion earlier! Zhejiang Univer some key facts and figures about the university:	ersity (ZJU) is a prestigious public research universit	y located in Hangzhou,	Zhejiang	Province	, China	. Here ar	e	
 History: ZJU was founded in 1897 as a veterinary school engineering natural sciences humanities and social 	ool, and it has since evolved into a comprehensive sciences	research university with	a wide ra	nge of d	iscipline	es, includ	ding	
Type a message							提交	
Retry	D Undo				Clear			

4. 参考链接

- 1. <u>最详细的ubuntu 安装 docker教程 知乎 (zhihu.com)</u>
- 2. 使用 Docker 快速上手官方版 LLaMA2 开源大模型 soulteary的博客-CSDN博客
- 3. 【nvidia-smi报错】Failed to initialize NVML: Driver/library version mismatch-CSDN博客
- 4. [docker gpu报错Error response from daemon: could not select device driver "" with capabilities: [[gpu]]_Adenialzz的博客-CSDN博客](<u>https://blog.csdn.net/weixin 44966641/art icle/details/123760614</u>)
- 5. <u>hf-mirror.com Huggingface 镜像站</u>

written by Hongye Fu 2023.11.24, 感谢soulteary前辈的文档和仓库