# [玩转AIGC]如何训练LLaMA2（模型训练、推理、代码讲解，并附可直接运行的kaggle连接）

原创　六五酥　☑ 已于 2024-02-26 09:53:14 修改　◉ 阅读量1.4w　★ 收藏　40　👍 点赞数 11　　　　　　版权

分类专栏　　于 2023-07-31 14:42:51 首次发布　GC　　python　　机器学习　　人工智能

G　**GitCode 开源社区** 文章已被社区收录

　　**玩转AIGC** 专栏收录该内容

**目录**

Llama 2，基于优化的 Transformer 架构，是Meta AI正式发布的最新一代开源大模型，一系列模型（7b、13b、70b）均开源可商用，效果直逼gpt3.5。

下面我们来介绍如何使用Llama 2来训练一个故事 生成模型　。

如果迫不及待想爽一把先，请直接跳到这里，可直接运行：llama2-c，学习不就是先让自己爽起来，而后才有欲望去深究为什么！！！

# 一、clone仓库

首先从 github　将仓库拉到本地：

```
1   git clone https://github.com/karpathy/llama2.c
```

进入目录

```
1   cd llama2.c
```

## 二、数据集下载与处理

> 注：具体内容可见[玩转AIGC]sentencepiece训练一个Tokenizer(标记器)这篇文章

### 1、数据集下载

```
1   python tinystories.py download
```

### 2、数据集tokenize预处理（耗时较长）

```
1   python tinystories.py pretokenize
```

若运行到中途卡死了，可将并行运行的线程改小，打开tinystories.py，找到def pretolenize()函数，修改线程数max_workers=4：

```
1       # process all the shards in a threadpool
2       with ThreadPoolExecutor(max_workers=4) as executor:
3           executor.map(process_shard, shard_filenames)
```

注：在运行文末的kaggle代码时，不需要运行这一步，我已经提前编好的，只要运行 !cp
/kaggle/input/llama2tranningdatabin/databin/* /kaggle/working/llama2.c/data/TinyStories_all_data 就可把文件拷贝到对应
目录下，因为编号太久了，而且在kaggle里面容易挂掉

## 三、修改配置

然后修改一些配置：

`train.py` 里面有几个参数要修改
batch_size改小一点，否则会报'CUDA out of memory' 的错误（土豪卡多随意，不介意的话送我一张）
dtype要改为"float16"，否则会报'Current CUDA Device does not support bfloat16'的错误
compile要改为False，否则会报CUDA Capability过低或complex64不支持的错误

```
1   batch_size = 64
2   dtype = "float16"
3   compile = False
```

可选改的参数：
max_iters：是迭代次数，可改小一点，在kaggle里面运行一代是1.2s，运行100000代的话，大概还得34小时，而kaggle免费时长

最长也才30小时，且一次运行最长12小时，所以要改小（还是那句话，土豪有卡的随意）

warmup_iters：是热身的迭代次数，主要是为了确定合适得学习率，卡有限的话可改小一些。

```
1   max_iters = 100000
2   warmup_iters = 1000
```

注：在运行文末的kaggle代码时，这里不用自己修改了，直接把我修改好的复制过来即可，运行 !cp /kaggle/input/llama2trainpy4/train.py ./train.py 即可复制，因为kaggle里面也手动改不了，能改也是很麻烦。

## 四、开始训练

```
1   python train.py
```

```
False
tokens per iteration will be: 65,536
breaks down as: 4 grad accum steps * 1 processes * 64 batch size * 256 max seq len
Initializing a new model from scratch
num decayed parameter tensors: 43, with 15,187,968 parameters
num non-decayed parameter tensors: 13, with 3,744 parameters
using fused AdamW: True
Created a PretokDataset with rng seed 42
Created a PretokDataset with rng seed 42
Created a PretokDataset with rng seed 42
step 0: train loss 10.4168, val loss 10.4172
0 | loss 10.4041 | lr 0.000000e+00 | 23329.64ms | mfu -100.00%
1 | loss 10.4179 | lr 5.000000e-07 | 1274.02ms | mfu -100.00%
2 | loss 10.4191 | lr 1.000000e-06 | 1274.26ms | mfu -100.00%
3 | loss 10.4062 | lr 1.500000e-06 | 1274.07ms | mfu -100.00%
4 | loss 10.4040 | lr 2.000000e-06 | 1272.92ms | mfu -100.00%
5 | loss 10.3932 | lr 2.500000e-06 | 1274.56ms | mfu 1.59%
6 | loss 10.3963 | lr 3.000000e-06 | 1274.45ms | mfu 1.59%
7 | loss 10.3781 | lr 3.500000e-06 | 1273.72ms | mfu 1.59%
8 | loss 10.3614 | lr 4.000000e-06 | 1274.24ms | mfu 1.59%
9 | loss 10.3502 | lr 4.500000e-06 | 1276.00ms | mfu 1.59%
10 | loss 10.3284 | lr 5.000000e-06 | 1274.97ms | mfu 1.59%
11 | loss 10.3042 | lr 5.500000e-06 | 1274.46ms | mfu 1.59%
12 | loss 10.3015 | lr 6.000000e-06 | 1274.41ms | mfu 1.59%
13 | loss 10.2579 | lr 6.500000e-06 | 1274.68ms | mfu 1.59%
14 | loss 10.2374 | lr 7.000000e-06 | 1274.30ms | mfu 1.59%
15 | loss 10.2206 | lr 7.500000e-06 | 1273.54ms | mfu 1.59%
```

CSDN @六五酥

训练后可在out文件夹下看到以下文件：

**Output (10.1GB / 19.5GB)**

▼ 📁 /kaggle/working
  ▼ 📁 llama2.c
    ▼ 📁 out
        📄 ckpt.pt
        📄 model.bin

CSDN @六五酥

当然还可以在命令行中指定训练参数(单GPU上训练)：

```
1   python -m train.py --compile=False --eval_iters=10 --batch_size=8
```

如果是多GPU，可采用分布式训练，例如采用DDP 在1个node，4个 gpu 上训练：

```
1   torchrun --standalone --nproc_per_node=4 train.py
```

采用DDP 在2个node，4个 gpu 上训练：

```
1   torchrun --nproc_per_node=8 --nnodes=2 --node_rank=0 --master_addr=123.456.123.456 --master_port=1234 train.
```

注：如果觉得训练太久了，可以跑我训练出来的模型，拷贝过来就可以了：

```
1   !cp /kaggle/input/llama2-out-model/* ./out
```

如果你没运行过train.py，那就是不会有llama2/out这个目录，这时候要先创建好out这个目录：

```
1   mkdir out
```

## 五、模型推理

先编译run（只需要运行一次）

```
1   make run
```

运行模型推理

```
1   ./run out/model.bin
```

输出如下：

> Once upon a time there was a little bear named manner. He was very playful, and he loved to jump around. One day, he saw a big pillow in his den. It was bright blue, just like the sky, and it seemed just right for him. He decided he would try to jump on it.
> He pASHed into the air, feeling the pillow with his fluffy paws and carried on the giant. He bounced up and down, dodgingages and twists. He felt so free and happy as he jumped higher and higher.
> But then he started to struggle. The pillow was too hard and he couldn't jump on it! He tried and tried and even though he was tired and frustrated, he couldn't take a kick. Finally, he gave up. But he still felt free, like nothing could stop him.
> Close was tired, panting and relieved at its same time. He jumped and flew away, still feeling happy after exploring the sky. He curled up on the pillow to take his nap and finally fell asleep, dreaming of himself jumping across the sky.

感兴趣的话，也可以跑其它模型试试：

```
1   wget https://huggingface.co/karpathy/tinyllamas/resolve/main/stories15M.bin
2   ./run stories15M.bin
3
4   wget https://huggingface.co/karpathy/tinyllamas/resolve/main/stories42M.bin
5   ./run stories42M.bin
```

当然还可指定参数运行：

```
1    ./run stories42M.bin 1.0 256 "One day, Lily met a Shoggoth"
```

注：在本小节末尾的kaggle代码时，可能会报找不到triton模块，此时运行 `!pip3 install triton`，进行安装即可。

以上可在kaggle运行：llama2-c

# 六、train.py训练代码讲解

注：为了方便讲解，已经把部分ddp（分布式训练，多gpu可用，我的kaggle是单gpu，所以去掉了）跟wandb_log（记录日志）去掉了，整体代码在**train.py**里面，此处不再赘述了。

## 1、导包

```python
1    import math
2    import os
3    import time
4    from contextlib import nullcontext
5    from datetime import datetime
6    from functools import partial
7
8    import torch
9    from model import Transformer, ModelArgs # 在model.py里的模型
10   from torch.distributed import destroy_process_group, init_process_group
11   from torch.nn.parallel import DistributedDataParallel as DDP
12   #import torch._dynamo
13
14   from tinystories import Task
```

## 2、定义模型训练参数与相关设置

```python
1
2    # -----------------------------------------------------------------------------
3    # I/O
     out_dir = "out"
4    eval_interval = 2000
5    log_interval = 1
6    eval_iters = 100
7    eval_only = False  # if True, script exits right after the first eval
8    always_save_checkpoint = False # if True, always save a checkpoint after each eval
9    init_from = "scratch"  # 'scratch' or 'resume'
10   # wandb logging
11   wandb_log = False  # disabled by default
12   wandb_project = "llamac"
13   wandb_run_name = "run" + datetime.now().strftime("%Y_%m_%d_%H_%M_%S")
14   # data
15   batch_size = 64  # if gradient_accumulation_steps > 1, this is the micro-batch size
16   max_seq_len = 256
17   # model
18   dim = 288
19   n_layers = 6
20   n_heads = 6
21   multiple_of = 32
```

```
22   dropout = 0.0
23   # adamw optimizer
24   gradient_accumulation_steps = 4   # used to simulate larger batch sizes
25   learning_rate = 5e-4   # max learning rate
26   max_iters = 5000   # total number of training iterations
27   weight_decay = 1e-1
28   beta1 = 0.9
29   beta2 = 0.95
30   grad_clip = 1.0   # clip gradients at this value, or disable if == 0.0
31   # learning rate decay settings
32   decay_lr = True   # whether to decay the learning rate
33   warmup_iters = 100   # how many steps to warm up for
34   # system
35   device = "cuda"   # examples: 'cpu', 'cuda', 'cuda:0', 'cuda:1' etc., or try 'mps' on macbooks
36   dtype = "float16"   # float32|bfloat16|float16
37   compile = False # use PyTorch 2.0 to compile the model to be faster
38   # --------------------------------------------------------------------------
39
```

下面是各个参数的解释：

```
1
2    I/O
3
4        out_dir: 模型训练输出路径。
5        eval_interval: 多少个训练步骤后进行一次模型评估。
6        log_interval: 多少个训练步骤后进行一次日志记录。
7        eval_iters: 在进行模型评估时，评估器将处理多少个数据集条目。
8        eval_only: 如果为 True，则仅进行一次模型评估并退出脚本。
9        always_save_checkpoint: 如果为 True，则在每次模型评估后始终保存一个检查点。
10       init_from: 模型初始化方法，可以是 'scratch'（从头开始训练）或 'resume'（从之前的检查点恢复）。
11
12   wandb 日志记录
13
14       wandb_log: 是否启用 wandb 日志记录。
15       wandb_project: wandb 项目名称。
16       wandb_run_name: wandb 运行名称。
17
18   数据
19
20       batch_size: 训练的批次大小。
21       max_seq_len: 输入序列的最大长度。
22
23   模型
24
25       dim: Transformer 模型中的隐藏层维度。
26       n_layers: Transformer 模型中的层数。
27       n_heads: Transformer 模型中的多头注意力头数。
28       multiple_of: 批次大小必须是此值的倍数。在 TPU 上训练时，批次大小必须是 8 的倍数。
29       dropout: 模型中的 dropout 概率。
30
31   adamw 优化器
32
33       gradient_accumulation_steps: 用于模拟较大批次大小的梯度累积步骤数。
34       learning_rate: 最大学习率。
```

```
35        max_iters: 总训练迭代次数。
36        weight_decay: AdamW 优化器中的权重衰减系数。
37        beta1: AdamW 优化器的 beta1 超参数。
38        beta2: AdamW 优化器的 beta2 超参数。
39        grad_clip: 梯度裁剪值，如果设为 0.0 则不进行梯度裁剪。
40
41   学习率衰减设置
42
43        decay_lr: 是否对学习率进行衰减。
44        warmup_iters: 学习率预热步骤数。
45
46   系统
47
48        device: 训练设备，可以是 'cpu'、'cuda'、'cuda:0'、'cuda:1' 等，或在 MacBook 上尝试 'mps'.
49        dtype: 训练数据类型，可以是 'float32'、'bfloat16' 或 'float16'.
50        compile: 是否使用 PyTorch 2.0 编译模型以提高速度。
```

## 3、加载模型配置

```python
1
2    # -----------------------------------------------------------------------------
3    config_keys = [
4        k
5        for k, v in globals().items()
6        if not k.startswith("_") and isinstance(v, (int, float, bool, str))
7    ]
8    exec(open("configurator.py").read())  # overrides from command line or config file
9    config = {k: globals()[k] for k in config_keys}  # will be useful for logging
10   # -----------------------------------------------------------------------------
11
12   # fixing some hyperparams to sensible defaults
13   lr_decay_iters = max_iters  # should be ~= max_iters per Chinchilla
     min_lr = 0.0  # minimum learning rate, should be ~= learning_rate/10 per Chinchilla
14
15   master_process = True
16   seed_offset = 0
17   ddp_world_size = 1
18
19   tokens_per_iter = gradient_accumulation_steps * ddp_world_size * batch_size * max_seq_len
20   if master_process:
21       print(f"tokens per iteration will be: {tokens_per_iter:,}")
22       print(f"breaks down as: {gradient_accumulation_steps} grad accum steps * {ddp_world_size} processes * {b
23
24   if master_process:
25       os.makedirs(out_dir, exist_ok=True)
26   torch.manual_seed(1337 + seed_offset)
27   torch.backends.cuda.matmul.allow_tf32 = True  # allow tf32 on matmul
28   torch.backends.cudnn.allow_tf32 = True  # allow tf32 on cudnn
29   device_type = "cuda" if "cuda" in device else "cpu"  # for later use in torch.autocast
30   # note: float16 data type will automatically use a GradScaler
31   ptdtype = {"float32": torch.float32, "bfloat16": torch.bfloat16, "float16": torch.float16}[dtype]
32   ctx = (
33       nullcontext()
34       if device_type == "cpu"
```

```
35        else torch.amp.autocast(device_type=device_type, dtype=ptdtype)
36    )
```

## 4、迭代生成数据

用于将输入数据集划分为若干个大小相同的批次，并对每个批次进行预处理和编码，以便于送入模型进行训练。在模型训练过程中，需要不断地生成训练数据批次，然后将这些批次送入模型进行训练。因此，这个新函数 iter_batches 就是用于迭代生成数据批次的迭代器函数，它会在每个迭代步骤中调用 Task.iter_batches() 函数生成一个数据批次，并返回给调用者。这样，就可以通过简单的 for 循环来不断地生成数据批次，然后将这些批次送入模型进行训练。

```
1    # task-specific setup
2    iter_batches = partial(
3        Task.iter_batches,
4        batch_size=batch_size,
5        max_seq_len=max_seq_len,
6        device=device,
7        num_workers=0,
8    )
```

## 5、模型初始化

iter_num: 当前迭代次数

best_val_loss:最佳验证集损失值

以上两个值用于记录模型训练的状态。这些变量的初始值为 0 和一个较大的数值，表示模型训练尚未开始，最佳验证集损失值尚未确定。

```
1    # init these up here, can override if init_from='resume' (i.e. from a checkpoint)
2    iter_num = 0
3    best_val_loss = 1e9
```

使用 ModelArgs 类来设置模型参数 model_args

```
1    # model init
2    model_args = dict(
3        dim=dim,
4        n_layers=n_layers,
5        n_heads=n_heads,
6        n_kv_heads=n_heads,
7        vocab_size=32000,
8        multiple_of=multiple_of,
9        max_seq_len=max_seq_len,
10       dropout=dropout,
11   )  # start with model_args from command line
```

然后，根据 init_from 参数的取值来初始化模型。如果 init_from 的值为 "scratch"，则表示从头开始训练一个新模型，如果 init_from 的值为 "resume"，则表示从之前的训练中恢复训练。

先来看：init_from == "scratch"

直接根据model_args创建参数类，然后通过 Transformer 类来创建一个新的 Transformer 模型 model。

```
1   if init_from == "scratch":
2       # init a new model from scratch
3       print("Initializing a new model from scratch")
4       gptconf = ModelArgs(**model_args)
5       model = Transformer(gptconf)
```

再来看：init_from == "resume"

1）先从checkpoint加载模型（"ckpt.pt"）

2）然后根据model_args创建参数类，然后通过 Transformer 类来创建一个新的 Transformer 模型 model。

3）最后修复模型状态字典中的键名，因为有时候在保存模型的时候会在键名前面加上不必要的前缀，导致在加载模型时出现问题。坦白说我们也不清楚这个前缀是怎么产生的，需要进一步进行调试和研究。

```
1    elif init_from == "resume":
2        print(f"Resuming training from {out_dir}")
3        # resume training from a checkpoint.
4        ckpt_path = os.path.join(out_dir, "ckpt.pt")
5        checkpoint = torch.load(ckpt_path, map_location=device)
6        checkpoint_model_args = checkpoint["model_args"]
7        # force these config attributes to be equal otherwise we can't even resume training
8        # the rest of the attributes (e.g. dropout) can stay as desired from command line
9        for k in ["dim", "n_layers", "n_heads", "n_kv_heads", "vocab_size", "multiple_of", "max_seq_len"]:
10           model_args[k] = checkpoint_model_args[k]
11       # create the model
12       gptconf = ModelArgs(**model_args)
13       model = Transformer(gptconf)
14       state_dict = checkpoint["model"]
15       # fix the keys of the state dictionary :(
16       # honestly no idea how checkpoints sometimes get this prefix, have to debug more
17       unwanted_prefix = "_orig_mod."
18       for k, v in list(state_dict.items()):
19           if k.startswith(unwanted_prefix):
20               state_dict[k[len(unwanted_prefix) :]] = state_dict.pop(k)
21       model.load_state_dict(state_dict)
22       iter_num = checkpoint["iter_num"]
23       best_val_loss = checkpoint["best_val_loss"]
```

无论是从头开始训练新模型还是从之前的训练中恢复模型，最后都将模型 model 移动到指定的设备 device 上进行训练。

```
1   model.to(device)
```

## 6、设置自动混合精度与优化函数

scaler：设置自动混合精度

optimizer：优化函数

如果是恢复训练的，并且checkpoint里面存在optimizer，那么optimizer会从checkpoint里面加载

```
1
2   # initialize a GradScaler. If enabled=False scaler is a no-op
3   scaler = torch.cuda.amp.GradScaler(enabled=(dtype == "float16"))
4
5   # optimizer
```

```
6    optimizer = model.configure_optimizers(weight_decay, learning_rate, (beta1, beta2), device_type)
7    if init_from == "resume" and "optimizer" in checkpoint:
8        optimizer.load_state_dict(checkpoint["optimizer"])
     checkpoint = None  # free up memory
```

## 7、损失评估与学习率获取

**estimate_loss()** ：

使用多个数据批次来计算一个任意精度的损失值，用于评估模型在指定数据集上的性能。该函数会将模型置为评估模式（调用 model.eval() 函数），然后对训练集和验证集分别进行评估。对于每个数据集，该函数会使用 iter_batches() 函数生成一个数据批次迭代器 batch_iter，然后对迭代器中的每个数据批次进行前向传播和损失计算，得到一个损失值。这里使用 eval_iters 次评估来对整个数据集进行评估，从而得到一个更准确的损失值。

最后，该函数返回一个字典 out，其中包含了训练集和验证集的平均损失值。需要注意的是，评估完成后，该函数将模型重新置为训练模式（调用 model.train() 函数），以便于之后的模型训练。

**get_lr(it)** ：

这段代码定义了一个函数 get_lr()，用于根据当前迭代次数来动态计算学习率。具体来说，该函数使用了一种余弦退火（cosine annealing）的学习率调度策略，即先进行一定步数的学习率线性预热（linear warmup），然后使用余弦函数进行学习率退火，直到学习率降到最小值为止。

该函数的输入参数是当前迭代次数 it，输出参数是当前迭代次数下的学习率。具体来说，该函数首先判断当前迭代次数是否小于预热步数 warmup_iters，如果是，则按照线性预热的方式逐步增加学习率，直到达到 learning_rate。如果当前迭代次数已经超过了学习率退火的步数 lr_decay_iters，则直接返回最小学习率 min_lr。如果当前迭代次数在预热步数和退火步数之间，则使用余弦函数计算当前迭代次数下的学习率，将学习率从初始值 learning_rate 逐渐下降到最小学习率 min_lr。

需要注意的是，在函数中使用了一个系数 coeff，它的值在 0 和 1 之间，用于调整学习率的下降速度。这个系数是通过余弦函数计算得到的，随着迭代次数的增加，系数的值逐渐从 1 下降到 0，从而实现学习率的下降。

```
1
2    # helps estimate an arbitrarily accurate loss over either split using many batches
3    @torch.no_grad()
4    def estimate_loss():
5        out = {}
6        model.eval()
7        for split in ["train", "val"]:
8            batch_iter = iter_batches(split)
9            losses = torch.zeros(eval_iters)  # keep on CPU
10           for k in range(eval_iters):
11               X, Y = next(batch_iter)
12               with ctx:
13                   logits, loss = model(X, Y)
14               losses[k] = loss.item()
15           out[split] = losses.mean()
16       model.train()
17       return out
18
19   # learning rate decay scheduler (cosine with warmup)
20   def get_lr(it):
21       # 1) linear warmup for warmup_iters steps
22       if it < warmup_iters:
23           return learning_rate * it / warmup_iters
24       # 2) if it > lr_decay_iters, return min learning rate
25       if it > lr_decay_iters:
26           return min_lr
```

```
26      # 3) in between, use cosine decay down to min learning rate
27      decay_ratio = (it - warmup_iters) / (lr_decay_iters - warmup_iters)
28      assert 0 <= decay_ratio <= 1
29      coeff = 0.5 * (1.0 + math.cos(math.pi * decay_ratio))  # coeff ranges 0..1
30      return min_lr + coeff * (learning_rate - min_lr)
```

## 8、日志保存初始化

没什么可说的，就是采用wandb，并进行初始化

```
1   # logging
2   if wandb_log and master_process:
3       import wandb
4       wandb.init(project=wandb_project, name=wandb_run_name, config=config)
```

## 9、循环训练

这部分主要实现了模型的训练循环。具体来说，该循环会在训练集上进行多次迭代，每次迭代会使用一个数据批次进行前向传播、反向传播和参数更新。在每次迭代过程中，该循环会根据当前迭代次数来动态调整学习率，并使用学习率和损失函数对模型参数进行更新。此外，该循环还会定期对训练集和验证集进行损失评估，并保存模型checkpoint。

开始时，该循环会使用 iter_batches() 函数生成一个训练集数据批次迭代器 train_batch_iter，然后从迭代器中获取第一个数据批次 X, Y。接着，该循环会不断迭代，直到达到指定的最大迭代次数 max_iters 为止。

在**while True**的每次迭代中：该循环会先根据当前迭代次数调整学习率，然后使用优化器对模型参数进行更新。如果当前迭代次数是评估间隔 eval_interval 的倍数，并且当前进程是主进程，那么该循环会对训练集和验证集进行损失评估，并将评估结果保存到 losses 中。如果当前验证集的损失值比之前的最佳损失值要小，或者总是保存检查点的标志 always_save_checkpoint 被设置为真，那么该循环会将模型参数和优化器状态保存到检查点文件中。

接着，该循环会使用当前数据批次 X, Y 进行前向传播和反向传播，并根据梯度累积步数 gradient_accumulation_steps 进行梯度平均。如果在训练过程中使用了混合精度训练，那么该循环还会使用 scaler 对梯度进行缩放。然后，该循环会使用裁剪梯度的方法对梯度进行限幅，并使用优化器对模型参数进行更新。最后，该循环会将梯度清零，以便于下一次迭代。

在训练过程中，该循环会定期打印当前迭代次数、损失值、学习率、迭代时间和模型的内存使用率等信息。此外，该循环会根据指定的最大迭代次数 max_iters 判断是否终止训练。

```
2
3   # training loop
    train_batch_iter = iter_batches("train")
4
    X, Y = next(train_batch_iter)  # fetch the very first batch
5
    t0 = time.time()
6
    local_iter_num = 0  # number of iterations in the lifetime of this process
7
    raw_model = model.module if ddp else model  # unwrap DDP container if needed
8
    running_mfu = -1.0
9
    while True:
10
        # determine and set the learning rate for this iteration
11
        lr = get_lr(iter_num) if decay_lr else learning_rate
12
        for param_group in optimizer.param_groups:
13
            param_group["lr"] = lr
14
15
        # evaluate the loss on train/val sets and write checkpoints
16
        if iter_num % eval_interval == 0 and master_process:
17
            losses = estimate_loss()
```

```
18              print(f"step {iter_num}: train loss {losses['train']:.4f}, val loss {losses['val']:.4f}")
19          if losses["val"] < best_val_loss or always_save_checkpoint:
20              best_val_loss = losses["val"]
21              if iter_num > 0:
22                  checkpoint = {
23                      "model": raw_model.state_dict(),
24                      "optimizer": optimizer.state_dict(),
25                      "model_args": model_args,
26                      "iter_num": iter_num,
27                      "best_val_loss": best_val_loss,
28                      "config": config,
29                  }
30                  print(f"saving checkpoint to {out_dir}")
31                  torch.save(checkpoint, os.path.join(out_dir, "ckpt.pt"))
32                  raw_model.export(os.path.join(out_dir, "model.bin"))
33      if iter_num == 0 and eval_only:
34          break
35
36      # forward backward update, with optional gradient accumulation to simulate larger batch size
37      # and using the GradScaler if data type is float16
38      for micro_step in range(gradient_accumulation_steps):
39          with ctx:
40              logits, loss = model(X, Y)
41              loss = loss / gradient_accumulation_steps
42          # immediately async prefetch next batch while model is doing the forward pass on the GPU
43          X, Y = next(train_batch_iter)
44          # backward pass, with gradient scaling if training in fp16
45          scaler.scale(loss).backward()
46      # clip the gradient
47      if grad_clip != 0.0:
48          scaler.unscale_(optimizer)
49          torch.nn.utils.clip_grad_norm_(model.parameters(), grad_clip)
50      # step the optimizer and scaler if training in fp16
51      scaler.step(optimizer)
52      scaler.update()
53      # flush the gradients as soon as we can, no need for this memory anymore
54      optimizer.zero_grad(set_to_none=True)
55
56      # timing and logging
57      t1 = time.time()
58      dt = t1 - t0
59      t0 = t1
60      if iter_num % log_interval == 0 and master_process:
61          # get loss as float, scale up due to the divide above. note: this is a CPU-GPU sync point
62          lossf = loss.item() * gradient_accumulation_steps
63          if local_iter_num >= 5:  # let the training loop settle a bit
64              mfu = raw_model.estimate_mfu(batch_size * gradient_accumulation_steps, dt)
65              running_mfu = mfu if running_mfu == -1.0 else 0.9 * running_mfu + 0.1 * mfu
66          print(
67              f"{iter_num} | loss {lossf:.4f} | lr {lr:e} | {dt*1000:.2f}ms | mfu {running_mfu*100:.2f}%"
68          )
69      iter_num += 1
70      local_iter_num += 1
71
72      # termination conditions
```

```
73        if iter_num > max_iters:
              break
```

## 七、run.c推理代码讲解

注：整体代码在**run.c**里面。

## 1、结构及内存管理

定义了一个基于Transformer架构的神经网络模型的配置参数、运行状态和内存管理。
**具体来说**：

- Config结构体定义了Transformer模型的一些超参数，例如Transformer的维度、隐藏层维度、层数、查询头数、键/值头数、词汇表大小、最大序列长度等等。
- TransformerWeights结构体定义了Transformer模型的所有权重矩阵，包括Token嵌入表、RMSNorm的权重、矩阵乘积的权重、前馈网络的权重、RoPE相对位置编码的频率矩阵以及用于分类的权重（可选）。
- RunState结构体定义了Transformer模型每个时间步的状态，例如当前激活值、残差分支内的激活值、FFN中的隐藏层激活值、查询、键、值、注意力得分、输出概率分布等等，同时还包括用于缓存键/值的缓存矩阵。函数malloc_run_state用于为RunState结构体中的各个数组分配内存空间。
- 函数free_run_state用于释放RunState结构体中各个数组占用的内存空间。

在代码中，这些数组的内存分配和释放都使用了calloc和free函数。calloc函数与malloc函数类似，但它会在分配内存后将其初始化为零，避免出现未初始化的内存访问问题。在分配内存后，代码还会检查是否分配成功以避免出现内存分配失败的情况。

```
2
3    // ----------------------------------------------------------------------------
4    // Transformer and RunState structs, and related memory management
5    typedef struct {
6        int dim; // transformer dimension
7        int hidden_dim; // for ffn layers
8        int n_layers; // number of layers
9        int n_heads; // number of query heads
10       int n_kv_heads; // number of key/value heads (can be < query heads because of multiquery)
11       int vocab_size; // vocabulary size, usually 256 (byte-level)
12       int seq_len; // max sequence length
13   } Config;
14
15   typedef struct {
16       // token embedding table
17       float* token_embedding_table;    // (vocab_size, dim)
18       // weights for rmsnorms
19       float* rms_att_weight; // (layer, dim) rmsnorm weights
20       float* rms_ffn_weight; // (layer, dim)
21       // weights for matmuls
22       float* wq; // (layer, dim, dim)
23       float* wk; // (layer, dim, dim)
24       float* wv; // (layer, dim, dim)
25       float* wo; // (layer, dim, dim)
26       // weights for ffn
27
```

```c
28        float* w1; // (layer, hidden_dim, dim)
29        float* w2; // (layer, dim, hidden_dim)
30        float* w3; // (layer, hidden_dim, dim)
31        // final rmsnorm
32        float* rms_final_weight; // (dim,)
33        // freq_cis for RoPE relatively positional embeddings
34        float* freq_cis_real; // (seq_len, dim/2)
35        float* freq_cis_imag; // (seq_len, dim/2)
36        // (optional) classifier weights for the logits, on the last layer
37        float* wcls;
38    } TransformerWeights;
39
40    typedef struct {
41        // current wave of activations
42        float *x; // activation at current time stamp (dim,)
43        float *xb; // same, but inside a residual branch (dim,)
44        float *xb2; // an additional buffer just for convenience (dim,)
45        float *hb; // buffer for hidden dimension in the ffn (hidden_dim,)
46        float *hb2; // buffer for hidden dimension in the ffn (hidden_dim,)
47        float *q; // query (dim,)
48        float *k; // key (dim,)
49        float *v; // value (dim,)
50        float *att; // buffer for scores/attention values (n_heads, seq_len)
51        float *logits; // output logits
52        // kv cache
53        float* key_cache;   // (layer, seq_len, dim)
54        float* value_cache; // (layer, seq_len, dim)
55    } RunState;
56
57    void malloc_run_state(RunState* s, Config* p) {
58        // we calloc instead of malloc to keep valgrind happy
59        s->x = calloc(p->dim, sizeof(float));
60        s->xb = calloc(p->dim, sizeof(float));
61        s->xb2 = calloc(p->dim, sizeof(float));
62        s->hb = calloc(p->hidden_dim, sizeof(float));
63        s->hb2 = calloc(p->hidden_dim, sizeof(float));
64        s->q = calloc(p->dim, sizeof(float));
65        s->k = calloc(p->dim, sizeof(float));
66        s->v = calloc(p->dim, sizeof(float));
67        s->att = calloc(p->n_heads * p->seq_len, sizeof(float));
68        s->logits = calloc(p->vocab_size, sizeof(float));
69        s->key_cache = calloc(p->n_layers * p->seq_len * p->dim, sizeof(float));
70        s->value_cache = calloc(p->n_layers * p->seq_len * p->dim, sizeof(float));
71        // ensure all mallocs went fine
72        if (!s->x || !s->xb || !s->xb2 || !s->hb || !s->hb2 || !s->q
73         || !s->k || !s->v || !s->att || !s->logits || !s->key_cache
74         || !s->value_cache) {
75            printf("malloc failed!\n");
76            exit(1);
77        }
78    }
79
80    void free_run_state(RunState* s) {
81        free(s->x);
82        free(s->xb);
```

```
83        free(s->xb2);
84        free(s->hb);
85        free(s->hb2);
86        free(s->q);
87        free(s->k);
88        free(s->v);
89        free(s->att);
90        free(s->logits);
91        free(s->key_cache);
          free(s->value_cache);
    }
```

## 2、模型初始化：读取checkpoint

具体来说，该函数的输入包括：

- 一个指向TransformerWeights结构体的指针w，用于存储从文件中读取的权重；
- 一个指向Config结构体的指针p，包含Transformer模型的超参数；
- 一个指向float类型的指针f，指向从文件中读取的所有模型权重的连续内存块；
- 一个布尔值shared_weights，指示是否共享Token嵌入表和分类器权重。

该函数的实现方式是，首先将指向连续内存块的指针f按照顺序分别指向各个数组的起始位置，然后使用指针算术运算逐个填充这些数组。填充顺序与数组在TransformerWeights结构体中的定义顺序相同。

需要注意的是，最后一个权重数组wcls的起始位置可能与前面的数组不同，这取决于shared_weights的值。如果shared_weights为true，则Token嵌入表和分类器权重共享相同的内存空间，此时wcls指向Token嵌入表的起始位置；否则，wcls指向连续内存块的当前位置。

该函数的作用是将从文件中读取的权重初始化到内存中，以便后续的Transformer模型推理过程中使用。

```
1
2   // ---------------------------------------------------------------------------
3   // initialization: read from checkpoint
4
5   void checkpoint_init_weights(TransformerWeights *w, Config* p, float* f, int shared_weights) {
6       float* ptr = f;
7       w->token_embedding_table = ptr;
8       ptr += p->vocab_size * p->dim;
9       w->rms_att_weight = ptr;
10      ptr += p->n_layers * p->dim;
11      w->wq = ptr;
12      ptr += p->n_layers * p->dim * p->dim;
13      w->wk = ptr;
14      ptr += p->n_layers * p->dim * p->dim;
15      w->wv = ptr;
16      ptr += p->n_layers * p->dim * p->dim;
17      w->wo = ptr;
18      ptr += p->n_layers * p->dim * p->dim;
19      w->rms_ffn_weight = ptr;
20      ptr += p->n_layers * p->dim;
21      w->w1 = ptr;
22      ptr += p->n_layers * p->dim * p->hidden_dim;
```

```
23        w->w2 = ptr;
24        ptr += p->n_layers * p->hidden_dim * p->dim;
25        w->w3 = ptr;
26        ptr += p->n_layers * p->dim * p->hidden_dim;
27        w->rms_final_weight = ptr;
28        ptr += p->dim;
29        w->freq_cis_real = ptr;
30        int head_size = p->dim / p->n_heads;
31        ptr += p->seq_len * head_size / 2;
32        w->freq_cis_imag = ptr;
33        ptr += p->seq_len * head_size / 2;
34        w->wcls = shared_weights ? w->token_embedding_table : ptr;
     }
```

## 3、神经网络模块

这些方法用于实现神经网络模型，具体来说是Transformer模型。以下是每个方法的简要描述：

1）accum(float *a, float b, int size)：此方法将数组b的每个元素加到数组a的相应元素上。

2）rmsnorm(float o, float x, float* weight, int size)：此方法使用作为参数提供的权重向量，在输入向量x上执行RMS归一化。结果存储在输出向量o中。

3）softmax(float* x, int size)：此方法对输入向量x应用softmax函数，该函数是类别集上的概率分布。结果是相同类别集上的概率分布，输出向量中的每个元素表示相应类别的概率。

4）matmul(float* xout, float* x, float* w, int n, int d)：此方法在输入向量x和权重矩阵w之间执行矩阵乘法，得到输出向量xout。矩阵和向量的维度由参数n和d指定。

5）transformer(int token, int pos, Config* p, RunState* s, TransformerWeights* w)：此方法为给定的输入token和位置实现Transformer模型。它使用提供的权重计算给定输入的输出向量，并将结果存储在RunState对象的输出向量中。

6）sample(float* probabilities, int n)：此方法从由输入向量probabilities指定的概率分布中抽取一个元素。向量的长度由整数参数n给出。

7）argmax(float* v, int n)：此方法返回输入向量v中最大元素的索引。向量的长度由整数参数n给出。

```
 4   // ---------------------------------------------------------------------------
 5   // neural net blocks
 6
 7   void accum(float *a, float *b, int size) {
 8       for (int i = 0; i < size; i++) {
 9           a[i] += b[i];
10       }
11   }
12
13   void rmsnorm(float* o, float* x, float* weight, int size) {
14       // calculate sum of squares
15       float ss = 0.0f;
16       for (int j = 0; j < size; j++) {
17           ss += x[j] * x[j];
18       }
19       ss /= size;
20       ss += 1e-5f;
```

```c
21        ss = 1.0f / sqrtf(ss);
22        // normalize and scale
23        for (int j = 0; j < size; j++) {
24            o[j] = weight[j] * (ss * x[j]);
25        }
26    }
27
28    void softmax(float* x, int size) {
29        // find max value (for numerical stability)
30        float max_val = x[0];
31        for (int i = 1; i < size; i++) {
32            if (x[i] > max_val) {
33                max_val = x[i];
34            }
35        }
36        // exp and sum
37        float sum = 0.0f;
38        for (int i = 0; i < size; i++) {
39            x[i] = expf(x[i] - max_val);
40            sum += x[i];
41        }
42        // normalize
43        for (int i = 0; i < size; i++) {
44            x[i] /= sum;
45        }
46    }
47
48    void matmul(float* xout, float* x, float* w, int n, int d) {
49        // W (d,n) @ x (n,) -> xout (d,)
50        // by far the most amount of time is spent inside this little function
51        int i;
52        #pragma omp parallel for private(i)
53        for (i = 0; i < d; i++) {
54            float val = 0.0f;
55            for (int j = 0; j < n; j++) {
56                val += w[i * n + j] * x[j];
57            }
58            xout[i] = val;
59        }
60    }
61
62    void transformer(int token, int pos, Config* p, RunState* s, TransformerWeights* w) {
63
64        // a few convenience variables
65        float *x = s->x;
66        int dim = p->dim;
67        int hidden_dim =  p->hidden_dim;
68        int head_size = dim / p->n_heads;
69
70        // copy the token embedding into x
71        float* content_row = &(w->token_embedding_table[token * dim]);
72        memcpy(x, content_row, dim*sizeof(*x));
73
74        // pluck out the "pos" row of freq_cis_real and freq_cis_imag
75        float* freq_cis_real_row = w->freq_cis_real + pos * head_size / 2;
76        float* freq_cis_imag_row = w->freq_cis_imag + pos * head_size / 2;
```

```c
76
77        // forward all the layers
78        for(int l = 0; l < p->n_layers; l++) {
79
80            // attention rmsnorm
81            rmsnorm(s->xb, x, w->rms_att_weight + l*dim, dim);
82
83            // qkv matmuls for this position
84            matmul(s->q, s->xb, w->wq + l*dim*dim, dim, dim);
85            matmul(s->k, s->xb, w->wk + l*dim*dim, dim, dim);
86            matmul(s->v, s->xb, w->wv + l*dim*dim, dim, dim);
87
88            // apply RoPE rotation to the q and k vectors for each head
89            for (int h = 0; h < p->n_heads; h++) {
90                // get the q and k vectors for this head
91                float* q = s->q + h * head_size;
92                float* k = s->k + h * head_size;
93                // rotate q and k by the freq_cis_real and freq_cis_imag
94                for (int i = 0; i < head_size; i+=2) {
95                    float q0 = q[i];
96                    float q1 = q[i+1];
97                    float k0 = k[i];
98                    float k1 = k[i+1];
99                    float fcr = freq_cis_real_row[i/2];
100                   float fci = freq_cis_imag_row[i/2];
101                   q[i]   = q0 * fcr - q1 * fci;
102                   q[i+1] = q0 * fci + q1 * fcr;
103                   k[i]   = k0 * fcr - k1 * fci;
104                   k[i+1] = k0 * fci + k1 * fcr;
105               }
106           }
107
108           // save key,value at this time step (pos) to our kv cache
109           int loff = l * p->seq_len * dim; // kv cache layer offset for convenience
110           float* key_cache_row = s->key_cache + loff + pos * dim;
111           float* value_cache_row = s->value_cache + loff + pos * dim;
112           memcpy(key_cache_row, s->k, dim*sizeof(*key_cache_row));
113           memcpy(value_cache_row, s->v, dim*sizeof(*value_cache_row));
114
115           // multihead attention. iterate over all heads
116           int h;
117           #pragma omp parallel for private(h)
118           for (h = 0; h < p->n_heads; h++) {
119               // get the query vector for this head
120               float* q = s->q + h * head_size;
121               // attention scores for this head
122               float* att = s->att + h * p->seq_len;
123               // iterate over all timesteps, including the current one
124               for (int t = 0; t <= pos; t++) {
125                   // get the key vector for this head and at this timestep
126                   float* k = s->key_cache + loff + t * dim + h * head_size;
127                   // calculate the attention score as the dot product of q and k
128                   float score = 0.0f;
129                   for (int i = 0; i < head_size; i++) {
130                       score += q[i] * k[i];
```

```
131                    }
132                    score /= sqrtf(head_size);
133                    // save the score to the attention buffer
134                    att[t] = score;
135                }
136
137                // softmax the scores to get attention weights, from 0..pos inclusively
138                softmax(att, pos + 1);
139
140                // weighted sum of the values, store back into xb
141                float* xb = s->xb + h * head_size;
142                memset(xb, 0, head_size * sizeof(float));
143                for (int t = 0; t <= pos; t++) {
144                    // get the value vector for this head and at this timestep
145                    float* v = s->value_cache + loff + t * dim + h * head_size;
146                    // get the attention weight for this timestep
147                    float a = att[t];
148                    // accumulate the weighted value into xb
149                    for (int i = 0; i < head_size; i++) {
150                        xb[i] += a * v[i];
151                    }
152                }
153            }
154
155            // final matmul to get the output of the attention
156            matmul(s->xb2, s->xb, w->wo + l*dim*dim, dim, dim);
157
158            // residual connection back into x
159            accum(x, s->xb2, dim);
160
161            // ffn rmsnorm
162            rmsnorm(s->xb, x, w->rms_ffn_weight + l*dim, dim);
163
164            // Now for FFN in PyTorch we have: self.w2(F.silu(self.w1(x)) * self.w3(x))
165            // first calculate self.w1(x) and self.w3(x)
166            matmul(s->hb, s->xb, w->w1 + l*dim*hidden_dim, dim, hidden_dim);
167            matmul(s->hb2, s->xb, w->w3 + l*dim*hidden_dim, dim, hidden_dim);
168
169            // F.silu; silu(x)=x*σ(x),where σ(x) is the logistic sigmoid
170            for (int i = 0; i < hidden_dim; i++) {
171                s->hb[i] = s->hb[i] * (1.0f / (1.0f + expf(-s->hb[i])));
172            }
173
174            // elementwise multiply with w3(x)
175            for (int i = 0; i < hidden_dim; i++) {
176                s->hb[i] = s->hb[i] * s->hb2[i];
177            }
178
179            // final matmul to get the output of the ffn
180            matmul(s->xb, s->hb, w->w2 + l*dim*hidden_dim, hidden_dim, dim);
181
182            // residual connection
183            accum(x, s->xb, dim);
184        }
185
```

```
186        // final rmsnorm
187        rmsnorm(x, x, w->rms_final_weight, dim);
188
189        // classifier into logits
190        matmul(s->logits, x, w->wcls, p->dim, p->vocab_size);
191    }
192
193    int sample(float* probabilities, int n) {
194        // sample index from probabilities, they must sum to 1
195        float r = (float)rand() / (float)RAND_MAX;
196        float cdf = 0.0f;
197        for (int i = 0; i < n; i++) {
198            cdf += probabilities[i];
199            if (r < cdf) {
200                return i;
201            }
202        }
203        return n - 1; // in case of rounding errors
204    }
205
206    int argmax(float* v, int n) {
207        // return argmax of v in elements 0..n
208        int max_i = 0;
209        float max_p = v[0];
210        for (int i = 1; i < n; i++) {
211            if (v[i] > max_p) {
212                max_i = i;
213                max_p = v[i];
214            }
215        }
        return max_i;
    }
```

## 4、main函数入口

程序通过读取已经训练好的模型和分词器，生成指定长度的文本。下面是主函数的作用：

1）解析命令行参数，包括模型文件路径、温度和步数等，如果缺少必要的参数，则打印出用法说明并返回1。

2）读取模型文件，并将其映射到内存中，同时读取词汇表文件。

3）初始化运行状态，包括分配内存和初始化状态。

4）生成文本，通过不断调用Transformer模型并使用softmax函数和采样方法获取下一个词语，直到生成指定长度的文本。

5）打印生成的文本，并计算生成速度。

6）释放分配的内存，关闭文件句柄。

程序主要的工作是在while循环中完成的，每次循环生成一个词语，并更新运行状态。生成的词语通过printf函数打印出来，最后计算生成速度并输出。

```
1
```

```c
int main(int argc, char *argv[]) {

    // poor man's C argparse
    char *checkpoint = NULL;  // e.g. out/model.bin
    float temperature = 0.9f; // e.g. 1.0, or 0.0
    int steps = 256;          // max number of steps to run for, 0: use seq_len
    // 'checkpoint' is necessary arg
    if (argc < 2) {
        printf("Usage: %s <checkpoint_file> [temperature] [steps]\n", argv[0]);
        return 1;
    }
    if (argc >= 2) {
        checkpoint = argv[1];
    }
    if (argc >= 3) {
        // optional temperature. 0.0 = (deterministic) argmax sampling. 1.0 = baseline
        temperature = atof(argv[2]);
    }
    if (argc >= 4) {
        steps = atoi(argv[3]);
    }

    // seed rng with time. if you want deterministic behavior use temperature 0.0
    srand((unsigned int)time(NULL));

    // read in the model.bin file
    Config config;
    TransformerWeights weights;
    int fd = 0;         // file descriptor for memory mapping
    float* data = NULL; // memory mapped data pointer
    long file_size;     // size of the checkpoint file in bytes
    {
        FILE *file = fopen(checkpoint, "rb");
        if (!file) { printf("Couldn't open file %s\n", checkpoint); return 1; }
        // read in the config header
        if(fread(&config, sizeof(Config), 1, file) != 1) { return 1; }
        // negative vocab size is hacky way of signaling unshared weights. bit yikes.
        int shared_weights = config.vocab_size > 0 ? 1 : 0;
        config.vocab_size = abs(config.vocab_size);
        // figure out the file size
        fseek(file, 0, SEEK_END); // move file pointer to end of file
        file_size = ftell(file); // get the file size, in bytes
        fclose(file);
        // memory map the Transformer weights into the data pointer
        fd = open(checkpoint, O_RDONLY); // open in read only mode
        if (fd == -1) { printf("open failed!\n"); return 1; }
        data = mmap(NULL, file_size, PROT_READ, MAP_PRIVATE, fd, 0);
        if (data == MAP_FAILED) { printf("mmap failed!\n"); return 1; }
        float* weights_ptr = data + sizeof(Config)/sizeof(float);
        checkpoint_init_weights(&weights, &config, weights_ptr, shared_weights);
    }
    // right now we cannot run for more than config.seq_len steps
    if (steps <= 0 || steps > config.seq_len) { steps = config.seq_len; }
```

```c
57          // read in the tokenizer.bin file
58      char** vocab = (char**)malloc(config.vocab_size * sizeof(char*));
59      {
60          FILE *file = fopen("tokenizer.bin", "rb");
61          if (!file) { printf("Couldn't load tokenizer.bin\n"); return 1; }
62          int len;
63          for (int i = 0; i < config.vocab_size; i++) {
64              if(fread(&len, sizeof(int), 1, file) != 1) { return 1; }
65              vocab[i] = (char *)malloc(len + 1);
66              if(fread(vocab[i], len, 1, file) != 1) { return 1; }
67              vocab[i][len] = '\0'; // add the string terminating token
68          }
69          fclose(file);
70      }
71
72      // create and init the application RunState
73      RunState state;
74      malloc_run_state(&state, &config);
75
76      // the current position we are in
77      long start = 0; // used to time our code, only initialized after first iteration
78      int next;
79      int token = 1; // 1 = BOS token in Llama-2 sentencepiece
80      int pos = 0;
81      printf("<s>\n"); // explicit print the initial BOS token (=1), stylistically symmetric
82      while (pos < steps) {
83
84          // forward the transformer to get logits for the next token
85          transformer(token, pos, &config, &state, &weights);
86
87          // sample the next token
88          if(temperature == 0.0f) {
89              // greedy argmax sampling
90              next = argmax(state.logits, config.vocab_size);
91          } else {
92              // apply the temperature to the logits
93              for (int q=0; q<config.vocab_size; q++) { state.logits[q] /= temperature; }
94              // apply softmax to the logits to get the probabilities for next token
95              softmax(state.logits, config.vocab_size);
96              // we now want to sample from this distribution to get the next token
97              next = sample(state.logits, config.vocab_size);
98          }
99          // following BOS token (1), sentencepiece decoder strips any leading whitespace (see PR #89)
100         char *token_str = (token == 1 && vocab[next][0] == ' ') ? vocab[next]+1 : vocab[next];
101         printf("%s", token_str);
102         fflush(stdout);
103
104         // advance forward
105         token = next;
106         pos++;
107         // init our timer here because the first iteration is slow due to memmap
108         if (start == 0) { start = time_in_ms(); }
109     }
110
111     // report achieved tok/s
```

```
112        long end = time_in_ms();
113        printf("\nachieved tok/s: %f\n", (steps-1) / (double)(end-start)*1000);
114
115        // memory and file handles cleanup
116        free_run_state(&state);
117        for (int i = 0; i < config.vocab_size; i++) { free(vocab[i]); }
118        free(vocab);
119        if (data != MAP_FAILED) munmap(data, file_size);
120        if (fd != -1) close(fd);
           return 0;
       }
```

至此，你已经基本知道怎么训练llama2了，并且还知道训练跟推理大致是怎么实现的，虽然没有对代码逐一细讲，但也觉得没必要一一细讲，只要把大致流程说清楚，其余部分自然也就通了，LLM是不是其实也并不难，跟着教程玩起来吧！！！

参考文献：

> https://zhuanlan.zhihu.com/p/634267984