

使用 Megatron 预训练 GPT2

1. 下载代码

```
1 git clone https://github.com/NVIDIA/Megatron-LM.git
2 cd Megatron-LM
3 git checkout 23.06
```

2. 下载 gpt2 数据集

```
1 wget https://huggingface.co/bigscience/misc-test-data/resolve/main/stas/oscar-1GB.jsonl.xz
2 wget https://s3.amazonaws.com/models.huggingface.co/bert/gpt2-vocab.json
3 wget https://s3.amazonaws.com/models.huggingface.co/bert/gpt2-merges.txt
4 xz -d oscar-1GB.jsonl.xz
```

3. 数据预处理

```
1 jsonfile="/workspace/dataset/oscar-1GB.jsonl"
2 vocabfile="/workspace/dataset/gpt2-vocab.json"
3 mergefile="/workspace/dataset/gpt2-merges.txt"
4 prefix="my_gpt2"
5
6 python tools/preprocess_data.py \
7     --input $jsonfile \
8     --output-prefix $prefix \
9     --vocab-file $vocabfile \
10    --merge-file $mergefile \
11    --dataset-impl mmap \
12    --tokenizer-type GPT2BPETokenizer \
13    --workers 32 \
14    --log-interval 500
```

4. 可以根据情况, 修改下 examples/pretrain_gpt_distributed_with_mp.sh, 如修改模型配置参数, 张量并行度, pipeline 并行度等。

```
1  #!/bin/bash
2  export CUDA_DEVICE_MAX_CONNECTIONS=1
3
4  GPUS_PER_NODE=8
5  # Change for multinode config
6  MASTER_ADDR=localhost
7  MASTER_PORT=6000
8  NNODES=1
9  NODE_RANK=0
10 WORLD_SIZE=$((($GPUS_PER_NODE*$NNODES))
11
12 CHECKPOINT_PATH=/workspace/checkpoint
13 VOCAB_FILE=/workspace/dataset/gpt2-vocab.json
14 MERGE_FILE=/workspace/dataset/gpt2-merges.txt
15 DATA_PATH=/workspace/Megatron-LM/my_gpt2_text_document
16
17 DISTRIBUTED_ARGS="
18     --nproc_per_node $GPUS_PER_NODE \
19     --nnodes $NNODES \
20     --node_rank $NODE_RANK \
21     --master_addr $MASTER_ADDR \
22     --master_port $MASTER_PORT
23 "
24
25 GPT_ARGS="
26     --tensor-model-parallel-size 1 \
27     --pipeline-model-parallel-size 8 \
28     --num-layers 48 \
29     --hidden-size 2048 \
30     --num-attention-heads 32 \
31     --seq-length 1024 \
32     --max-position-embeddings 1024 \
33     --micro-batch-size 4 \
34     --global-batch-size 16 \
35     --lr 0.00015 \
36     --train-iters 500000 \
37     --lr-decay-iters 32000 \
38     --lr-decay-style cosine \
39     --min-lr 1.0e-5 \
40     --weight-decay 1e-2 \
41     --lr-warmup-fraction .01 \
42     --clip-grad 1.0 \
43     --fp16
44 "
45
46 DATA_ARGS="
47     --data-path $DATA_PATH \
```

```
48     --vocab-file $VOCAB_FILE \  
49     --merge-file $MERGE_FILE \  
50     --data-impl mmap \  
51     --split 949,50,1  
52     "  
53  
54     OUTPUT_ARGS="  
55         --log-interval 100 \  
56         --save-interval 100000 \  
57         --eval-interval 100000 \  
58         --eval-iters 10  
59     "  
60  
61     torchrun $DISTRIBUTED_ARGS pretrain_gpt.py \  
62         $GPT_ARGS \  
63         $DATA_ARGS \  
64         $OUTPUT_ARGS \  
65         --distributed-backend nccl \  
66         --save $CHECKPOINT_PATH \  
67         --load $CHECKPOINT_PATH
```

5. 开始训练

```
1 ./examples/pretrain_gpt_distributed_with_mp.sh
```