

MoE 优化

Mixtral-8x7B 模型优化

该模型有8个专家，每个专家7B规模，模型总大小约为 96G, 因此用 2 卡推理较为合适。

从组网看，该模型只是在 llama 的基础上，将 mlp 部分 变成了 moe 部分，moe 组网也较为简单，核心代码如下：

```
1 class FeedForward(nn.Module):
2     def __init__(
3         self,
4         config
5     ):
6         """
7         Initialize the FeedForward module.
8
9         Args:
10             dim (int): Input dimension.
11             hidden_dim (int): Hidden dimension of the feedforward layer.
12             multiple_of (int): Value to ensure hidden dimension is a multiple
13                 of this value.
14             ffn_dim_multiplier (float, optional): Custom multiplier for
15                 hidden dimension. Defaults to None.
16
17         Attributes:
18             w1 (ColumnParallelLinear): Linear transformation for the first
19                 layer.
20             w2 (RowParallelLinear): Linear transformation for the second
21                 layer.
22             w3 (ColumnParallelLinear): Linear transformation for the third
23                 layer.
24
25         """
26         super().__init__()
27         self.w1 = nn.Linear(
28             config.hidden_size, config.intermediate_size, bias=False
29         )
30         self.w2 = nn.Linear(
31             config.intermediate_size, config.hidden_size, bias=False
32         )
33         self.w3 = nn.Linear(
```

```

30         config.hidden_size, config.intermediate_size, bias=False
31     )
32
33     def forward(self, x):
34         return self.w2(F.silu(self.w1(x)) * self.w3(x))
35
36 class MoE(nn.Module):
37     def __init__(
38         self,
39         config,
40     ):
41         super().__init__()
42         self.config = config
43         num_experts = config.num_experts
44         self.experts = nn.ModuleList([FeedForward(config) for i in
range(num_experts)])
45         self.gate = nn.Linear(config.hidden_size, num_experts, bias=False)
46         self.num_experts_per_token = config.num_experts_per_token
47
48     def forward(self, x):
49         orig_shape = x.shape #(batch, seq_len, d_model)
50         x = x.view(-1, x.shape[-1]) #(num_tokens, d_model)
51
52         scores = self.gate(x) #(num_tokens, num_experts)
53
54         expert_weights, expert_indices = torch.topk(scores,
self.num_experts_per_token, dim=-1) #(num_tokens, k)
55
56         expert_weights = expert_weights.softmax(dim=-1) #(num_tokens,
num_experts_per_token)
57         flat_expert_indices = expert_indices.view(-1) #(num_tokens *
num_experts_per_token)
58
59         x = x.repeat_interleave(self.num_experts_per_token, dim=0) #
(num_tokens * num_experts_per_token, d_model) -->
60
61         y = torch.empty_like(x)
62         for i, expert in enumerate(self.experts):
63             print(f"flat_expert_indices == i == {flat_expert_indices == i}")
64             y[flat_expert_indices == i] = expert(x[flat_expert_indices == i])
65         y = (y.view(*expert_weights.shape, -1) *
expert_weights.unsqueeze(-1)).sum(dim=1) #(num_tokens *
num_experts_per_token, d_model) --> (num_tokens, num_experts_per_token,
d_model)
66
67         # --> (num_tokens, 1, d_model)
68         return y.view(*orig_shape)

```

注意这里对 gating 后的分数 scores 的处理方式：

top_k + softmax 归一化

```
expert_weights, expert_indices = torch.topk(scores, self.num_experts_per_token,
expert_weights = expert_weights.softmax(dim=-1) #(num_tokens, num_experts_per_to
```

Vllm 推理

Vllm 较新的代码已经支持了该模型的部署，但是使用的是 EP 并行，forward 部分，采用 for 循环，性能可能较差，下面是具体的测试数据：

A100，双卡

elapsed	Prompt = 1	Prompt = 1000	Prompt = 2000	Prompt = 4000
Batch = 1	39ms	42ms	44ms	48ms

注意 vllm 的实现，对 gating 后的分数（router_logits）处理方式不同

Softmax + topk + 除以sum归一化

```
routing_weights = F.softmax(router_logits, dim=1, dtype=torch.float)
routing_weights, selected_experts = torch.topk(routing_weights,
self.top_k,
dim=-1)
routing_weights /= routing_weights.sum(dim=-1, keepdim=True)
```

https://github.com/vllm-project/vllm/blob/main/vllm/model_executor/models/mixtral.py#L143

乍一看，代码实现上确实不同，但是简单地做下数学推导，就会发现，是一致的！

使用 group gemm 优化

该 moe 没有 EP 均衡措施，所以无法使用 BMM 替代 for 循环进行加速计算，因此可以使用 group gemm

参考 trt-llm 的 moe 实现，将 trt-llm 中的 moe kernel 抽离出来然后修改，编译为单独的算子库，方便即插即用！

未来可以继续拓展、维护，变成支持更多算子的算子库。

<https://git.singularity-ai.com/shengying.wei/MoE>

运行里面的单测，可以检测算子的正确性。

使用vllm 框架的最新代码，重新组网 mixtral ，利用我们的 MoE 算子库加速

如上所说，vllm 中 mixtral 使用的 ep 并行，重新组网使用 TP 并行即可：

```
class MixtralMoE(nn.Module):
    def __init__(
        self,
        config: MixtralConfig,
        linear_method: Optional[LinearMethodBase] = None,
    ) -> None:
        super().__init__()
        self.tp_size = get_tensor_model_parallel_world_size()
        self.hidden_size = config.hidden_size
        self.inter_size = config.intermediate_size // self.tp_size
        self.top_k = config.num_experts_per_tok
        self.num_experts = config.num_local_experts
        params_dtype = torch.get_default_dtype()
        self.gate = ReplicatedLinear(config.hidden_size,
                                     self.num_experts,
                                     bias=False,
                                     linear_method=None)

        from torch.nn.parameter import Parameter
        self.moe_up_proj_weight = Parameter(torch.empty((self.num_experts, self.hidden_size, 2 * self.inter_size),
                                                         device=torch.cuda.current_device(),
                                                         dtype=params_dtype))

        self.moe_down_proj_weight = Parameter(torch.empty((self.num_experts, self.inter_size, self.hidden_size),
                                                         device=torch.cuda.current_device(),
                                                         dtype=params_dtype))

        import moe_ops
        self.run_moe_fc = moe_ops.run_moe_fc

    def forward(self, hidden_states):
        batch_size, sequence_length, hidden_dim = hidden_states.shape
        hidden_states = hidden_states.view(-1, hidden_dim) # (num_tokens, hidden_dim)
        num_rows = hidden_states.size(0)

        gating_output, _ = self.gate(hidden_states) # (num_tokens, n_experts)

        moe_output = self.run_moe_fc(hidden_states, gating_output, self.moe_up_proj_weight,
                                     "Swiglu", self.moe_down_proj_weight, num_rows, self.top_k)

        if self.tp_size == 1:
            return moe_output
        torch.distributed.all_reduce(moe_output, group=get_tensor_model_parallel_group())
        return moe_output.view(batch_size, sequence_length, hidden_dim)
```

如图，我们直接把 8 个专家合并在一起（self.moe_up_proj_weight, self.moe_down_proj_weight），然后进行切分，forward 函数中，使用我们 moe 算子库中封装好的支持 torch tensor 的 api，组网变得非常简洁！

最后，因为组网的改变，把 8 个专家合并在一起，所以模型切分逻辑稍微复杂一些，切分专家的主要代码如下：

```
1  if ("block_sparse_moe.experts." in name): # like
    "model.layers.31.block_sparse_moe.experts.7.w2.weight"
2      substrings = name.split('.')
3      layer_num = int(substrings[2])
4      expert_num = int(substrings[5])
5      if("w1" in name or "w3" in name): # 列切
6          loaded_weight = loaded_weight.T # (hidden_size, inter_size)
7          shard_size = loaded_weight.shape[-1] // tp_size
```

```

8         loaded_weight_slice = loaded_weight[:, tp_rank * shard_size :
          (tp_rank + 1) * shard_size] # (hidden_size, inter_size // tp_size)
9         coresponded_param_name = f"model.layers.
          {layer_num}.block_sparse_moe.moe_up_proj_weight"
10        param = params_dict[coresponded_param_name] #(num_expert, hidden_size,
          2 * inter_size // tp_size)
11        if ("w1" in name):
12            param_slice = param[expert_num, :, shard_size :]
13            assert loaded_weight_slice.shape == param_slice.shape
14            param_slice.data.copy_(loaded_weight_slice)
15            continue
16        else: # w3
17            param_slice = param[expert_num, :, : shard_size]
18            assert loaded_weight_slice.shape == param_slice.shape
19            param_slice.data.copy_(loaded_weight_slice)
20            continue
21        elif ("w2" in name): # 行切
22            loaded_weight = loaded_weight.T # (inter_size, hidden_size)
23            shard_size = loaded_weight.shape[0] // tp_size
24            loaded_weight_slice = loaded_weight[tp_rank * shard_size : (tp_rank +
          1) * shard_size, :] # (inter_size // tp_size, hidden_size)
25            coresponded_param_name = f"model.layers.
          {layer_num}.block_sparse_moe.moe_down_proj_weight"
26            param = params_dict[coresponded_param_name] #(num_expert, inter_size
          // tp_size, hidden_size, )
27            param_slice = param[expert_num]
28            assert loaded_weight_slice.shape == param_slice.shape
29            param_slice.data.copy_(loaded_weight_slice)
30            continue
31        else:
32            raise ValueError(f"invalid params: {name}")
33        continue

```

完整组网和切分代码见附录。

性能测试

Mixtral-8x7B，平均每个token耗时，包含首token

Vllm 原生推理（双卡），EP 并行

elapsed	Prompt = 1	Prompt = 1000	Prompt = 2000	Prompt = 4000
Batch = 1	39 ms	42 ms	44 ms	48 ms

使用 group gemm 优化后，TP 并行

elapsed	Prompt = 1	Prompt = 1000	Prompt = 2000	Prompt = 4000
Batch = 1	18.64 ms	19.82 ms	21.70 ms	24.00 ms

Vllm 最新的代码支持 cuda-gprah，开启后：

首token / 生成耗时	Prompt = 1	Prompt = 1000	Prompt = 2000	Prompt = 4000
Batch = 1	26 ms 16.64 ms	129 ms 18.23 ms	240 ms 19.3 ms	464 ms 21.79 ms

cuda-graph 约有 5% 到 10% 的性能收益。

手动配置 cutlass gemm config，避免内部自己遍历寻优

首token / 生成耗时	Prompt = 1	Prompt = 1000	Prompt = 2000	Prompt = 4000
Batch = 1	25.6 ms / 16.00 ms	106 ms 17.40 ms	192 ms 18.4 ms	372 ms 20.35 ms

生成性能提升不明显：约 4%

首token 性能有较大提升：在 prompt 长度1000 以上，有 20%+ 的提升

四卡推理：

elapsed	Prompt = 1	Prompt = 1000	Prompt = 2000	Prompt = 4000
Batch = 1	11.61 ms	12.61 ms	13.45 ms	14.75 ms

精度已对齐。

```
1 prompt = "The capital of America is"
```

```
2
3 vllm 原生推理输出:
4 ' a city that is full of history and culture. It is also a city that is full
  of things to do. If you are looking for a place to stay in Washington DC,
  then you should consider staying in a hotel. There are many hotels in
  Washington DC, and each one has its own unique features.\n\n## The Best
  Hotels in Washington DC\n\nThere are many hotels in Washington DC, and each
  one has its own unique features. Some of the best hotels in Washington DC
  include'
5
6 优化后输出:
7 ' a city that is full of history and culture. It is also a city that is full
  of things to do. If you are looking for a place to stay in Washington DC,
  then you should consider staying in a hotel. There are many hotels in
  Washington DC, and each one has its own unique features.\n\n## The Best
  Hotels in Washington DC\n\nThere are many hotels in Washington DC, and each
  one has its own unique features. Some of the best hotels in Washington DC
  include'
```

部署细节

vllm: <https://git.singularity-ai.com/shengying.wei/vllm-moe> moe_develop_0117 分支

Moe 算子库:<https://git.singularity-ai.com/shengying.wei/MoE> master 分支

Docker: gpu-image-cn-shanghai.cr.volces.com/gpu-train/skylm2:v0.1

代码细节

浅谈 moe 中 kernel 的设计思想

了解 kernel 的算法设计, 才能吃透代码, 感受 cuda 的艺术。

topkGatingSoftmax

该 kernel 完成对 gating 之后的分数的处理, 对应的 python 逻辑如下:

```
1 expert_scales = F.softmax(input_dict["gating_output"], dim=-1)
2 expert_scales, experts_for_row = torch.topk(expert_scales, k, dim=-1)
3 expert_scales /= expert_scales.sum(dim=-1, keepdim=True) #除和归一化
```

该 kernel 处理的任务规模不大, 输入的 shape 为 (num_tokens, num_experts), 在 mixtral 中, 专家个数为 8。

1. 向量化读取 (提升内存读取效率):

- a. 该kernel约定，每个线程每次加载的字节个数为：BYTES_PER_LDG，那么每次可以加载的元素个数为：ELTS_PER_LDG = BYTES_PER_LDG / sizeof(T)，因此，把ELTS_PER_LDG个元素看成一个向量，线程只需要一个内存加载操作即可完成该向量（ELTS_PER_LDG个元素）的加载，ELTS_PER_LDG也即向量的长度。

2. 算法设计 (thread block的协同)

- a. 每一行都有 num_experts 个元素（专家），因此至少用一个线程束 warp 来负责一行，那么每个线程负责的元素个数为：EXPERTS / WARP_SIZE，即 EXPERTS / (ELTS_PER_LDG * WARP_SIZE) 个向量。如果 num_experts 较少，则**每个线程只负责一个向量**，一个线程束可以处理多行。代码如下：

```
1 static constexpr int VECs_PER_THREAD = std::max(1, EXPERTS / (ELTS_PER_LDG * WARP_SIZE));
2 static constexpr int VPT = VECs_PER_THREAD * ELTS_PER_LDG; //This gives the total number of elements (or vectors, depending on how you interpret it) processed by a single thread.
```

- b. 进一步便可以求出，一行需要多少个线程来处理，以及一个线程束可以处理多少行

```
1 static constexpr int THREADS_PER_ROW = EXPERTS / VPT; // 每一行需要 THREADS_PER_ROW 个 thread 来处理
2 static constexpr int ROWS_PER_WARP = WARP_SIZE / THREADS_PER_ROW; // 一个线程束 wrap 可以处理 ROWS_PER_WARP 行
```

3. 线程组织方式 (grid_dim and block_dim)

- a. 上述计算出，一个线程束 warp 可以处理 ROWS_PER_WARP 行，那么对于输入的数据规模，则可以计算出一共需要多少个线程束，每次 thread block 最多有 WARPS_PER_TB 个线程束，因此一共需要的 block 数目也可随之计算出来：

```
1 const int num_warps = (num_rows + ROWS_PER_WARP - 1) / ROWS_PER_WARP; //一行需要 ROWS_PER_WARP 个 warp 来处理，所以一共需要 num_warps 个 wrap
2 const int num_blocks = (num_warps + WARPS_PER_TB - 1) / WARPS_PER_TB; //一个 thread block 中最多有 WARPS_PER_TB 个 warp，所以需要 num_blocks 个 thread block
3 dim3 block_dim(WARP_SIZE, WARPS_PER_TB); //一个 thread block 中有 WARPS_PER_TB 个 warp
```

4. 线程束 warp 内部的合作方式

- a. 一个 warp 可能会负责多行，所以一个 warp 中的32个线程，也会被分割为多个sub-group，每个子部分 sub-group 处理一行。所以需要计算出，每个线程需要处理哪一行。当然，专家的个数可能是 32 的倍数，那么一个 warp 可能只处理一行，这时候，每个线程需要负责多个向量。

```
1 // The threads in a warp are split into sub-groups that will work on a row.
2 // We compute row offset for each thread sub-group
3 const int thread_row_in_warp = threadIdx.x / THREADS_PER_ROW;
4 const int thread_row = warp_base_row + thread_row_in_warp; // 当前线程需要处理的对应行
```

b. 合并内存访问

正如上面所说，假如 warp 内的每个线程需要处理 2 个向量，那么处理方式如下：

[0][1][2][3].....[31][0][1][2].....[31]

而不是：

[0][0][1][1].....[31][31]

对应代码部分：LDG_PER_THREAD 即每个线程需要处理 LDG_PER_THREAD 个向量，注意代码中的跳跃式加载。

```
1 #pragma unroll
2 for (int ii = 0; ii < LDG_PER_THREAD; ++ii) // 将全局内存中的 input 数据，加载到 row_chunk_input 来
3 {
4     row_chunk_vec_ptr[ii] = vec_thread_read_ptr[ii * THREADS_PER_ROW];
5 }
```

5. Shared memory

使用 cutlass 中提供的静态数组，将 global memory 中的input数据，加载到 shared memory中

```
1 cutlass::Array<T, VPT> row_chunk_input; // 每个线程需要处理 VPT 个元素
2 cutlass::Array<ComputeType, VPT> row_chunk =
   compute_type_converter(row_chunk_input); // 把加载进来的输入数据，转换为 float
```

6. 计算精度

- a. 由于有 softmax，需要求和以及指数运算，因此计算精度最好是 float，这要把输入的 half 或者 bfloat16 转换为 float。

```

1  using ComputeType = float;
2  using Converter = cutlass::NumericArrayConverter<ComputeType, T, VPT>;
3  Converter compute_type_converter;
4  cutlass::Array<ComputeType, VPT> row_chunk =
    compute_type_converter(row_chunk_input); // 把加载进来的输入数据, 转换为 float

```

7. 其他cuda kernel常用变成技巧

- a. 线程束内通信: __shfl_xor_sync、butterfly reduce 等。
- b. 强制精度类型转换

更多详细细节可以查看代码和相关注释。

单测代码

```

1  from vllm import LLM, SamplingParams
2  import time
3  prompt = [
4      "The capital of America is"
5  ]
6  max_tokens = 100
7  sampling_params = SamplingParams(temperature=0.8, top_p=0.95, top_k = 1,
    max_tokens = max_tokens, ignore_eos=True)
8
9  llm = LLM(model="/mnt/infra/weishengying/model/Mixtral-8x7B-v0.1",
    max_num_batched_tokens=32768, tensor_parallel_size=2, trust_remote_code=True,
    enforce_eager=False)
10
11  # while True:
12      # prompt = input("User:\t")
13      # if prompt == 'quit':
14          # break
15  warm_time = 2
16  repeat_time = 5
17
18  for _ in range(warm_time):
19      llm.generate(prompt, sampling_params)
20
21  start = time.time()
22  for _ in range(repeat_time):
23      outputs = llm.generate(prompt, sampling_params)
24  end = time.time()
25  print(f"total elapsed time: {(end-start)*1000/repeat_time}ms")
26  # Print the outputs.
27  for output in outputs:

```

```

28     prompt = output.prompt
29     generated_text = output.outputs[0].text
30     print("generate len: ", len(output.outputs[0].token_ids))
31     print(f"Generated text: {generated_text!r}")

```

Mixtral 组网和切分代码:

```

1  # coding=utf-8
2  # Adapted from
3  #
4  # https://github.com/huggingface/transformers/blob/v4.28.0/src/transformers/models/llama/modeling\_llama.py
5  # Copyright 2023 The vLLM team.
6  # Copyright 2022 EleutherAI and the HuggingFace Inc. team. All rights
7  # reserved.
8  #
9  # This code is based on EleutherAI's GPT-NeoX library and the GPT-NeoX
10 # and OPT implementations in this library. It has been modified from its
11 # original forms to accommodate minor architectural differences compared
12 # to GPT-NeoX and OPT used by the Meta AI team that trained the model.
13 #
14 # Licensed under the Apache License, Version 2.0 (the "License");
15 # you may not use this file except in compliance with the License.
16 # You may obtain a copy of the License at
17 #
18 # http://www.apache.org/licenses/LICENSE-2.0
19 #
20 # Unless required by applicable law or agreed to in writing, software
21 # distributed under the License is distributed on an "AS IS" BASIS,
22 # WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
23 # See the License for the specific language governing permissions and
24 # limitations under the License.
25 """Inference-only Mixtral model."""
26 from typing import List, Optional, Tuple
27
28 import numpy as np
29
30 import torch
31 import torch.nn.functional as F
32
33 from torch import nn
34 from transformers import MixtralConfig
35
36 from vllm.model_executor.input_metadata import InputMetadata
37 from vllm.model_executor.layers.attention import PagedAttention

```

```

36 from vllm.model_executor.layers.layernorm import RMSNorm
37 from vllm.model_executor.layers.linear import (LinearMethodBase,
38                                                ReplicatedLinear,
39                                                QKVParallelLinear,
40                                                RowParallelLinear)
41 from vllm.model_executor.layers.rotary_embedding import get_rope
42 from vllm.model_executor.layers.sampler import Sampler
43 from vllm.model_executor.layers.vocab_parallel_embedding import (
44     VocabParallelEmbedding, ParallelLMHead)
45 from vllm.model_executor.parallel_utils.communication_op import (
46     tensor_model_parallel_all_reduce)
47 from vllm.model_executor.parallel_utils.parallel_state import (
48     get_tensor_model_parallel_rank, get_tensor_model_parallel_world_size)
49 from vllm.model_executor.sampling_metadata import SamplingMetadata
50 from vllm.model_executor.weight_utils import (default_weight_loader,
51                                                hf_model_weights_iterator)
52 from vllm.sequence import SamplerOutput
53
54 KVCache = Tuple[torch.Tensor, torch.Tensor]
55
56 from vllm.model_executor.parallel_utils.parallel_state import (
57     get_tensor_model_parallel_world_size,
58     get_tensor_model_parallel_group,
59 )
60 from vllm.utils import NcclAllReduce
61
62 class MixtralMoE(nn.Module):
63     def __init__(
64         self,
65         config: MixtralConfig,
66         linear_method: Optional[LinearMethodBase] = None,
67     ) -> None:
68         super().__init__()
69         self.tp_size = get_tensor_model_parallel_world_size()
70         self.hidden_size = config.hidden_size
71         self.inter_size = config.intermediate_size
72         self.top_k = config.num_experts_per_tok
73         self.num_experts = config.num_local_experts
74         params_dtype = torch.get_default_dtype()
75         self.gate = ReplicatedLinear(config.hidden_size,
76                                     self.num_experts,
77                                     bias=False,
78                                     linear_method=None)
79
80         from torch.nn.parameter import Parameter
81         self.moe_up_proj_weight = Parameter(torch.empty((self.num_experts,
82                                                         self.hidden_size, 2 * self.inter_size // self.tp_size),

```

```

82     device=torch.cuda.current_device(),
83     dtype=params_dtype))
84
85     self.moe_down_proj_weight = Parameter(torch.empty((self.num_experts,
86     self.inter_size // self.tp_size, self.hidden_size),
87     device=torch.cuda.current_device(),
88     dtype=params_dtype))
89     import moe_ops
90     self.run_moe_fc = moe_ops.run_moe_fc
91
92     def forward(self, hidden_states):
93         batch_size, sequence_length, hidden_dim = hidden_states.shape
94         hidden_states = hidden_states.view(-1, hidden_dim) # (num_tokens,
95         hidden_dim)
96         num_rows = hidden_states.size(0)
97
98         gating_output, _ = self.gate(hidden_states) # (num_tokens, n_experts)
99
100        moe_output = self.run_moe_fc(hidden_states, gating_output,
101        self.moe_up_proj_weight,
102        "Swiglu", self.moe_down_proj_weight,
103        num_rows, self.top_k )
104        if self.tp_size == 1:
105            return moe_output
106        # torch.distributed.all_reduce(moe_output,
107        group=get_tensor_model_parallel_group())
108        moe_output = NcclAllReduce.all_reduce(moe_output)
109        return moe_output.view(batch_size, sequence_length, hidden_dim)
110
111    class MixtralAttention(nn.Module):
112
113    def __init__(self,
114        hidden_size: int,
115        num_heads: int,
116        num_kv_heads: int,
117        max_position: int = 4096 * 32,
118        rope_theta: float = 10000,
119        linear_method: Optional[LinearMethodBase] = None,
120        sliding_window: Optional[int] = None) -> None:
121        super().__init__()
122        self.hidden_size = hidden_size
123        tp_size = get_tensor_model_parallel_world_size()
124        self.total_num_heads = num_heads

```

```

120     assert self.total_num_heads % tp_size == 0
121     self.num_heads = self.total_num_heads // tp_size
122     self.total_num_kv_heads = num_kv_heads
123     if self.total_num_kv_heads >= tp_size:
124         # Number of KV heads is greater than TP size, so we partition
125         # the KV heads across multiple tensor parallel GPUs.
126         assert self.total_num_kv_heads % tp_size == 0
127     else:
128         # Number of KV heads is less than TP size, so we replicate
129         # the KV heads across multiple tensor parallel GPUs.
130         assert tp_size % self.total_num_kv_heads == 0
131     self.num_kv_heads = max(1, self.total_num_kv_heads // tp_size)
132     self.head_dim = hidden_size // self.total_num_heads
133     self.q_size = self.num_heads * self.head_dim
134     self.kv_size = self.num_kv_heads * self.head_dim
135     self.scaling = self.head_dim**-0.5
136     self.rope_theta = rope_theta
137     self.sliding_window = sliding_window
138
139     self.qkv_proj = QKVParallelLinear(
140         hidden_size,
141         self.head_dim,
142         self.total_num_heads,
143         self.total_num_kv_heads,
144         bias=False,
145         linear_method=linear_method,
146     )
147     self.o_proj = RowParallelLinear(
148         self.total_num_heads * self.head_dim,
149         hidden_size,
150         bias=False,
151         linear_method=linear_method,
152     )
153     self.rotary_emb = get_rope(
154         self.head_dim,
155         rotary_dim=self.head_dim,
156         max_position=max_position,
157         base=int(self.rope_theta),
158         is_neox_style=True,
159     )
160     self.attn = PagedAttention(
161         self.num_heads,
162         self.head_dim,
163         self.scaling,
164         num_kv_heads=self.num_kv_heads,
165         sliding_window=self.sliding_window,
166     )

```

```

167
168     def forward(
169         self,
170         positions: torch.Tensor,
171         hidden_states: torch.Tensor,
172         kv_cache: KVCache,
173         input_metadata: InputMetadata,
174     ) -> torch.Tensor:
175         qkv, _ = self.qkv_proj(hidden_states)
176         q, k, v = qkv.split([self.q_size, self.kv_size, self.kv_size], dim=-1)
177         q, k = self.rotary_emb(positions, q, k)
178         k_cache, v_cache = kv_cache
179         attn_output = self.attn(q, k, v, k_cache, v_cache, input_metadata)
180         output, _ = self.o_proj(attn_output)
181         return output
182
183 class MixtralDecoderLayer(nn.Module):
184
185     def __init__(
186         self,
187         config: MixtralConfig,
188         linear_method: Optional[LinearMethodBase] = None,
189     ) -> None:
190         super().__init__()
191         self.hidden_size = config.hidden_size
192         # Requires transformers > 4.32.0
193         rope_theta = getattr(config, "rope_theta", 10000)
194         self.self_attn = MixtralAttention(
195             hidden_size=self.hidden_size,
196             num_heads=config.num_attention_heads,
197             max_position=config.max_position_embeddings,
198             num_kv_heads=config.num_key_value_heads,
199             rope_theta=rope_theta,
200             sliding_window=config.sliding_window,
201             linear_method=linear_method)
202         self.block_sparse_moe = MixtralMoE(config=config,
203                                           linear_method=linear_method)
204         self.input_layernorm = RMSNorm(config.hidden_size,
205                                       eps=config.rms_norm_eps)
206         self.post_attention_layernorm = RMSNorm(config.hidden_size,
207                                                eps=config.rms_norm_eps)
208
209     def forward(
210         self,
211         positions: torch.Tensor,
212         hidden_states: torch.Tensor,
213         kv_cache: KVCache,

```



```

214         input_metadata: InputMetadata,
215         residual: Optional[torch.Tensor],
216     ) -> torch.Tensor:
217         # Self Attention
218         if residual is None:
219             residual = hidden_states
220             hidden_states = self.input_layernorm(hidden_states)
221         else:
222             hidden_states, residual = self.input_layernorm(
223                 hidden_states, residual)
224         hidden_states = self.self_attn(
225             positions=positions,
226             hidden_states=hidden_states,
227             kv_cache=kv_cache,
228             input_metadata=input_metadata,
229         )
230
231         # Fully Connected
232         hidden_states, residual = self.post_attention_layernorm(
233             hidden_states, residual)
234         hidden_states = self.block_sparse_moe(hidden_states)
235         return hidden_states, residual
236
237 class MixtralModel(nn.Module):
238
239     def __init__(
240         self,
241         config: MixtralConfig,
242         linear_method: Optional[LinearMethodBase] = None,
243     ) -> None:
244         super().__init__()
245         self.padding_idx = config.pad_token_id
246         self.vocab_size = config.vocab_size
247
248         self.embed_tokens = VocabParallelEmbedding(
249             config.vocab_size,
250             config.hidden_size,
251         )
252         self.layers = nn.ModuleList([
253             MixtralDecoderLayer(config, linear_method=linear_method)
254             for _ in range(config.num_hidden_layers)
255         ])
256         self.norm = RMSNorm(config.hidden_size, eps=config.rms_norm_eps)
257
258     def forward(
259         self,
260         input_ids: torch.Tensor,

```

```

261         positions: torch.Tensor,
262         kv_caches: List[KVCache],
263         input_metadata: InputMetadata,
264     ) -> torch.Tensor:
265         hidden_states = self.embed_tokens(input_ids)
266         residual = None
267         for i in range(len(self.layers)):
268             layer = self.layers[i]
269             hidden_states, residual = layer(positions, hidden_states,
270                                             kv_caches[i], input_metadata,
271                                             residual)
272         hidden_states, _ = self.norm(hidden_states, residual)
273         return hidden_states
274
275 class MixtralForCausalLM(nn.Module):
276
277     def __init__(
278         self,
279         config: MixtralConfig,
280         linear_method: Optional[LinearMethodBase] = None,
281     ) -> None:
282         super().__init__()
283         self.config = config
284         self.linear_method = linear_method
285         self.model = MixtralModel(config, linear_method)
286         self.lm_head = ParallelLMHead(config.vocab_size, config.hidden_size)
287         self.sampler = Sampler(config.vocab_size)
288         NcclAllReduce.init_comm()
289
290     def forward(
291         self,
292         input_ids: torch.Tensor,
293         positions: torch.Tensor,
294         kv_caches: List[KVCache],
295         input_metadata: InputMetadata,
296     ) -> torch.Tensor:
297         hidden_states = self.model(input_ids, positions, kv_caches,
298                                   input_metadata)
299         return hidden_states
300
301     def sample(
302         self,
303         hidden_states: Optional[torch.Tensor],
304         sampling_metadata: SamplingMetadata,
305     ) -> Optional[SamplerOutput]:
306         next_tokens = self.sampler(self.lm_head.weight, hidden_states,
307                                   sampling_metadata)

```

```

308         return next_tokens
309
310     def load_weights(self,
311                     model_name_or_path: str,
312                     cache_dir: Optional[str] = None,
313                     load_format: str = "auto",
314                     revision: Optional[str] = None):
315         stacked_params_mapping = [
316             # (param_name, shard_name, shard_id)
317             ("qkv_proj", "q_proj", "q"),
318             ("qkv_proj", "k_proj", "k"),
319             ("qkv_proj", "v_proj", "v"),
320         ]
321         params_dict = dict(self.named_parameters())
322         tp_size = get_tensor_model_parallel_world_size()
323         tp_rank = get_tensor_model_parallel_rank()
324         assert tp_size == 2
325         for name, loaded_weight in hf_model_weights_iterator(
326             model_name_or_path,
327             cache_dir,
328             load_format,
329             revision,
330             fall_back_to_pt=False):
331             if "rotary_emb.inv_freq" in name:
332                 continue
333             for (param_name, weight_name, shard_id) in stacked_params_mapping:
334                 if weight_name not in name:
335                     continue
336                 name = name.replace(weight_name, param_name)
337                 # Skip loading extra bias for GPTQ models.
338                 if name.endswith(".bias") and name not in params_dict:
339                     continue
340                 param = params_dict[name]
341                 weight_loader = param.weight_loader
342                 weight_loader(param, loaded_weight, shard_id)
343                 break
344             else:
345                 # Skip loading extra bias for GPTQ models.
346                 if name.endswith(".bias") and name not in params_dict:
347                     continue
348                 # Skip experts that are not assigned to this worker.
349                 if ("block_sparse_moe.experts." in name): # like
350                     "model.layers.31.block_sparse_moe.experts.7.w2.weight"
351                     substrings = name.split('.')
352                     layer_num = int(substrings[2])
353                     expert_num = int(substrings[5])
354                     if ("w1" in name or "w3" in name): # 列切

```

```

354         loaded_weight = loaded_weight.T # (hidden_size,
        inter_size)
355         shard_size = loaded_weight.shape[-1] // tp_size
356         loaded_weight_slice = loaded_weight[:, tp_rank *
shard_size : (tp_rank + 1) * shard_size] # (hidden_size, inter_size //
tp_size)
357         coresponded_param_name = f"model.layers.
{layer_num}.block_sparse_moe.moe_up_proj_weight"
358         param = params_dict[coresponded_param_name] #
(num_expert, hidden_size, 2 * inter_size // tp_size)
359         if ("w1" in name):
360             param_slice = param[expert_num, :, shard_size : ]
361             assert loaded_weight_slice.shape ==
param_slice.shape
362             param_slice.data.copy_(loaded_weight_slice)
363             continue
364         else: # w3
365             param_slice = param[expert_num, :, : shard_size]
366             assert loaded_weight_slice.shape ==
param_slice.shape
367             param_slice.data.copy_(loaded_weight_slice)
368             continue
369         elif ("w2" in name): # 行切
370             loaded_weight = loaded_weight.T # (inter_size,
hidden_size)
371             shard_size = loaded_weight.shape[0] // tp_size
372             loaded_weight_slice = loaded_weight[tp_rank *
shard_size : (tp_rank + 1) * shard_size, :] # (inter_size // tp_size,
hidden_size)
373             coresponded_param_name = f"model.layers.
{layer_num}.block_sparse_moe.moe_down_proj_weight"
374             param = params_dict[coresponded_param_name] #
(num_expert, inter_size // tp_size, hidden_size, )
375             param_slice = param[expert_num]
376             assert loaded_weight_slice.shape == param_slice.shape
377             param_slice.data.copy_(loaded_weight_slice)
378             continue
379         else:
380             raise ValueError(f"invalid params: {name}")
381             continue
382         else:
383             param = params_dict[name]
384             weight_loader = getattr(param, "weight_loader",
default_weight_loader)
385             weight_loader(param, loaded_weight)
386

```

