

MoE

1. MoE

1.1 什么是 MoE

MoE，全称 Mixture of Experts，混合专家模型。

什么是 MoE？从 Transformer 模型的角度来说，MoE 包含两个要素：**稀疏 MoE 层和门控网络**。

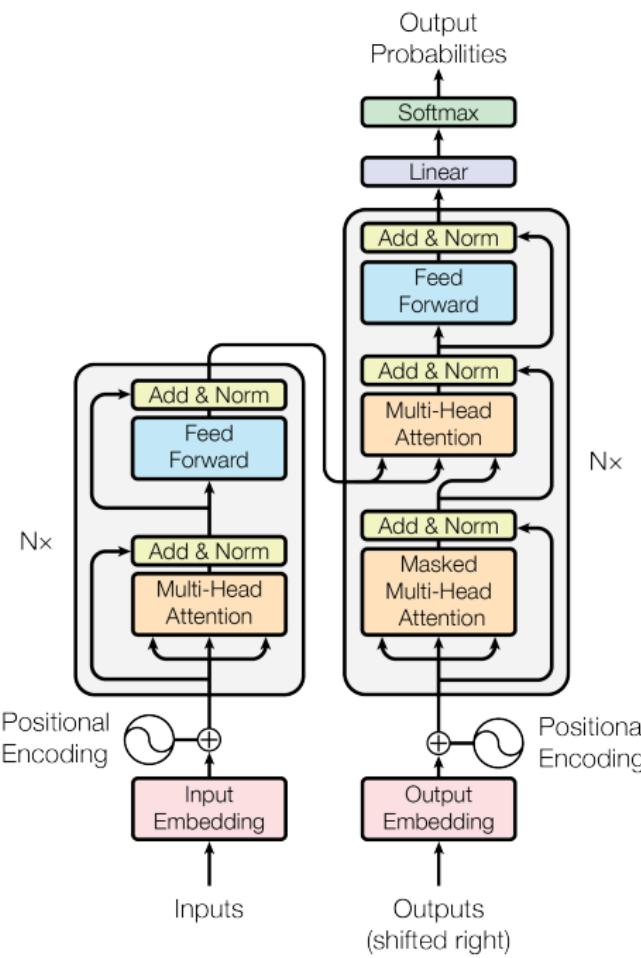
简单理解，**稀疏 MoE 层就是专家层，由 N 个专家网络构成，门控网络用于将 token 分发给不同的专家**。在 MoE 模型运行过程中，在每个时间步由门控网络选择 n 个 ($n < N$) 专家执行计算，也就是每个时间步仅使用一部分参数执行计算。

我们之前学习的 Transformer 系列模型都是**稠密模型**，在处理数据时，**所有的参数都会对所有输入数据进行处理**。

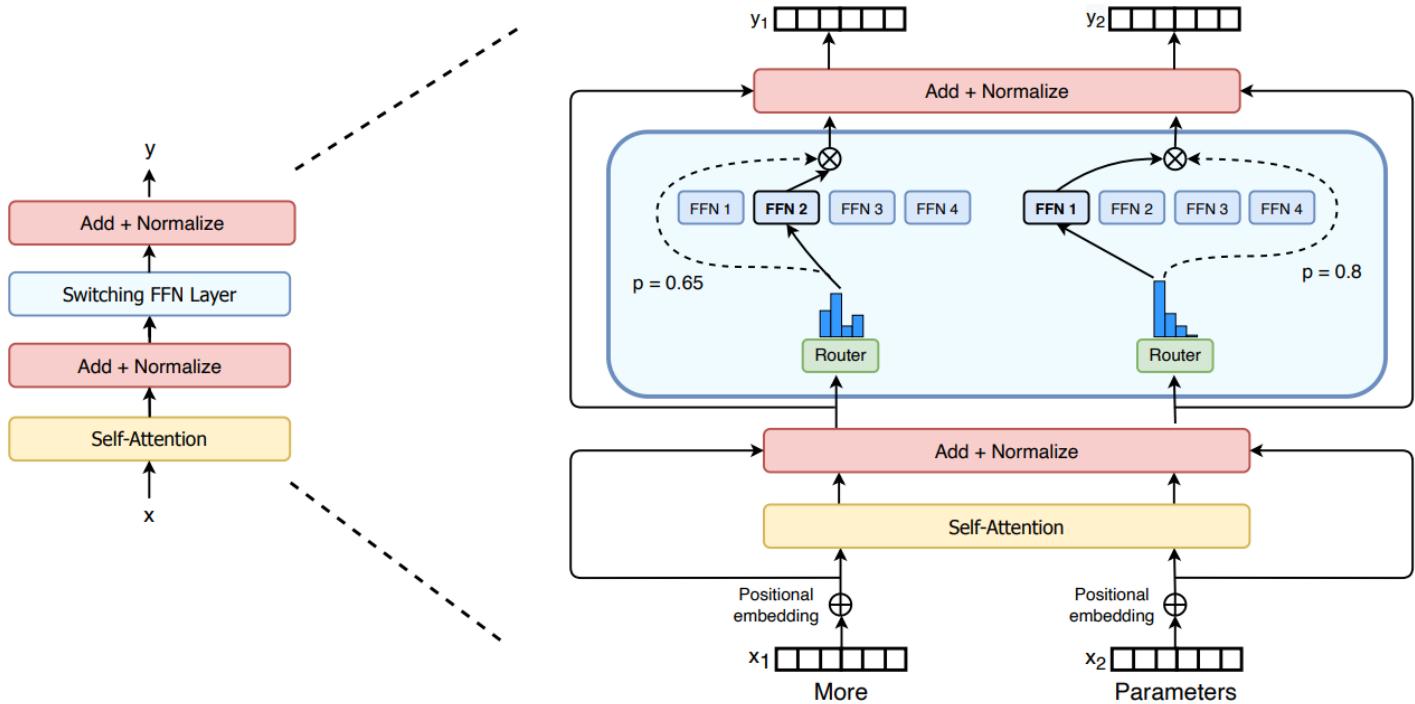
稀疏性的概念采用了**条件计算**的思想。相比于稠密模型，**稀疏性允许我们仅针对整个系统的某些特定部分执行计算**。这意味着**并非所有参数都会在处理每个输入时被激活或使用，而是根据输入的特定特征或需求，只有部分参数集合被调用和运行**。

1.2 MoE 的结构

Transformer 架构：



典型MoE架构：



MoE 基于 Transformer 架构，主要由两部分组成：

- **门控网络/路由**：这个部分用于决定哪些 token 被发送到哪个专家。token 的路由方式是 MoE 使用中的一个关键点，因为路由器由学习的参数组成，并且与网络的其他部分一同进行预训练。

- **稀疏 MoE 层**: 这些层代替了传统 Transformer 模型中的前馈网络 (FFN) 层。MoE 层包含若干“专家”(例如 8 个)，每个专家本身是一个独立的神经网络。在实际应用中，这些专家通常是前馈网络 (FFN)，但它们也可以是更复杂的网络结构。

1.3 不同专家传递了什么信息

MoE 通过门控网络，在每个时间步，将 input 中的 token 分发给不同的专家。

1.4 专家是否真的实现了“分工”

以 ST MoE 为例，在 ST MoE 的研究中，发现有专门处理标点、连词和冠词、动词、视觉描述、专有名词、计数和数字的专家，同样也有对特殊 token (例如 <extra id x>) 进行操作的专家。

Expert specialization	Expert position	Routed tokens
Sentinel tokens	Layer 1	been <extra_id_4><extra_id_7>floral to <extra_id_10><extra_id_12><extra_id_15> <extra_id_17><extra_id_18><extra_id_19>...
	Layer 4	<extra_id_0><extra_id_1><extra_id_2> <extra_id_4><extra_id_6><extra_id_7> <extra_id_12><extra_id_13><extra_id_14>...
	Layer 6	<extra_id_0><extra_id_4><extra_id_5> <extra_id_6><extra_id_7><extra_id_14> <extra_id_16><extra_id_17><extra_id_18>...
Punctuation	Layer 2	, , , , , , - , , , , .)
	Layer 6	, , , , : : , & , & & ? & - , , ? , , , <extra_id_27>
Conjunctions and articles	Layer 3	The the the the the the the the the the the The the the the
	Layer 6	a and and and and and and or and a and . the the if ? a designed does been is not
Verbs	Layer 1	died falling identified fell closed left posted lost felt left said read miss place struggling falling signed died falling designed based disagree submitted develop
Visual descriptions <i>color, spatial position</i>	Layer 0	her over her know dark upper dark outer center upper blue inner yellow raw mama bright bright over open your dark blue
Proper names	Layer 1	A Mart Gr Mart Kent Med Cor Tri Ca Mart R Mart Lorraine Colin Ken Sam Ken Gr Angel A Dou Now Ga GT Q Ga C Ko C Ko Ga G
Counting and numbers <i>written and numerical forms</i>	Layer 1	after 37 19. 6. 27 I I Seven 25 4, 54 I two dead we Some 2012 who we few lower each

此外，研究者们还对 ST MoE 进行了多语言训练，尽管预期每个专家处理一种特定语言，但实际上并非如此，由于令牌路由和负载均衡的机制，没有任何专家被特定配置以专门处理某一特定语言。

1.5 专家根据什么进行“分工”

以 Mixtral 8x7B 为例，下图显示了不同领域（Python 代码、数学和英语）的文本示例，其中每个 token 都用与其所选专家相对应的背景色突出显示：

Layer 0	Layer 15	Layer 31
<pre>class Moelayer(nn.Module): def __init__(self, experts: List[nn.Module], gate_fn): super().__init__() assert len(experts) > 0 self.experts = nn.ModuleList(experts) self.gate = gate_fn self.args = moe_args def forward(self, inputs: torch.Tensor): inputs_squashed = inputs.view(-1, inputs.size(1)) gate_logits = self.gate(inputs_squashed) weights, selected_experts = torch.topk(gate_logits, self.args.num_experts_per_step) weights = nn.functional.softmax(weights, dim=1, dtype=torch.float,).type_as(inputs) results = torch.zeros_like(inputs_squashed) for i, expert in enumerate(self.experts): batch_idx, nth_expert = torch.where(s results[batch_idx] += weights[batch_idx] * expert(inputs_squashed[batch_idx])) return results.view_as(inputs)</pre> <pre>Question: Solve -42*r + 27*c = -1167 and 130*r Answer: 4</pre> <pre>Question: Calculate -841880142.544 + 411127. Answer: -841469015.544</pre> <pre>Question: Let x(g) = 9*g + 1. Let q(c) = 2*c + Answer: 54*a - 30</pre> <pre>A model airplane flies slower when flying into the wind and faster with wind at its back. When launching right angles to the wind, a cross wind, its ground speed compared with flying in still air is (A) the same (B) greater (C) less (D) either greater or less depending on wind speed</pre>	<pre>class Moelayer(nn.Module): def __init__(self, experts: List[nn.Module], gate_fn): super().__init__() assert len(experts) > 0 self.experts = nn.ModuleList(experts) self.gate = gate_fn self.args = moe_args def forward(self, inputs: torch.Tensor): inputs_squashed = inputs.view(-1, inputs.size(1)) gate_logits = self.gate(inputs_squashed) weights, selected_experts = torch.topk(gate_logits, self.args.num_experts_per_step) weights = nn.functional.softmax(weights, dim=1, dtype=torch.float,).type_as(inputs) results = torch.zeros_like(inputs_squashed) for i, expert in enumerate(self.experts): batch_idx, nth_expert = torch.where(s results[batch_idx] += weights[batch_idx] * expert(inputs_squashed[batch_idx])) return results.view_as(inputs)</pre> <pre>Question: Solve -42*r + 27*c = -1167 and 130*r Answer: 4</pre> <pre>Question: Calculate -841880142.544 + 411127. Answer: -841469015.544</pre> <pre>Question: Let x(g) = 9*g + 1. Let q(c) = 2*c + Answer: 54*a - 30</pre> <pre>A model airplane flies slower when flying into the wind and faster with wind at its back. When launching right angles to the wind, a cross wind, its ground speed compared with flying in still air is (A) the same (B) greater (C) less (D) either greater or less depending on wind speed</pre>	<pre>class Moelayer(nn.Module): def __init__(self, experts: List[nn.Module], gate_fn): super().__init__() assert len(experts) > 0 self.experts = nn.ModuleList(experts) self.gate = gate_fn self.args = moe_args def forward(self, inputs: torch.Tensor): inputs_squashed = inputs.view(-1, inputs.size(1)) gate_logits = self.gate(inputs_squashed) weights, selected_experts = torch.topk(gate_logits, self.args.num_experts_per_step) weights = nn.functional.softmax(weights, dim=1, dtype=torch.float,).type_as(inputs) results = torch.zeros_like(inputs_squashed) for i, expert in enumerate(self.experts): batch_idx, nth_expert = torch.where(s results[batch_idx] += weights[batch_idx] * expert(inputs_squashed[batch_idx])) return results.view_as(inputs)</pre> <pre>Question: Solve -42*r + 27*c = -1167 and 130*r Answer: 4</pre> <pre>Question: Calculate -841880142.544 + 411127. Answer: -841469015.544</pre> <pre>Question: Let x(g) = 9*g + 1. Let q(c) = 2*c + Answer: 54*a - 30</pre> <pre>A model airplane flies slower when flying into the wind and faster with wind at its back. When launching right angles to the wind, a cross wind, its ground speed compared with flying in still air is (A) the same (B) greater (C) less (D) either greater or less depending on wind speed</pre>

在 Mixtral 8x7B 中，**专家的选择似乎更符合语法规而非领域**，尤其是在初始层和最终层。

目前，“**专家专门化**”是MoE技术的重要研究方向，在陈丹琦团队今年5月份发布的MoE架构Lory中，通过在训练数据中分组相似文档，鼓励“专家专门化”。

2. MoE的代码实现原理

2.1 门控网络

门控网络，也称为路由，确定哪个专家网络接收来自多头注意力的 token 的输出。

2.1.1 基础门控网络

举个例子解释路由的机制，假设有 4 个专家，token 需要被路由到前 2 个专家中。首先需要通过线性层将 token 输入到门控网络中。该层将对应于 [batch_size, max_seq, model_dim] 的输入张量从 [2, 4, 32] 维度，投影到对应于 [batch_size, max_seq, num_expert] 的新形状：[2, 4, 4]。其中 n_embed 是输入的通道维度，num_experts 是专家网络的计数。

```
1 # 定义模型超参数
2 batch_size=2
3 context_length=4
```

```

4 n_embed=32
5 # 创建随机向量表示MHA输出
6 mh_output = torch.randn(batch_size, context_length, n_embed)
7 # 定义线性变换层
8 topkgate_linear = nn.Linear(n_embed, numExperts)
9 # [batch_size, max_seq, num_expert]
10 logits = topkgate_linear(mh_output)
11 # 获取TopK的Logits和对应的Indices
12 top_k_logits, top_k_indices = logits.topk(top_k, dim=-1)
13 top_k_logits, top_k_indices

```

```

1 # output
2 (tensor([[[ 0.0246, -0.0190],
3          [ 0.1991,  0.1513],
4          [ 0.9749,  0.7185],
5          [ 0.4406, -0.8357]],
6         [[ 0.6206, -0.0503],
7          [ 0.8635,  0.3784],
8          [ 0.6828,  0.5972],
9          [ 0.4743,  0.3420]]])
10 tensor([[[[2, 3],
11           [2, 1],
12           [3, 1],
13           [2, 1]],
14           [[[0, 2],
15             [0, 3],
16             [3, 2],
17             [3, 0]]]])

```

获得稀疏门控的输出后，用 `-inf`（负无穷）填充其余部分，再使用 `softmax` 激活函数，目的是经过 `softmax` 处理时负无穷会被映射至零，而最大的前两个值会更加突出，且和为 1。要求和为 1 是为了对专家输出的内容进行加权。

```

1 # 使用 -inf (负无穷) 填充其余部分，恢复 num_experts 维度
2 zeros = torch.full_like(logits, float('-inf'))
3 sparse_logits = zeros.scatter(-1, top_k_indices, top_k_logits)
4 sparse_logits

```

```

1 #output
2 tensor([[[ -inf,    -inf,  0.0246, -0.0190],
3          [ -inf,   0.1513,  0.1991,    -inf],
4          [ -inf,    -inf,  0.9749,  0.7185],
5          [ -inf,    -inf,  0.4406, -0.8357]])

```

```
4      [-inf,  0.7185,    -inf,  0.9749],
5      [-inf, -0.8357,  0.4406,    -inf]], 
6      [[ 0.6206,     -inf, -0.0503,    -inf],
7      [ 0.8635,     -inf,     -inf,  0.3784],
8      [-inf,     -inf,  0.5972,  0.6828],
9      [ 0.3420,     -inf,     -inf,  0.4743]]], grad_fn=<ScatterBackward0>)
```

进行softmax处理，Top-2的权重被归一化，其余专家对应权重被归0：

```
1 gating_output= F.softmax(sparse_logits, dim=-1)
2 gating_output
```

```
1 #output
2 tensor([[ [0.0000, 0.0000, 0.5109, 0.4891],
3           [0.0000, 0.4881, 0.5119, 0.0000],
4           [0.0000, 0.4362, 0.0000, 0.5638],
5           [0.0000, 0.2182, 0.7818, 0.0000]],
6           [[0.6617, 0.0000, 0.3383, 0.0000],
7           [0.6190, 0.0000, 0.0000, 0.3810],
8           [0.0000, 0.0000, 0.4786, 0.5214],
9           [0.4670, 0.0000, 0.0000, 0.5330]]], grad_fn=<SoftmaxBackward0>)
```

综上，基础的门控机制如下：

```
1 # First define the top k router module
2 class TopkRouter(nn.Module):
3     def __init__(self, n_embed, num_experts, top_k):
4         super(TopkRouter, self).__init__()
5         self.top_k = top_k
6         self.linear = nn.Linear(n_embed, num_experts)
7     def forward(self, mh_output):
8         # [batch_size, max_seq, model_dim] -> [batch_size, max_seq, num_expert]
9         logits = self.linear(mh_output)
10        # 获取 Top-K
11        top_k_logits, indices = logits.topk(self.top_k, dim=-1)
12        # 填充 "-inf"
13        zeros = torch.full_like(logits, float('-inf'))
14        sparse_logits = zeros.scatter(-1, indices, top_k_logits)
15        # 生成权重
16        router_output = F.softmax(sparse_logits, dim=-1)
17        return router_output, indices
```

测试结果：

```
1 batch_size=2
2 context_length=4
3 n_embed=32
4 mh_output = torch.randn(batch_size, context_length, n_embed)
5 # 示例输出
6 inputtop_k_gate = TopkRouter(n_embed, num_experts, top_k)
7 gating_output, indices = top_k_gate(mh_output)
8 gating_output.shape, gating_output, indices
```

```
1 #output
2 (torch.Size([2, 4, 4]),
3  tensor([[ [0.5284, 0.0000, 0.4716, 0.0000],
4           [0.0000, 0.4592, 0.0000, 0.5408],
5           [0.0000, 0.3529, 0.0000, 0.6471],
6           [0.3948, 0.0000, 0.0000, 0.6052]],
7           [[0.0000, 0.5950, 0.4050, 0.0000],
8           [0.4456, 0.0000, 0.5544, 0.0000],
9           [0.7208, 0.0000, 0.0000, 0.2792],
10          [0.0000, 0.0000, 0.5659, 0.4341]]], grad_fn=<SoftmaxBackward0>),
11 tensor([[[[0, 2],
12           [3, 1],
13           [3, 1],
14           [3, 0]],
15           [[1, 2],
16           [2, 0],
17           [0, 3],
18           [2, 3]]]))
```

2.1.2 带噪声的门控网络

在通常的混合专家模型 (MoE) 训练中，**门控网络往往倾向于主要激活相同的几个专家。这种情况可能会自我加强，因为受欢迎的专家训练得更快，因此它们更容易被选择。**但是，如果所有的令牌都被发送到只有少数几个受欢迎的专家，那么训练效率将会降低。因此，为了负载平衡，从门控的线性层向 logits 激活函数添加标准正态噪声。

```
1 class NoisyTopkRouter(nn.Module):
2     def __init__(self, n_embed, num_experts, top_k):
3         super(NoisyTopkRouter, self).__init__()
```

```

4     self.top_k = top_k
5     self.topkroute_linear = nn.Linear(n_embed, num_experts)
6     self.noise_linear = nn.Linear(n_embed, num_experts)
7     def forward(self, mh_output):
8         # dim -> num_experts
9         logits = self.topkroute_linear(mh_output)
10        # 生成noise
11        noise_logits = self.noise_linear(mh_output)
12        # 生成高斯白噪声
13        noise = torch.randn_like(logits)*F.softplus(noise_logits)
14        # 加入高斯白噪声
15        noisy_logits = logits + noise
16        # 获取top_k
17        top_k_logits, indices = noisy_logits.topk(self.top_k, dim=-1)
18        # 填充 -inf
19        zeros = torch.full_like(noisy_logits, float('-inf'))
20        sparse_logits = zeros.scatter(-1, indices, top_k_logits)
21        # 权重结果
22        router_output = F.softmax(sparse_logits, dim=-1)
23        return router_output, indices

```

测试结果：

```

1 batch_size=2
2 context_length=4
3 n_embed=32
4 num_experts=8
5
6 mh_output = torch.randn(2, 4, n_embed)
7 # 示例输出
8 noisy_top_k_gate = NoisyTopkRouter(n_embed, num_experts, top_k)
9 gating_output, indices = noisy_top_k_gate(mh_output)
10 gating_output.shape, gating_output, indices

```

```

1 #output
2 (torch.Size([2, 4, 8]),
3 tensor([[ [0.4181, 0.0000, 0.5819, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000],
4           [0.4693, 0.5307, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000],
5           [0.0000, 0.4985, 0.5015, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000],
```

```
6 [0.0000, 0.0000, 0.0000, 0.2641, 0.0000, 0.7359, 0.0000, 0.0000]],  
7 [[0.0000, 0.0000, 0.0000, 0.6301, 0.0000, 0.3699, 0.0000, 0.0000],  
8 [0.0000, 0.0000, 0.0000, 0.4766, 0.0000, 0.0000, 0.0000, 0.5234],  
9 [0.0000, 0.0000, 0.0000, 0.6815, 0.0000, 0.0000, 0.3185, 0.0000],  
10 [0.4482, 0.5518, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000]]],  
grad_fn=<SoftmaxBackward0>),  
11 tensor([[2, 0],  
12 [1, 0],  
13 [2, 1],  
14 [5, 3],  
15 [[3, 5],  
16 [7, 3],  
17 [3, 6],  
18 [1, 0]]]))
```

除了为路由增加噪声，还可以通过随机路由、专家容量、辅助损失等策略实现 token 负载均衡。

2.2 混合专家模块

2.2.1 专家定义

每个专家都是一个**简单的多层感知器**。

```
1 class Expert(nn.Module):  
2     """ An MLP is a simple linear layer followed by a non-linearity i.e. each  
3     Expert """  
4     def __init__(self, n_embd):  
5         super().__init__()  
6         # 创建多层容器  
7         self.net = nn.Sequential(  
8             nn.Linear(n_embd, 4 * n_embd),  
9             nn.ReLU(),  
10            nn.Linear(4 * n_embd, n_embd),  
11            nn.Dropout(dropout),  
12        )  
13    def forward(self, x):  
14        return self.net(x)
```

2.2.2 混合专家模块

在获得门控网络的输出结果之后，对于给定的 token，将前 k 个值选择性地与来自相应的前 k 个专家的输出相乘。这种选择性乘法的结果是一个加权和，该加权和构成 SparseMoe 模块的输出。

这个过程的关键和难点是避免不必要的乘法运算，只为前 k 名专家进行正向转播。为每个专家执行前向传播将破坏使用稀疏 MoE 的目的，因为这个过程将不再是稀疏的。

```
1 class SparseMoE(nn.Module):
2     def __init__(self, n_embed, num_experts, top_k):
3         super(SparseMoE, self).__init__()
4         self.router = NoisyTopkRouter(n_embed, num_experts, top_k)
5         self.experts = nn.ModuleList([Expert(n_embed) for _ in
6             range(num_experts)])
7         self.top_k = top_k
8
9     def forward(self, x):
10        # x [4, 8, 16]
11        # 每个token被分发给2个专家，所以 indices [4, 8, 2]
12        # gating_output 是每个token选定专家的加权权重 [4, 8, 8]
13        gating_output, indices = self.router(x)
14        final_output = torch.zeros_like(x)
15        # flat_x [32, 16]
16        flat_x = x.view(-1, x.size(-1))
17        # flat_gating_output [32, 8]
18        flat_gating_output = gating_output.view(-1, gating_output.size(-1))
19        # 执行专家处理逻辑
20        for i, expert in enumerate(self.experts):
21            # indices与当前expert的i一致的位置设置为True，代表该token会被分发给当前
22            # expert
23            # expert_mask 代表输入的所有token是否要经过该expert处理
24            # expert_mask [4, 8]
25            expert_mask = (indices == i).any(dim=-1)
26            # 将 Mask 拉直为扁平向量 ([4, 8] -> [32])
27            flat_mask = expert_mask.view(-1)
28            # 如果有token需要被expert处理
29            if flat_mask.any():
30                # 取要被该expert处理的token
31                # expert_input [N, 16]
32                expert_input = flat_x[flat_mask]
33                # 执行expert处理逻辑
34                # expert_output [N, 16]
35                expert_output = expert(expert_input)
36                # 获取每个token的权重
37                # gating_scores [N, 1]
38                gating_scores = flat_gating_output[flat_mask, i].unsqueeze(1)
39                # 计算加权值
40                # weighted_output [N, 16]
```

```
39         weighted_output = expert_output * gating_scores
40         # 将当前加权值加入总值中
41         # final_output [4, 8, 16]
42         final_output[expert_mask] += weighted_output.squeeze(1)
43         print('*****')
44     return final_output
```

代码测试：

```
1 import torchimport torch.nn as nn
2 batch_size=4
3 context_length=8
4 n_embed=16
5
6 mh_output = torch.randn(batch_size, context_length, n_embed)
7 # 创建MoE
8 outputsparse_moe = SparseMoE(n_embed, num_experts, top_k)
9 # 获取MoE结果
10 final_output = sparse_moe(mh_output)
11 print("Shape of the final output:", final_output.shape)
```

```
1 Shape of the final output: torch.Size([4, 8, 16])
```

3. 总结

目前，MoE 技术在大模型中的应用还面临一些挑战，如模型训练复杂、权重分配不平衡、分布式训练成本高等。但随着技术的不断进步，这些问题将逐渐得到解决。未来，MoE 技术有望在更多领域得到应用，为人工智能的发展带来新的机遇。