

# Deepseek-v3 Op List

## vLLM

共计是99个算子，需要算子实现的是54个

需要实现的和llama2重复的有39个标为红色

独有15个，标为蓝色

加粗字体的是和llama2 相同的算子（包含需要实现和不需要实现）

Operation Name	op 文档	Need Impl	Init	Prefill
<b>_C::fused_add_rms_norm</b>	OP的开发文档	X	X	
<b>_C::rms_norm</b>	OP 的开发文档	X	X	
<b>_C::silu_and_mul</b>	OP 的开发文档	X	X	
<b>_C_cache_ops::concat_and_cache_mla</b>	OP 的开发文档	X		
<b>_moe_C::moe_align_block_size</b>	OP 的开发文档	X	X	X
<b>_moe_C::moe_sum</b>	OP 的开发文档	X	X	X
<b>_vllm_fa2_C::varlen_fwd</b>	OP 的开发文档	X	X	X
<b>aten::_log_softmax</b>	OP的开发文档	X	X	
<b>aten::_reshape_alias</b>	OP 的开发文档	X	X	X
<b>aten::_softmax</b>	OP的开发文档	X	X	
<b>aten::_to_copy</b>	<b>CompositeExplicitAutograd</b>		X	X
<b>aten::_unsafe_view</b>	CompositeExplicitAutograd		X	X
<b>aten::add.Tensor</b>	DONE	X	X	X
<b>aten::alias</b>	<b>CompositeExplicitAutograd</b>		X	
<b>aten::arange</b>	OP的开发文档	X	X	
<b>aten::arange.start_out</b>	OP的开发文档	X	X	
<b>aten::arange.start_step</b>	OP 的开发文档	X	X	
<b>aten::argmax</b>	OP的开发文档	X	X	X
<b>aten::as_strided</b>	OP的开发文档	X	X	X
<b>aten::bitwise_not</b>	OP 的开发文档	X		X
<b>aten::bmm</b>	OP 的开发文档	X	X	
<b>aten::cat</b>	OP的开发文档	X	X	X
<b>aten::chunk</b>	CompositeImplicitAutograd		X	X
<b>aten::clamp</b>	OP 的开发文档	X	X	
<b>aten::clone</b>	CompositeExplicitAutograd			
<b>aten::constant_pad_nd</b>	CompositeExplicitAutograd		X	X
<b>aten::contiguous</b>	CompositeImplicitAutograd		X	X
<b>aten::copy_</b>	<b>CompositeExplicitAutograd</b>		X	X
<b>aten::cos</b>	OP的开发文档	X	X	
<b>aten::cumsum</b>	OP的开发文档	X	X	
<b>aten::detach</b>	<b>CompositeExplicitAutograd</b>		X	
<b>aten::detach_</b>	<b>CompositeExplicitAutograd</b>		X	X
<b>aten::div.Tensor</b>	OP的开发文档	X	X	
<b>aten::div_.Tensor</b>	OP 的开发文档		X	X
<b>aten::embedding</b>	<b>CompositeExplicitAutograd</b>		X	X
<b>aten::empty.memory_format</b>	DONE	X	X	

<b>aten::empty_like</b>	<b>CompositeExplicitAutograd</b>		X	X
<b>aten::empty_strided</b>	<b>DONE</b>	X	X	X
aten::expand	CompositeExplicitAutograd			
<b>aten::exponential_</b>	<a href="#">OP的开发文档</a>	X	X	
<b>aten::fill_Scalar</b>	<a href="#">OP的开发文档</a>	X	X	X
<b>aten::fill_Tensor</b>	<a href="#">OP的开发文档</a>	X	X	
<b>aten::flatten.using_ints</b>	<b>CompositeImplicitAutograd</b>			
aten::full	CompositeExplicitAutograd		X	
<b>aten::gather</b>	<a href="#">OP的开发文档</a>	X	X	
<b>aten::index.Tensor</b>	<a href="#">OP的开发文档</a>	X	X	
<b>aten::index_select</b>	<a href="#">OP的开发文档</a>	X	X	
<b>aten::le.Tensor</b>	<a href="#">OP 的开发文档</a>	X	X	
<b>aten::lift_fresh</b>	<b>CompositeExplicitAutograd</b>		X	X
<b>aten::linear</b>	<b>CompositeImplicitAutograd</b>		X	
<b>aten::log_softmax.int</b>	<b>CompositeImplicitAutograd</b>		X	
<b>aten::lt.Tensor</b>	<a href="#">OP 的开发文档</a>	X	X	
aten::masked_fill.Scalar	CompositeExplicitAutograd			
<b>aten::masked_fill_Scalar</b>	<a href="#">OP的开发文档</a>	X	X	X
<b>aten::matmul</b>	<b>CompositeImplicitAutograd</b>		X	X
aten::max.dim	<a href="#">OP 的开发文档</a>	X		X
<b>aten::mm</b>	<a href="#">OP的开发文档</a>	X	X	
<b>aten::mul.Tensor</b>	<a href="#">OP的开发文档</a>	X	X	X
aten::narrow	CompositeImplicitAutograd		X	X
aten::neg	<b>DONE</b>	X	X	X
<b>aten::ones</b>	<b>CompositeExplicitAutograd</b>		X	
aten::pad	CompositeImplicitAutograd		X	X
<b>aten::permute</b>	<b>CompositeExplicitAutograd</b>		X	
<b>aten::pow.Scalar</b>	<a href="#">OP 的开发文档</a>	X	X	
<b>aten::pow.Tensor_Tensor_out</b>	<a href="#">OP 的开发文档</a>	X	X	
<b>aten::reciprocal</b>	<a href="#">OP 的开发文档</a>	X	X	
aten::repeat_interleave.self_int	<a href="#">OP 的开发文档</a>	X	X	X
<b>aten::reshape</b>	<b>CompositeImplicitAutograd</b>		X	X
<b>aten::resize_</b>	<a href="#">OP 的开发文档</a>	X	X	X
aten::result_type.Scalar	CompositeImplicitAutograd		X	
aten::result_type.Scalar_Tensor	CompositeImplicitAutograd		X	
<b>aten::rsub.Scalar</b>	<a href="#">OP 的开发文档</a>	X	X	
aten::scalar_tensor	CompositeExplicitAutograd		X	

aten::scatter_.src	OP 的开发文档		X	
aten::scatter_.value	scatter_.value			
aten::select.int	CompositeExplicitAutograd		X	X
aten::sin	OP 的开发文档	X	X	
aten::slice.Tensor	CompositeExplicitAutograd		X	X
aten::softmax.int	CompositeImplicitAutograd		X	X
aten::sort	OP 的开发文档	X	X	
aten::sort.stable	OP 的开发文档	X	X	
aten::split.Tensor	CompositeExplicitAutograd		X	X
aten::split_with_sizes	CompositeExplicitAutograd		X	
aten::squeeze.dim	CompositeExplicitAutograd			X
aten::stack	CompositeExplicitAutograd		X	X
aten::sub.Tensor	OP 的开发文档	X	X	X
aten::sum.IntList_out	OP 的开发文档	X	X	X
aten::t	CompositeExplicitAutograd		X	
aten::to.dtype	CompositeImplicitAutograd		X	
aten::to.dtype_layout	CompositeImplicitAutograd		X	
aten::topk	OP 的开发文档	X		
aten::transpose.int	CompositeExplicitAutograd		X	
aten::unsqueeze	CompositeExplicitAutograd		X	
aten::view	OP 的开发文档	X	X	X
aten::zero_	OP 的开发文档	X	X	X
aten::zeros	CompositeExplicitAutograd		X	X
aten::zeros_like	CompositeExplicitAutograd			
vllm::inplace_fused_experts	OP 的开发文档	X	X	X
vllm::unified_attention	OP 的开发文档	X	X	

## HF transformers

*attn\_implementation = "eager"* 共计是87个算子，需要算子实现的是55个 (Need Impl)

与llama2 op list 相比，需要算子实现的重复op有39个标为红色

Deepseek-v3 独有有16个 标为蓝色

*attn\_implementation="flash\_attention\_2"* 会有一个额外的需要impl算子

**flash\_attn::\_flash\_attn\_forward**

Operation Name	op 文档	Need Impl	Init	Prefil	Prefill F	Prefill
aten::_index_put_impl_	OP的开发文档	X			X	
aten::_local_scalar_dense	OP的开发文档	X			X	
aten::_softmax	OP的开发文档	X			X	
aten::add.Tensor	OP的开发文档	X			X	X
aten::any	OP的开发文档	X			X	
aten::any.dim	OP的开发文档	X			X	X
aten::arange	OP的开发文档	X			X	X
aten::arange.start_out	OP的开发文档	X			X	X
aten::arange.start_step	OP的开发文档	X			X	
aten::argmax	OP的开发文档	X				X
aten::as_strided	OP的开发文档	X			X	X
aten::bitwise_and.Tensor	OP的开发文档	X				X
aten::bitwise_not	OP的开发文档	X			X	X
aten::bitwise_or.Tensor	OP的开发文档	X				X
aten::bmm	OP的开发文档	X			X	
aten::cat	OP的开发文档	X			X	X
aten::cos	OP的开发文档	X			X	
aten::cumsum	OP的开发文档	X			X	
aten::div.Tensor	OP的开发文档	X			X	
aten::empty.memory_format	DONE	X			X	X
aten::empty_strided	DONE	X	X	X	X	X
aten::eq.Scalar	OP的开发文档	X			X	X
aten::eq.Tensor	OP的开发文档	X			X	X
aten::fill_.Scalar	OP的开发文档	X			X	X
aten::floor_divide	OP的开发文档	X			X	
aten::gather	OP的开发文档	X			X	
aten::index.Tensor	OP的开发文档	X			X	
aten::index_select	OP的开发文档	X			X	
aten::isin.Tensor_Tensor	OP的开发文档	X			X	X
aten::lt.Scalar	OP的开发文档	X			X	
aten::lt.Tensor	OP的开发文档	X			X	
aten::masked_fill_.Scalar	OP的开发文档	X			X	
aten::max	OP的开发文档	X				X
aten::mean.dim	OP的开发文档	X			X	
aten::mm	OP的开发文档	X			X	
aten::mul.Tensor	OP的开发文档	X			X	X

Operation Name	OP 文档	Need Impl	Init	Prefill Prepare	Prefill Forward	Prefill Postpr
aten::neg	OP 的开发文档	X			X	
aten::pow.Scalar	OP 的开发文档	X			X	
aten::pow.Tensor_Scalar	OP 的开发文档	X			X	
aten::pow.Tensor_Tensor_out	OP 的开发文档	X			X	
aten::reciprocal	OP 的开发文档	X			X	
aten::resize_	OP 的开发文档	X			X	X
aten::rsqrt	OP 的开发文档	X			X	
aten::rsub.Scalar	OP 的开发文档	X			X	X
aten::sigmoid	OP 的开发文档	X			X	
aten::silu	OP 的开发文档	X			X	
aten::sin	OP 的开发文档	X			X	
aten::sort	OP 的开发文档	X			X	
aten::sort.stable	OP 的开发文档	X			X	
aten::sub.Tensor	OP 的开发文档	X			X	X
aten::sum.dim_IntList	OP 的开发文档	X			X	
aten::topk	OP 的开发文档	X			X	
aten::view	DONE	X			X	X
aten::where.self	OP 的开发文档	X				X
aten::zero_	OP 的开发文档	X			X	

llama2 op独有的为三个

Operation Name	OP 文档	Need Impl	Init	Prefill Prepare	Prefill Forward	Prefill Postpr
aten::ge.Scalar	OP 的开发文档	X			X	
aten::gt.Tensor	OP 的开发文档	X			X	
aten::triu	OP 的开发文档	X			X	

## FP8 支持需求

通过DeepSeek开源的weights (<https://huggingface.co/deepseek-ai/DeepSeek-V3-Base/tree/main>) + 源码 (<https://github.com/deepseek-ai/DeepSeek-V3/blob/main/inference/kernel.py> <https://github.com/deepseek-ai/DeepSeek-V3/blob/main/inference/model.py>) 分析

### Embedding:

- model.embed\_tokens.weight → 使用 BF16, 不需要 FP8

## Multi-Headed Attention Layer (MLA):

- self\_attn.kv\_a\_proj\_with\_mqa.weight → FP8 (Linear -> fp8\_gemm)
- self\_attn.kv\_b\_proj.weight → FP8 (ColumnParallelLinear->Linear ->fp8\_gemm)
- self\_attn.o\_proj.weight → FP8 (RowParallelLinear ->Linear ->fp8\_gemm)
- self\_attn.q\_a\_proj.weight → FP8 (ColumnParallelLinear->Linear ->fp8\_gemm)
- self\_attn.q\_b\_proj.weight → FP8 (ColumnParallelLinear->Linear ->fp8\_gemm)

## FFN:

### • 非 MoE 和 MoE 层:

- mlp.down\_proj.weight → FP8 (ColumnParallelLinear / Linear -> linear->fp8\_gemm)
- mlp.gate\_proj.weight → FP8 (Linear -> linear->fp8\_gemm)
- mlp.up\_proj.weight → FP8 (RowParallelLinear / ColumnParallelLinear -> linear->fp8\_gemm)

涉及到使用FP8数据类型计算的kernel只有**fp8\_gemm\_kernel()** [https://github.com/deepseek-ai/DeepSeek-](https://github.com/deepseek-ai/DeepSeek-V3/blob/a878eada08ea6913f5a2ae80a43afeffdef082ef/inference/kernel.py#L115C5-L115C20)

[V3/blob/a878eada08ea6913f5a2ae80a43afeffdef082ef/inference/kernel.py#L115C5-L115C20](https://github.com/deepseek-ai/DeepSeek-V3/blob/a878eada08ea6913f5a2ae80a43afeffdef082ef/inference/kernel.py#L115C5-L115C20)

在**fp8\_gemm\_kernel**中关键kernel为**tl.dot()**，参考文档：<https://triton-lang.org/main/python-api/generated/triton.language.dot.html>

Deepseek开源的inference代码中的ColumnParallelLinear/RowParallelLinear/Linear 在调用fp8\_gemm\_kernel时，对应在pytorch框架中等价为**matmul(gemm)**，**mv(gemv)**，**bmm**三个算子

此外还需要实现**act\_quant()** **weight\_dequant()**量化和反量化操作<https://github.com/deepseek-ai/DeepSeek-V3/blob/a878eada08ea6913f5a2ae80a43afeffdef082ef/inference/kernel.py>

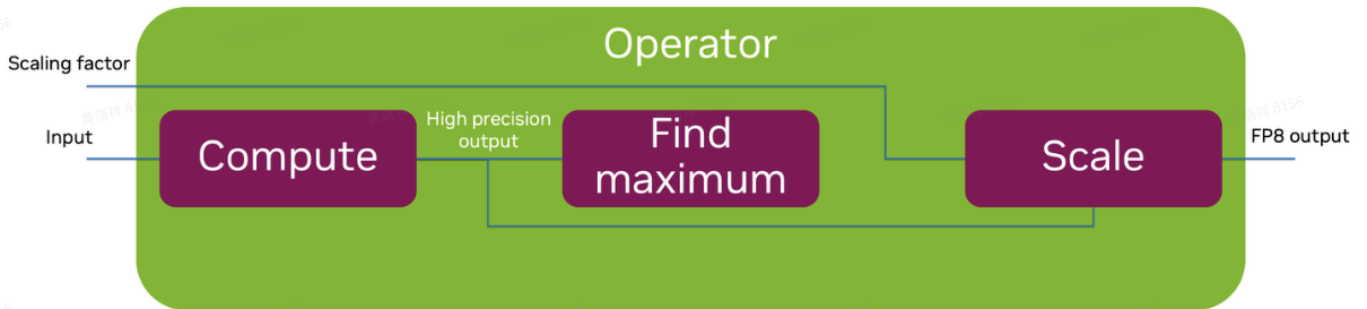
TE中的实现进行了算子融合，会涉及到更多的FP8情况

### 1. FP8 量化与反量化算子

[https://docs.nvidia.com/deeplearning/transformer-engine/user-guide/examples/fp8\\_primer.html#Mixed-precision-training-with-FP8](https://docs.nvidia.com/deeplearning/transformer-engine/user-guide/examples/fp8_primer.html#Mixed-precision-training-with-FP8)

- **Delayed scaling**, 即假设我们提前知道了 scale 值, 那么计算过程就可以在一个 kernel 上完成, 而 amax 的计算和 scale 值的更新与其独立, 不会中断计算进程, 因此这种方式能完全发挥 FP8 的性能。但“提前知道的” scale 值需要额外的空间来记录, 同时会引入一定的误差。

- To overcome that we need to know the scaling factor before seeing the output



对每个 GEMM 算子用到的 tensor 记录一个 amax history 数组, 当我们需要 scale 值时, 就从这个数组中取出最近一段时间窗口内 amax 的最大值, 以此近似现在这个 tensor 的 amax, 并默认用以下方式计算 scale 值:

代码块

```
1 FP8_MAX = maximum_representable_value(fp8_format)
2 new_scaling_factor = (FP8_MAX / amax) / (2 ^ margin)
```

## 2. FP8 Cast 与转置算子/padding结合

[https://github.com/NVIDIA/TransformerEngine/blob/db2aaa9e7765a3bcd1e652e3539d57e85cbbd5bc/transformer\\_engine/common/include/transformer\\_engine/transpose.h#L79](https://github.com/NVIDIA/TransformerEngine/blob/db2aaa9e7765a3bcd1e652e3539d57e85cbbd5bc/transformer_engine/common/include/transformer_engine/transpose.h#L79)

[https://github.com/NVIDIA/TransformerEngine/blob/db2aaa9e7765a3bcd1e652e3539d57e85cbbd5bc/transformer\\_engine/pytorch/module/fp8\\_padding.py#L69](https://github.com/NVIDIA/TransformerEngine/blob/db2aaa9e7765a3bcd1e652e3539d57e85cbbd5bc/transformer_engine/pytorch/module/fp8_padding.py#L69)

## 3. GEMM FP8 (矩阵乘法) 算子

## 4. GEMM+activation+quantize fusion

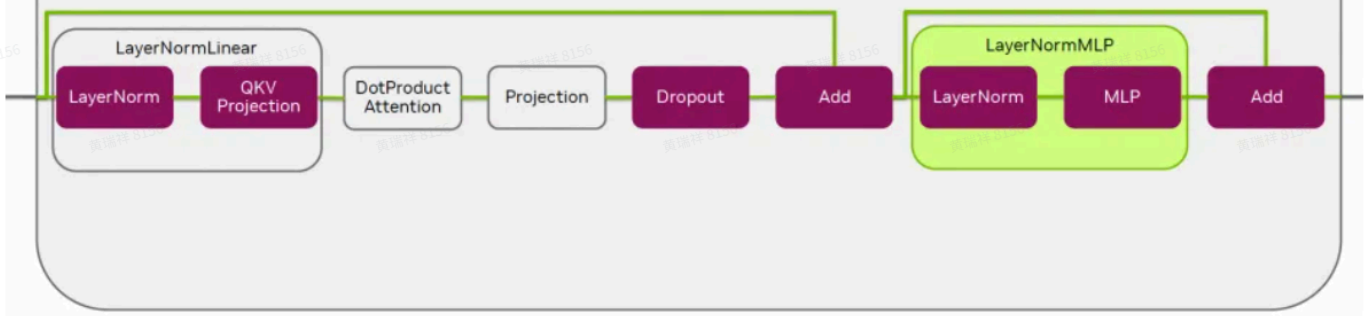
[https://github.com/NVIDIA/TransformerEngine/blob/db2aaa9e7765a3bcd1e652e3539d57e85cbbd5bc/transformer\\_engine/pytorch/module/layernorm\\_mlp.py#L715](https://github.com/NVIDIA/TransformerEngine/blob/db2aaa9e7765a3bcd1e652e3539d57e85cbbd5bc/transformer_engine/pytorch/module/layernorm_mlp.py#L715)

## 5. FP8 融合 Bias-Dropout-Add 算子, 将FP8 casting操作和Bias-Dropout-Add结合

# Transformer Engine

The API

## Transformer Layer



## Attention 推理

`precompute_freqs_cis + apply_rotary_emb` part (ROPE先放在cpu上运行)

1. `arange()` -> `aten::arange.start_out` (MID)
2. `ones()` -> `fill_Scalar` ((MID))
3. `clamp.out` (MID)
4. `data type cast` (MID)
  - a. Data type cast for `direct_copy_kernel_sipu` is not supported now

代码块

- 1 默认input 为bf16, 但是`view_as_complex/view_as_real`计算要求input为float32, 如果不添加`sipu`的type cast, 需要转到cpu上做

## 5. copy size restrict (BUG) @ yunyi WIP

代码块

- 1 `copy size > 8k` 会coredump, ds中会遇到`[32, 16384]` case
- 2 `t = fake_arange(0, seqlen, 1, dtype=torch.bfloat16)`
- 3 `ones = fake_ones(freqs.shape, dtype=torch.bfloat16)`

## 6. `polar.out` 或者 `cos.out + sin.out` (MID)

代码块

```

1 ROPE中位置计算需要
2 # torch.polar 的文档
3 # https://pytorch.org/docs/stable/generated/torch.polar.html
4 # 计算结果是个复数向量
5 # 假设 freqs = [x, y]
6 # 则 freqs_cis = [cos(x) + sin(x)i, cos(y) + sin(y)i]
7
8 freqs_cis = torch.polar(torch.ones_like(freqs), freqs)

```

7. ~~torch.view\_as\_complex (dont need kernel, done)~~

8. ~~torch.view\_as\_real (donot need kernel, done)~~

**RMSNorm (先放在cpu上运行, MID) 目前已有rmsnorm的整体实现, composition的算子优先级降低**

1. aten::pow.Tensor\_Scalar\_out

2. aten::result\_type.Scalar

3. aten::to.dtype

4. aten::mean.out

5. aten::sum.IntList\_out

6. aten::fill\_.Scalar

7. div\_.Scalar + div\_.Tensor

8. aten::rsqrt.out

**Linear**

1. ~~act\_quant (must) done~~

2. ~~fp8\_gemm (must) done~~

3. If world size > 1 还会涉及到 ~~dist.all\_reduce(tensor, op=ReduceOp.SUM)~~

**others**

1. **aten::zero\_ (mid)**

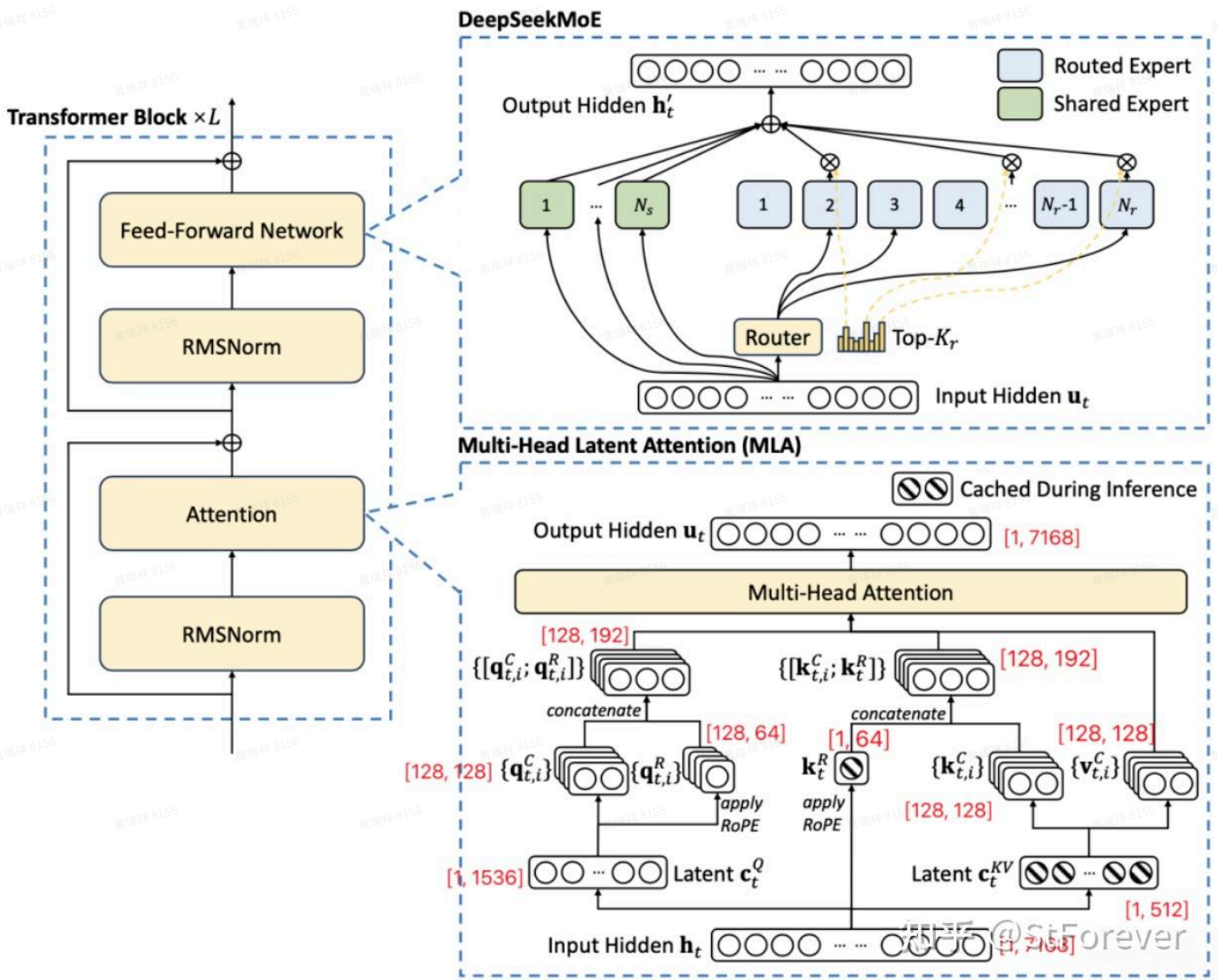
用来初始化kv cache, 可以使用aten::empty(sipu已经实现)代替

2. **aten::bmm.out (must)done**

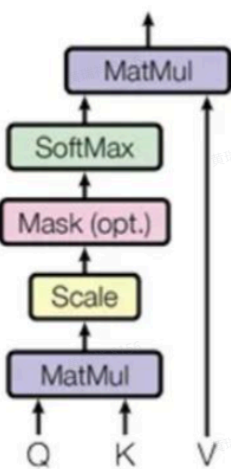
torch.einsum()会调用

3. slice操作

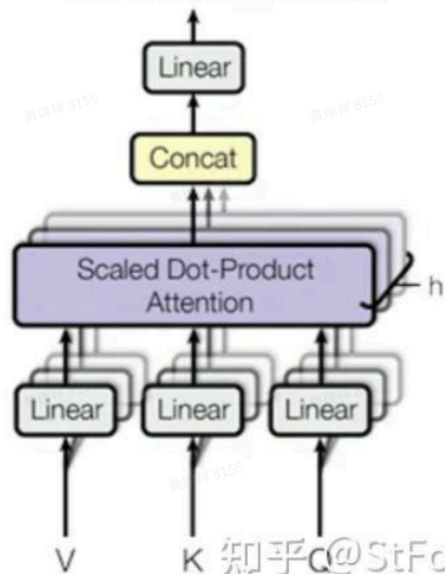
bug, 本身已经实现, 但是会有core dump



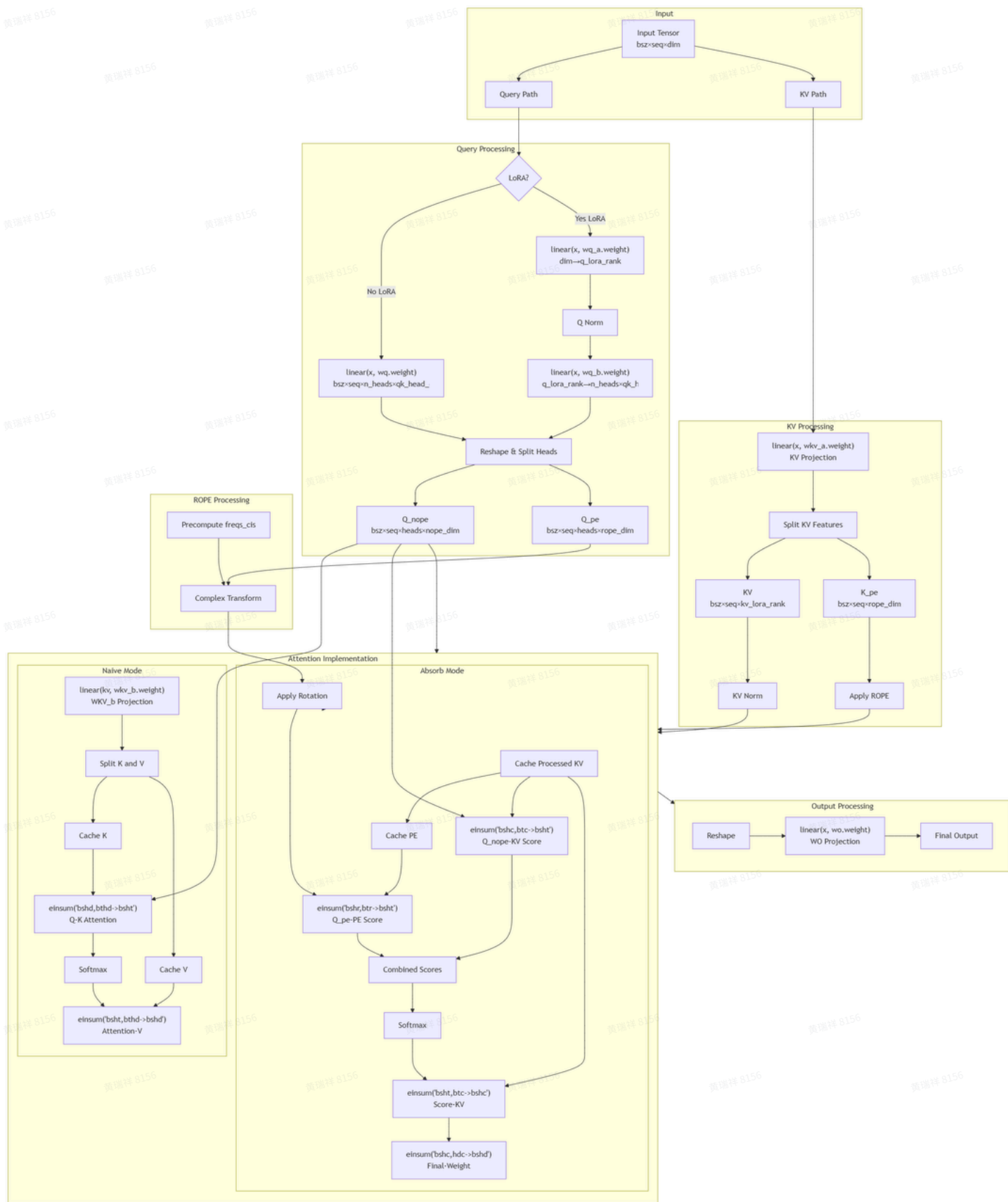
### Scaled Dot-Product Attention



### Multi-Head Attention



实现流程:



放在SIPU的模块为 attention implementation (包含bmm, softmax), linear(和RMS\_norm和ROPE穿插, 调小了模型, 目前可以在SIPU侧运行因为copy size restrict的问题, 目前无法放在SIPU侧, 单测bf16是ok的), RMS\_norm(SIPU未实现), ROPE(SIPU未实现)目前只能先放在CPU侧处理

## Embedding 推理

## F.embedding() trace

```
1  [call] op=[aten::embedding], key=[PythonTLSSnapshot]
2  [redispatchBoxed] op=[aten::embedding], key=[AutogradPrivateUse1]
3  [redispatch] op=[aten::embedding], key=[Python]
4  [callBoxed] op=[aten::embedding], key=[Python]
5  [callBoxed] op=[aten::embedding], key=[PrivateUse1]
6  [call] op=[aten::reshape], key=[PrivateUse1]
7  [call] op=[aten::view], key=[PrivateUse1]
8  [call] op=[aten::index_select], key=[PrivateUse1]
9  [call] op=[aten::_to_cpu], key=[PrivateUse1]
10 [call] op=[aten::to.dtype_layout], key=[BackendSelect]
11 [redispatch] op=[aten::to.dtype_layout], key=[PrivateUse1]
12 [call] op=[aten::_to_copy], key=[BackendSelect]
13 [redispatch] op=[aten::_to_copy], key=[PrivateUse1]
14 [call] op=[aten::empty_strided], key=[BackendSelect]
15 [redispatch] op=[aten::empty_strided], key=[CPU] (sipu_torch已接入)
16 [call] op=[aten::copy_], key=[PrivateUse1]
17 [call] op=[aten::_copy_from], key=[PrivateUse1]
18 [call] op=[aten::to.dtype_layout], key=[BackendSelect]
19 [redispatch] op=[aten::to.dtype_layout], key=[PrivateUse1]
20 [call] op=[aten::_to_copy], key=[BackendSelect]
21 [redispatch] op=[aten::_to_copy], key=[PrivateUse1]
22 [call] op=[aten::empty_strided], key=[BackendSelect]
23 [redispatch] op=[aten::empty_strided], key=[CPU] (sipu_torch已接入)
24 [call] op=[aten::copy_], key=[PrivateUse1]
25 [call] op=[aten::_copy_from], key=[PrivateUse1]
26 [redispatchBoxed] op=[aten::index_select], key=[CPU] (未实现)
27 [call] op=[aten::empty.memory_format], key=[BackendSelect]
28 [redispatch] op=[aten::empty.memory_format], key=[CPU] (sipu_torch已接
入)
29 [call] op=[aten::select.int], key=[CPU] (CompositeExplicitAutograd, 依赖
为as_strided, sipu_torch已接入)
30 [call] op=[aten::as_strided], key=[CPU] (sipu_torch已接入)
31 [call] op=[aten::select.int], key=[CPU] (CompositeExplicitAutograd, 依赖
为as_strided, sipu_torch已接入)
32 [call] op=[aten::as_strided], key=[CPU] (sipu_torch已接入)
33 [call] op=[aten::to.dtype_layout], key=[BackendSelect]
34 [redispatch] op=[aten::to.dtype_layout], key=[PrivateUse1]
35 [call] op=[aten::_to_copy], key=[BackendSelect]
36 [redispatch] op=[aten::_to_copy], key=[PrivateUse1]
37 [call] op=[aten::empty_strided], key=[BackendSelect]
38 [redispatch] op=[aten::empty_strided], key=[PrivateUse1]
39 [14:01:10.016][SIRT][W][2260291] [device.cpp:120] Legacy compatibility mode:
Converting SIPU to 1U16C
40 [call] op=[aten::copy_], key=[PrivateUse1]
41 [call] op=[aten::_copy_from], key=[PrivateUse1]
```

## 1. aten::select.int

- a. CompositeExplicitAutograd, 依赖算子为as\_strided, 已经实现, 经过测试可以完全在SIPU上运行

select trace

```

1  start select
2  [call] op=[aten::select.int], key=[PythonTLSSnapshot]
3  [redispatchBoxed] op=[aten::select.int], key=[AutogradPrivateUse1]
4  [redispatch] op=[aten::select.int], key=[ADInplaceOrView]
5  [redispatch] op=[aten::select.int], key=[Python]
6  [callBoxed] op=[aten::select.int], key=[Python]
7  [callBoxed] op=[aten::select.int], key=[PrivateUse1]
8  [call] op=[aten::as_strided], key=[PrivateUse1]

```

## 2. aten::index\_select 需要实现

## MLP/MoE(ffn)推理

MLP

```

1  class MLP(nn.Module):
2      """
3      Multi-Layer Perceptron (MLP) used as a feed-forward layer.
4
5      Attributes:
6          w1 (nn.Module): Linear layer for input-to-hidden transformation.
7          w2 (nn.Module): Linear layer for hidden-to-output transformation.
8          w3 (nn.Module): Additional linear layer for feature transformation.
9      """
10     def __init__(self, dim: int, inter_dim: int):
11         """
12         Initializes the MLP layer.
13
14         Args:
15             dim (int): Input and output dimensionality.
16             inter_dim (int): Hidden layer dimensionality.
17         """
18         super().__init__()
19         self.w1 = ColumnParallelLinear(dim, inter_dim)
20         self.w2 = RowParallelLinear(inter_dim, dim)
21         self.w3 = ColumnParallelLinear(dim, inter_dim)
22

```

```

23     def forward(self, x: torch.Tensor) -> torch.Tensor:
24         """
25         Forward pass for the MLP layer.
26
27         Args:
28             x (torch.Tensor): Input tensor.
29
30         Returns:
31             torch.Tensor: Output tensor after MLP computation.
32         """
33         return self.w2(F.silu(self.w1(x)) * self.w3(x))
34
35

```

涉及到的需要实现的算子

37 **aten::silu**

MoE

```

1     class Expert(nn.Module):
2         """
3         Expert layer for Mixture-of-Experts (MoE) models.
4         """
5         def __init__(self, dim: int, inter_dim: int):
6             super().__init__()
7             self.w1 = Linear(dim, inter_dim)
8             self.w2 = Linear(inter_dim, dim)
9             self.w3 = Linear(dim, inter_dim)
10
11        def forward(self, x: torch.Tensor) -> torch.Tensor:
12            """
13            Forward pass for the Expert layer.
14            """
15            return self.w2(F.silu(self.w1(x)) * self.w3(x))

```

涉及到的需要实现的算子与MLP一致

17 **aten::silu**

```

20    class Gate(nn.Module):
21        """
22        Gating mechanism for routing inputs in a mixture-of-experts (MoE) model.
23
24        Attributes:
25            dim (int): Dimensionality of input features.
26            topk (int): Number of top experts activated for each input.
27            n_groups (int): Number of groups for routing.
28            topk_groups (int): Number of groups to route inputs to.

```

```

29     score_func (str): Scoring function ('softmax' or 'sigmoid').
30     route_scale (float): Scaling factor for routing weights.
31     weight (torch.nn.Parameter): Learnable weights for the gate.
32     bias (Optional[torch.nn.Parameter]): Optional bias term for the gate.
33     """
34     def __init__(self, args: ModelArgs):
35         """
36         Initializes the Gate module.
37
38         Args:
39             args (ModelArgs): Model arguments containing gating parameters.
40         """
41         super().__init__()
42         self.dim = args.dim
43         self.topk = args.n_activated_experts
44         self.n_groups = args.n_expert_groups
45         self.topk_groups = args.n_limited_groups
46         self.score_func = args.score_func
47         self.route_scale = args.route_scale
48         self.weight = nn.Parameter(torch.empty(args.n_routed_experts,
args.dim))
49         self.bias = nn.Parameter(torch.empty(args.n_routed_experts)) if
self.dim == 7168 else None
50
51     def forward(self, x: torch.Tensor) -> Tuple[torch.Tensor, torch.Tensor]:
52         """
53         Forward pass for the gating mechanism.
54
55         Args:
56             x (torch.Tensor): Input tensor.
57
58         Returns:
59             Tuple[torch.Tensor, torch.Tensor]: Routing weights and selected
expert indices.
60         """
61         scores = linear(x, self.weight)
62         if self.score_func == "softmax":
63             scores = scores.softmax(dim=-1, dtype=torch.float32)
64         else:
65             scores = scores.sigmoid()
66         original_scores = scores
67         if self.bias is not None:
68             scores = scores + self.bias
69         if self.n_groups > 1:
70             scores = scores.view(x.size(0), self.n_groups, -1)
71             if self.bias is None:
72                 group_scores = scores.amax(dim=-1)

```

```

73         else:
74             group_scores = scores.topk(2, dim=-1)[0].sum(dim=-1)
75             indices = group_scores.topk(self.topk_groups, dim=-1)[1]
76             mask = scores.new_ones(x.size(0), self.n_groups,
dtype=bool).scatter_(1, indices, False)
77             scores = scores.masked_fill_(mask.unsqueeze(-1), float("-
inf")).flatten(1)
78             indices = torch.topk(scores, self.topk, dim=-1)[1]
79             weights = original_scores.gather(1, indices)
80             if self.score_func == "sigmoid":
81                 weights /= weights.sum(dim=-1, keepdim=True)
82             weights *= self.route_scale
83             return weights.type_as(x), indices

```

涉及到的需要实现的算子

aten::fill\_.Scalar

aten::sigmoid

aten::amax.out

aten::topk.values

aten::gather.out

aten::sum.IntList\_out

aten::div.out

to.dtype

Data type cast for direct\_copy\_kernel\_sipu is not supported now

```

97 class MoE(nn.Module):

```

```

98     """

```

```

99     Mixture-of-Experts (MoE) module.

```

```

100     Attributes:

```

```

101         dim (int): Dimensionality of input features.

```

```

102         n_routed_experts (int): Total number of experts in the model.

```

```

103         n_local_experts (int): Number of experts handled locally in
distributed systems.

```

```

104         n_activated_experts (int): Number of experts activated for each input.

```

```

105         gate (nn.Module): Gating mechanism to route inputs to experts.

```

```

106         experts (nn.ModuleList): List of expert modules.

```

```

107         shared_experts (nn.Module): Shared experts applied to all inputs.

```

```

108     """

```

```

109     def __init__(self, args: ModelArgs):

```

```

110         """

```

```

111         Initializes the MoE module.

```

```

112     Args:

```

```

113         args (ModelArgs): Model arguments containing MoE parameters.

```

```

114     """

```

```

117     super().__init__()
118     self.dim = args.dim
119     assert args.n_routed_experts % world_size == 0, f"Number of experts
must be divisible by world size (world_size={world_size})"
120     self.n_routed_experts = args.n_routed_experts
121     self.n_local_experts = args.n_routed_experts // world_size
122     self.n_activated_experts = args.n_activated_experts
123     self.experts_start_idx = rank * self.n_local_experts
124     self.experts_end_idx = self.experts_start_idx + self.n_local_experts
125     self.gate = Gate(args)
126     self.experts = nn.ModuleList([Expert(args.dim, args.moe_inter_dim) if
self.experts_start_idx <= i < self.experts_end_idx else None
127                                     for i in range(self.n_routed_experts)])
128     self.shared_experts = MLP(args.dim, args.n_shared_experts *
args.moe_inter_dim)
129
130     def forward(self, x: torch.Tensor) -> torch.Tensor:
131         """
132         Forward pass for the MoE module.
133
134         Args:
135             x (torch.Tensor): Input tensor.
136
137         Returns:
138             torch.Tensor: Output tensor after expert routing and computation.
139         """
140         shape = x.size()
141         x = x.view(-1, self.dim)
142         weights, indices = self.gate(x)
143         y = torch.zeros_like(x)
144         counts = torch.bincount(indices.flatten(),
minlength=self.n_routed_experts).tolist()
145         for i in range(self.experts_start_idx, self.experts_end_idx):
146             if counts[i] == 0:
147                 continue
148             expert = self.experts[i]
149             idx, top = torch.where(indices == i)
150             y[idx] += expert(x[idx]) * weights[idx, top, None]
151         z = self.shared_experts(x)
152         if world_size > 1:
153             dist.all_reduce(y)
154         return (y + z).view(shape)

```

156 涉及到的需要实现的算子

157 `aten::bincount`

158 `aten::eq.Scalar_out`

涉及到的需要sipu实现的算子

## MLP module

### 1. `aten::silu`

## Gate module

### 1. `aten::fill_.Scalar`

### 2. `aten::sigmoid`

### 3. `aten::amax.out` (new)

### 4. `aten::topk.values`

### 5. `aten::gather.out`

### 6. `aten::sum.IntList_out`

### 7. `aten::div.out`

### 8. `to.dtype:` ~~Data type cast for direct\_copy\_kernel\_sipu is not supported now~~

## MoE module

### 1. `aten::bincount` ~~(new)~~

### 2. `aten::eq.Scalar_out`

### 3. `aten::nonzero` (new)

## 630 DS 目前存在问题 (0623更新)

## Nonzero

### 1. `bool`类型的支持算子

代码块

```
1 RuntimeError: nonzero only supports int32 or int64 tensors on sipu, but get:
   bool
```

### 2. Running time问题 框架



```
34 File "/data/users/liuhao2/pytorch/torch/nn/modules/module.py", line 1736, in
   _wrapped_call_impl
35     return self._call_impl(*args, **kwargs)
36         ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
37 File "/data/users/liuhao2/pytorch/torch/nn/modules/module.py", line 1747, in
   _call_impl
38     return forward_call(*args, **kwargs)
39         ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
40 File
   "/data/users/liuhao2/torch_sipu/examples/run_deepseekv3_Transformer.py", line
   1238, in forward
41     h = layer(h, start_pos, freqs_cis, mask, device=device)
42         ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
43 File "/data/users/liuhao2/pytorch/torch/nn/modules/module.py", line 1736, in
   _wrapped_call_impl
44     return self._call_impl(*args, **kwargs)
45         ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
46 File "/data/users/liuhao2/pytorch/torch/nn/modules/module.py", line 1747, in
   _call_impl
47     return forward_call(*args, **kwargs)
48         ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
49 File
   "/data/users/liuhao2/torch_sipu/examples/run_deepseekv3_Transformer.py", line
   1174, in forward
50     x = x.cpu().bfloat16().to(device) + self.ffn(
51         ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
52 File "/data/users/liuhao2/pytorch/torch/nn/modules/module.py", line 1736, in
   _wrapped_call_impl
53     return self._call_impl(*args, **kwargs)
54         ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
55 File "/data/users/liuhao2/pytorch/torch/nn/modules/module.py", line 1747, in
   _call_impl
56     return forward_call(*args, **kwargs)
57         ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
58 File
   "/data/users/liuhao2/torch_sipu/examples/run_deepseekv3_Transformer.py", line
   628, in forward
59     print("idx", idx)
60 File "/data/users/liuhao2/pytorch/torch/_tensor.py", line 523, in __repr__
61     return torch._tensor_str._str(self, tensor_contents=tensor_contents)
62         ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
63 File "/data/users/liuhao2/pytorch/torch/_tensor_str.py", line 708, in _str
64     return _str_intern(self, tensor_contents=tensor_contents)
65         ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
66 File "/data/users/liuhao2/pytorch/torch/_tensor_str.py", line 625, in
   _str_intern
67     tensor_str = _tensor_str(self, indent)
```



```
代码块 group_scores = scores.amax(dim=-1, keepdim=True)
```

当expert groups>1的情况下，scores为3维，目前框架这边只能接受2维的输入，目前只对last dim求amax，应该可以在框架侧处理合并外层维度

## fill

进展：已解决

算子 RuntimeError: "fill\_func" not implemented for 'Bool' block by compiler

代码块

```
1 mask = scores.new_ones(x.size(0), self.n_groups, dtype=bool).scatter_(
2     1, indices, False
3     )
```

## tile\_format优化分析

### MLA

linear之间有rms norm和cat，flatten等操作，难以优化

### MLP

代码块

```
1 class MLP(nn.Module):
2     """
3     Multi-Layer Perceptron (MLP) used as a feed-forward layer.
4
5     Attributes:
6         w1 (nn.Module): Linear layer for input-to-hidden transformation.
7         w2 (nn.Module): Linear layer for hidden-to-output transformation.
8         w3 (nn.Module): Additional linear layer for feature transformation.
9     """
10
11     def __init__(self, dim: int, inter_dim: int, seed: Optional[int] = None):
12         """
```

```

13     Initializes the MLP layer.
14
15     Args:
16         dim (int): Input and output dimensionality.
17         inter_dim (int): Hidden layer dimensionality.
18         seed (Optional[int]): Seed for weight initialization to ensure
identical
19                                     weights across different devices. Default is
None.
20     """
21     super().__init__()
22     # Use predictable seeds for each layer to ensure identical
initialization
23     self.w1 = ColumnParallelLinear(
24         dim, inter_dim, seed=seed if seed is None else seed + 1
25     )
26     self.w2 = RowParallelLinear(
27         inter_dim, dim, seed=seed if seed is None else seed + 2
28     )
29     self.w3 = ColumnParallelLinear(
30         dim, inter_dim, seed=seed if seed is None else seed + 3
31     )
32
33     def forward(self, x: torch.Tensor, device: torch.device = "cpu") ->
torch.Tensor:
34     """
35     Forward pass for the MLP layer.
36
37     Args:
38         x (torch.Tensor): Input tensor.
39
40     Returns:
41         torch.Tensor: Output tensor after MLP computation.
42     """
43
44     res = self.w2(
45         F.silu(self.w1(x, device=device).to(device))
46         * self.w3(x, device=device).to(device),
47         device=device,
48     )
49     return res
50

```

linear1 (tile) -> silu(linear) -> linear2(tile)-> mul(linear)

^  
|  
linear3(tile)

tile -> linear 3 times

linear -> tile 1 times

如果silu 和 mul支持tileformat, MLP内部可以一直保持tileformat计算, 不会再涉及转换问题

## MoE & MoE中的Expert

代码块

```
1 class MoE(nn.Module):
2     """
3     Mixture-of-Experts (MoE) module.
4
5     Attributes:
6         dim (int): Dimensionality of input features.
7         n_routed_experts (int): Total number of experts in the model.
8         n_local_experts (int): Number of experts handled locally in
9         distributed systems.
10        n_activated_experts (int): Number of experts activated for each input.
11        gate (nn.Module): Gating mechanism to route inputs to experts.
12        experts (nn.ModuleList): List of expert modules.
13        shared_experts (nn.Module): Shared experts applied to all inputs.
14
15    """
16
17    def __init__(self, args: ModelArgs):
18        """
19        Initializes the MoE module.
20
21        Args:
22            args (ModelArgs): Model arguments containing MoE parameters.
23        """
24        super().__init__()
25        self.dim = args.dim
26        assert (
27            args.n_routed_experts % world_size == 0
28        ), f"Number of experts must be divisible by world size (world_size=
29        {world_size})"
30        self.n_routed_experts = args.n_routed_experts
31        self.n_local_experts = args.n_routed_experts // world_size
32        self.n_activated_experts = args.n_activated_experts
33        self.experts_start_idx = rank * self.n_local_experts
```

```

31     self.experts_end_idx = self.experts_start_idx + self.n_local_experts
32     self.gate = Gate(args)
33     self.experts = nn.ModuleList(
34         [
35             Expert(args.dim, args.moe_inter_dim)
36             if self.experts_start_idx <= i < self.experts_end_idx
37             else None
38             for i in range(self.n_routed_experts)
39         ]
40     )
41     self.shared_experts = MLP(args.dim, args.n_shared_experts *
args.moe_inter_dim)
42
43     def forward(self, x: torch.Tensor, device: torch.device = "cpu") ->
torch.Tensor:
44         """
45         Forward pass for the MoE module.
46
47         Args:
48             x (torch.Tensor): Input tensor.
49
50         Returns:
51             torch.Tensor: Output tensor after expert routing and computation.
52         """
53         shape = x.size()
54         x = x.view(-1, self.dim)
55         weights, indices = self.gate(x, device=device)
56         y = torch.zeros_like(x).to(device)
57         #breakpoint()
58         counts = torch.bincount(
59             indices.flatten(), minlength=self.n_routed_experts
60         ).tolist()
61         for i in range(self.experts_start_idx, self.experts_end_idx):
62             if counts[i] == 0:
63                 continue
64             expert = self.experts[i].to(device)
65             idx, top = torch.where(indices.cpu() == i)
66             print("indices.cpu().shape", indices.cpu().shape)
67
68             weights = weights.cpu()
69
70             y[idx] += expert(x[idx].to(device), device).to(device) * weights[
71                 idx, top, None
72             ].to(device)
73
74         z = self.shared_experts(x, device=device)
75         if world_size > 1:

```

```

76         dist.all_reduce(y)
77         return (y + z).view(shape)
78
79
80 class Expert(nn.Module):
81     """
82     Expert layer for Mixture-of-Experts (MoE) models.
83     """
84
85     def __init__(self, dim: int, inter_dim: int):
86         super().__init__()
87         self.w1 = Linear(dim, inter_dim)
88         self.w2 = Linear(inter_dim, dim)
89         self.w3 = Linear(dim, inter_dim)
90
91     def forward(self, x: torch.Tensor, device: torch.device = "cpu") ->
torch.Tensor:
92         """
93         Forward pass for the Expert layer.
94         """
95         return self.w2(
96             F.silu(self.w1(x, device=device)).to(device) * self.w3(x,
device=device).to(device),
97             device=device,
98         )
99

```

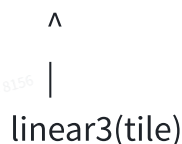
Expert和MLP本质上是一样的，只是feature dim不同，这里的情况和上面MLP中分析的一致  
 如果可以按照MLP中所说的那样进行优化，MoE的收益会更加明显

n\_shared\_experts: int = 2 #原始配置

n\_activated\_experts: int = 6 #原始配置

MoE中的每一个layer激活的路由专家（Expert模块）有6个，共享专家（MLP模块）2个

linear1 (tile) -> silu(linear) -> linear2(tile)-> mul(linear)



tile -> linear 3 times

linear -> tile 1 times

如果silu 和 mul支持tileformat，MLP内部可以一直保持tileformat计算，不会再涉及转换问题，

每一个layer的MoE模块可以减少  $(1+3) * 8 = 24$ 次format的转换

## Infer fallback 分析

除了init weight 和这个预计算ROPE的precompute\_freqs\_cis fallback到cpu来做，其余还fallback的还有

1. [redispatchBoxed] op=[aten::triu.out], key=[CPU]

代码块

```
1 mask = torch.full((seqlen, seqlen), float("-inf"),  
device=tokens.device).triu_(1)
```

2. [redispatchBoxed] op=[aten::zero\_], key=[CPU]

代码块

```
1 y = torch.zeros_like(x)
```

3. [redispatchBoxed] op=[aten::index.Tensor\_out], key=[CPU]

代码块

```
1 h = self.norm(h)[ :, -1]
```

4. [redispatchBoxed] op=[aten::\_index\_put\_impl\_], key=[CPU] 算子 宋涵 正在开发

代码块

```
1 mask = (x < self.vocab_start_idx) | (x >= self.vocab_end_idx)  
2 x = x - self.vocab_start_idx  
3 x[mask] = 0
```

5. [redispatchBoxed] op=[aten::reciprocal.out], key=[CPU] 应该是一些除法自动触发了优化，按照倒数的方式乘法计算 接入int

6. Add type

代码块

```

1 class ParallelEmbedding(nn.Module):输出结果为F32
2 在class Block(nn.Module):当中
3 x = x + self.attn(self.attn_norm(x), start_pos, freqs_cis, mask, device=device)
4 x = x + self.ffn(self.ffn_norm(x), device=device)
5 x 为 f32, self.attn, self.ffn为BF16, 直接相加 nan

```

## 7. 开源 infer code 修改 group\_scores = scores.amax(dim=-1, keepdim=True) fixed amax支持 keepdim=false

代码块

```

1 if self.n_groups > 1:
2     scores = scores.view(x.size(0), self.n_groups, -1)
3     if self.bias is None:
4         group_scores = scores.amax(dim=-1)
5     else:
6         group_scores = scores.topk(2, dim=-1)[0].sum(dim=-1)
7         indices = group_scores.topk(self.topk_groups, dim=-1)[1]
8         mask = scores.new_ones(x.size(0), self.n_groups, dtype=bool).scatter_(1,
indices, False)
9         scores = scores.masked_fill_(mask.unsqueeze(-1), float("-inf")).flatten(1)

```

改为

代码块

```

1 if self.n_groups > 1:
2     scores = scores.view(x.size(0), self.n_groups, -1)
3     if self.bias is None:
4         group_scores = scores.amax(dim=-1, keepdim=True)
5     else:
6         group_scores = scores.topk(2, dim=-1)[0].sum(dim=-1)
7         indices = group_scores.topk(self.topk_groups, dim=-1)[1]
8         mask = scores.new_ones(x.size(0), self.n_groups,
dtype=bool).to(device)
9         indices_for_scatter = indices.squeeze(-1) if indices.dim() > 2
else indices
10        mask = mask.scatter_(1, indices_for_scatter.to(device), False)
11        scores = scores.masked_fill_(mask.unsqueeze(-1), float("-
inf")).flatten(1)

```