

OP 的开发文档

Op接入资料&规范: [📖 开发备注](#)

_copy_from

Function 接口

子接口 1fused_add_rms_norm

contiguous, 也就是将 src tensor 中的内容重新排布后拷贝到 dst, 保证 dst 张量的元素在内存中的布局是线性、顺序的, 并且每个元素在内存中的地址都是连续的, 也就是 *contiguous* 的

一个 PyTorch 张量 T 是 *contiguous* 的, 当且仅当其步幅 (stride) 满足以下条件:

对于每个维度 i 来说, 步幅 $\text{stride}(i)$ 满足:

$$\text{stride}(i) = \prod_{j=i+1}^{n-1} \text{size}(j)$$

其中:

- n 是张量的维度数。
- $\text{size}(j)$ 是张量在第 j 个维度的大小 (shape)。
- $\text{stride}(i)$ 表示张量在维度 i 上相邻元素之间的内存步长。

这个表达式意味着元素在内存中按行优先顺序排列。如果不满足此条件, 张量就是非连续的。

接口调用方需要保证 src 与 dst tensor 的 shape 一致

代码块

```
1  /**
2   dims: dst tensor dims
3   sizes: output tensor sizes
4   strides: strides of nargs tensors, 0 means it's a broadcasted dim
5   element_size: element_size of tensor
6   **/
7  void fn(dev, const void* dst, const void* self, int dims, const int64_t*
8         sizes, const int64_t* const* strides, const int64_t element_size)
9  // 可选的形式(updated)
```

```
10 void fnv2(dev, const void* dst, const void* self, const int64_t self_ndims,
const int64_t* self_sizes, const int64_t* dst_strides, const int64_t*
self_strides, const int64_t element_size)
```

一些例子

代码块

```
1 // case 1
2 // python
3 b = torch.randn(4, 3).long().cuda()
4 b = b.t()
5 c = b.contiguous()
6 // kernel
7 copy_kernel_cuda: dst_device: cuda:0, src_device: cuda:0
8 copy_kernel_cuda: dst_sizes: [3, 4], src_sizes: [3, 4]
9 copy_kernel_cuda: dst_strides: [4, 1], src_strides: [1, 3]
10 copy_kernel_cuda: iter_shape: [4, 3]
11 copy_kernel_cuda: iter_dst_strides: [8, 32], iter_src_strides: [24, 8]
12
13 // case 2
14 // python
15 a = torch.randn(3).long().cuda()
16 b = torch.randn(4, 3).long().cuda()
17 a = a.expand_as(b).contiguous()
18 // kernel
19 copy_kernel_cuda: dst_device: cuda:0, src_device: cuda:0
20 copy_kernel_cuda: dst_sizes: [4, 3], src_sizes: [4, 3]
21 copy_kernel_cuda: dst_strides: [3, 1], src_strides: [0, 1]
22 copy_kernel_cuda: iter_shape: [3, 4]
23 copy_kernel_cuda: iter_dst_strides: [8, 24], iter_src_strides: [8, 0]
24
25 // case3
26 // python
27 a = torch.randn(4, 3, 2).long().cuda()
28 b = torch.randn(1, 3, 1).long().cuda()
29 a.copy_(b)
30 // kernel
31 copy_kernel_cuda: dst_device: cuda:0, src_device: cuda:0
32 copy_kernel_cuda: dst_sizes: [4, 3, 2], src_sizes: [1, 3, 1]
33 copy_kernel_cuda: dst_strides: [6, 2, 1], src_strides: [3, 1, 1]
34 copy_kernel_cuda: iter_shape: [2, 3, 4]
35 copy_kernel_cuda: iter_dst_strides: [8, 16, 48], iter_src_strides: [0, 8, 0]
36
37 // case4
38 // python
```

```
39     a = torch.randn(4, 3, 2).long().cuda()
40     b = torch.randn(3, 1).long().cuda()
41     a.copy_(b)
42     // kernel
43 copy_kernel_cuda: dst_device: cuda:0, src_device: cuda:0
44 copy_kernel_cuda: dst_sizes: [4, 3, 2], src_sizes: [3, 1]
45 copy_kernel_cuda: dst_strides: [6, 2, 1], src_strides: [1, 1]
46 copy_kernel_cuda: iter_shape: [2, 3, 4]
47 copy_kernel_cuda: iter_dst_strides: [8, 16, 48], iter_src_strides: [0, 8, 0]
```

torch 中的调用路径:

```
aten/src/ATen/native/Copy.cpp copy_impl -> copy_stub -> aten/src/ATen/native/cuda/Copy.cu
copy_kernel_cuda -> aten/src/ATen/native/cuda/Copy.cu copy_device_to_device ->
aten/src/ATen/native/cuda/Copy.cu direct_copy_kernel_cuda ->
torch/include/ATen/native/cuda/Loops.cuh gpu_kernel ->
torch/include/ATen/native/cuda/CUDALoops.cuh gpu_kernel_impl ->
torch/include/ATen/native/cuda/CUDALoops.cuh gpu_kernel_impl_nocast
```

ISSUE

dtype 可能不一致, 需要 kernel 做转换, 转换的参数如何定义

`_copy_from_and_resize`, `_copy_from_and_resize.out`

参考信息

采用CPUFallback机制实现, 通过 `copy_impl` 中的逻辑——先保证尺寸匹配、连续性, 再调用通用的 `copy_stub` 内核, 可以参考 `copy_from` 的实现 [OP 的开发文档](#)

ISSUE

`_softmax`

`_softmax.out`

参考信息

OP List 链接: <https://siorigin.feishu.cn/record/UC8trohn4eZvxBc4JAJcGsGAn6b>

Pytorch 官方文档参考: <https://pytorch.org/docs/stable/generated/torch.nn.Softmax.html>

Function 接口

代码块

1

```

2 void fn(dev, out*, input*, bool half_to_float, int64_t outer_size, dim_size,
  inner_size)
3 // (updated)
4 void fn(dev, const void* out, const void* self, bool half_to_float, int64_t
  outer_size, dim_size, inner_size)
5 // 其中dim_size是dim参数所在维度的size, outer_size是所在维度小于dim参数的所有维度大
  小:
6 // for (int64_t i = 0; i < dim; ++i)
7 //   outer_size *= input.size(i);
8 // 例如对于shape为[A,B,C,D,E]的Tensor在C维度做Softmax, 会view成[A]

```

- 其中 dim_size 是 dim 参数所在维度的 size, outer size 是所在维度小于 dim 参数的所有维度大小:
例如对于 shape 为[A, B, C, D, E]的 Tensor 在 C 维度做 Softmax, 会 view 成[A*B, C, D*E], 对应的 stride 为[C*D*E, D*E, 1], outer_size=A*B, inner_size=D*E, dim_size=C
- half_to_float 代表是否要将 HALF 的 input 转换成 FLOAT (可以不支持 BFLOAT16)
- Accumulate 的 Dtype 参考: [OP Summary](#)

ISSUE

Dtype 支持

float32

float16

bfloat16 目前缺少对应的指令

mxformat

非 contiguous 的是否能支持 (如果支持需要修改 Function 接口)

Memory 是否需要 Align, Align 到多少

目前需要软件将 dim_size padding 到 256 的倍数

需要优化

后期需要支持任意维度上做 softmax

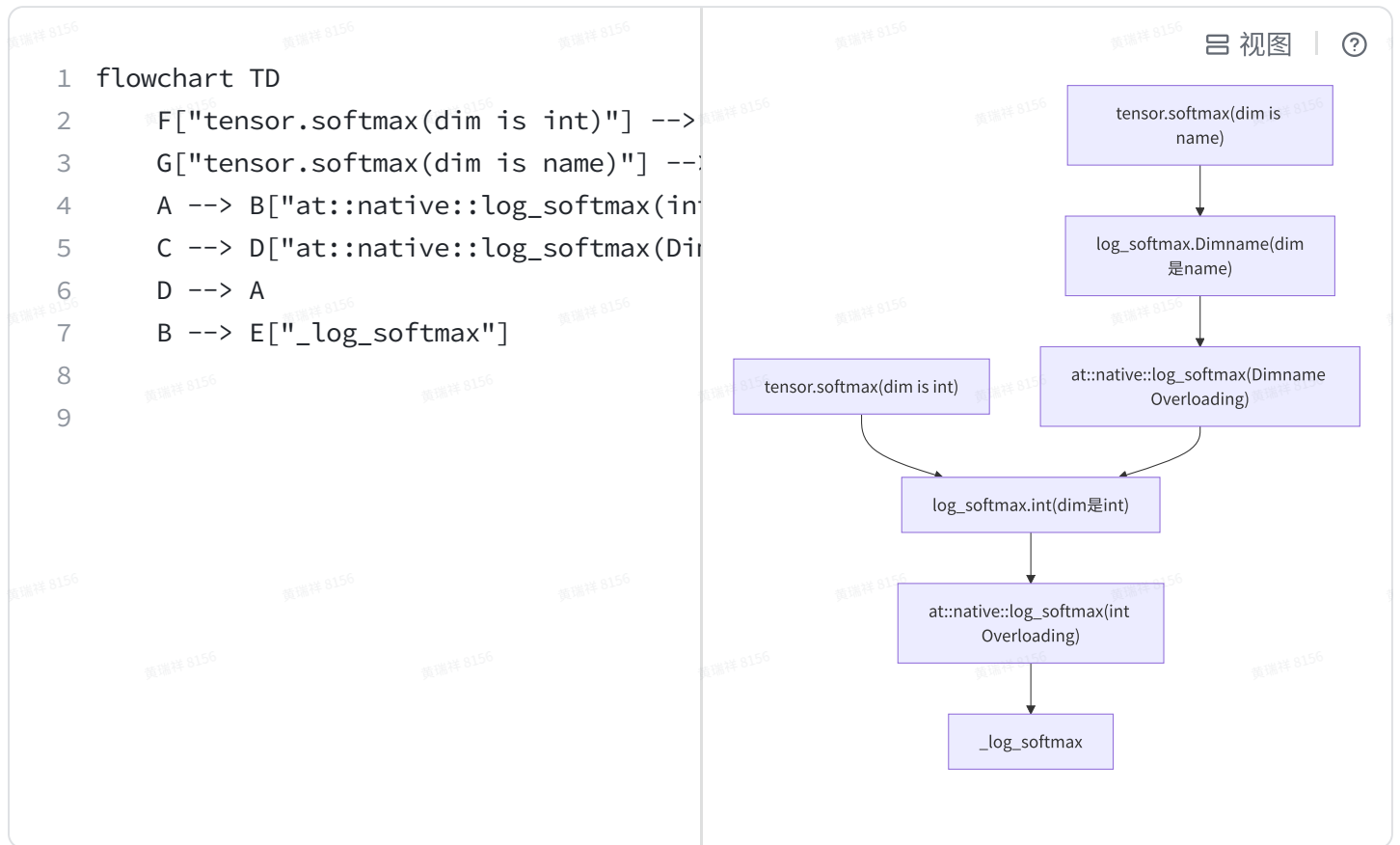
- a. 提供给算法的原型还是按照下面表格封装, 但是第一步先只支持在最 inner 维度做 softmax, 原因参考下面的分析。
- b. 第二步视情况考虑实现任意维度 softmax, 当前有个想法是在 RVV 上利用 index load/store 来做 gather, 将参与计算的元素都收集起来连续存放, 不用 tile core 的原因是 tile 的 index load/store 粒度太大, 最小也是 32B, 而这个场景的需求是按照元素来 index load/store. [softmax](#)

log_softmax.int, log_softmax.Dimname, _log_softmax

OP

Pytorch 官方文档参考:

<https://pytorch.org/docs/stable/generated/torch.nn.LogSoftmax.html#torch.nn.LogSoftmax>



_log_softmax

mathematically equivalent to $\log(\text{softmax}(x))$

参考信息

OP List 链接: <https://siorigin.feishu.cn/record/GXmArsJOSeXZepcWqO3clgQXnah>

Pytorch 官方文档参考:

<https://pytorch.org/docs/stable/generated/torch.nn.LogSoftmax.html#torch.nn.LogSoftmax>

Function 接口

```
代码块
1 void fn(dev, out*, input*, bool half_to_float, int64_t outer_size, dim_size,
  inner_size)
2 //(updated)
3 void fn(dev, const void* out, const void* self, bool half_to_float, int64_t
  outer_size, dim_size, inner_size)
4 // 其中dim_size是dim参数所在维度的size, outer size是所在维度小于dim参数的所有维度大
  小:
```

```
5 // for (int64_t i = 0; i < dim; ++i)
6 //     outer_size *= input.size(i);
7 // 例如对于shape为[A,B,C,D,E]的Tensor在C维度做Softmax, 会view成[A]
```

ISSUE

Dtype 支持

float32

float16

bfloat16 目前缺少对应的指令

mxformat

非 contiguous 的是否能支持 (如果支持需要修改 Function 接口)

Memory 是否需要 Align, Align 到多少

目前需要软件将 dim_size padding 到 256 的倍数

需要优化

后期需要支持任意维度上做 softmax

- 提供给算法的原型还是按照下面表格封装, 但是第一步先只支持在最 inner 维度做 softmax, 原因参考下面的分析。
- 第二步视情况考虑实现任意维度 softmax, 当前有个想法是在 RVV 上利用 index load/store 来做 gather, 将参与计算的元素都收集起来连续存放, 不用 tile core 的原因是 tile 的 index load/store 粒度太大, 最小也是 32B, 而这个场景的需求是按照元素来 index load/store。 [softmax](#)

_local_scalar_dense

在将 scalar 转换为 cpu scalar 时会调用, 应该直接使用 runtime API 就可以实现, 参考 CUDA 实现:

_local_scalar_dense_cuda

代码块

```
1 // schema: func:
  _local_scalar_dense(Tensor self) ->
  Scalar
```

```
In [12]: torch.randn(1).cuda().item()
[call] op=[aten::randn], key=[BackendSelect]
[redispatch] op=[aten::randn], key=[CPU]
[call] op=[aten::empty_memory_format], key=[BackendSelect]
[redispatch] op=[aten::empty_memory_format], key=[CPU]
[call] op=[aten::normal_], key=[CPU]
[call] op=[aten::to_dtype_layout], key=[AutogradCPU]
[call] op=[aten::_to_copy], key=[AutogradCPU]
[redispatch] op=[aten::_to_copy], key=[BackendSelect]
[redispatch] op=[aten::_to_copy], key=[CUDA]
[call] op=[aten::empty_strided], key=[BackendSelect]
[redispatch] op=[aten::empty_strided], key=[CUDA]
[call] op=[aten::copy_], key=[CUDA]
[call] op=[aten::item], key=[AutogradCUDA]
[call] op=[aten::_local_scalar_dense], key=[AutogradCUDA]
[redispatch] op=[aten::_local_scalar_dense], key=[CUDA]
Out [12]: -0.9774109721183777
```

add.Tensor

参考信息

只能使用 TileCore 实现，需要

按道理应该注册到 CompositImplicitAutograd，但是却直接注册到了 CUDA；看起来像是针对 tag 为 ufunc_inner_loop 的 op 专门在 gen.py 中做了处理：

```
for g in structured_native_functions:
    if not g.out.ufunc_inner_loop or not is_ufunc_dispatch_key(dispatch_key):
        continue
    name = g.functional.func.name.name
    if dispatch_key is DispatchKey.CPU:
        assert fm is cpu_fm
        fm.write_with_template(
            f"UfuncCPU_{name}.cpp",
            "UfuncCPU.cpp",
            lambda: {
                "meta_declaration": compute_meta_function_declaration(g),
                "native_declaration": dest.compute_native_function_declaration(
                    g, backend_indices[dispatch_key]
                ),
                "native_definitions": dest.compute_ufunc_cpu(g),
            },
        )
        cpu_vec_fm.write_with_template(
            f"UfuncCPUKernel_{name}.cpp",
            "UfuncCPUKernel.cpp",
            lambda: {
                "name": name,
                "native_definitions": dest.compute_ufunc_cpu_kernel(g),
            },
        )
    elif dispatch_key is DispatchKey.CUDA:
        cuda_headers = "#include <ATen/native/cuda/Loops.cuh>"
        if rocm:
            cuda_headers = "#include <ATen/native/hip/Loops.cuh>"
        fm.write_with_template(
            f"UfuncCUDA_{name}.cu",
            "UfuncCUDA.cu",
            lambda: {
                "name": name,
                "cuda_headers": cuda_headers,
                "meta_declaration": compute_meta_function_declaration(g),
                "native_declaration": dest.compute_native_function_declaration(
                    g, backend_indices[dispatch_key]
                ),
                "native_definitions": dest.compute_ufunc_cuda(g),
            },
        )
    else:
        raise AssertionError(f"unrecognized {dispatch_key} for ufunc")
```

参考信息

OP List 链接：<https://siorigin.feishu.cn/record/Hy3BrJkVel59HcGDWycyJ7Wnfg>

Pytorch 官方文档参考：<https://pytorch.org/docs/stable/generated/torch.add>

具体对应的 Op 参考：THPVariable_add

Cuda 实现参考：

- 源代码里没有对应实现，为自动生成，生成后在 build/aten/src/ATen/UfuncCUDA_add.cu 文件
- 直接调用封装好的 gpu_kernel，支持传入一个 Scalar + 一个 Tensor 或者两个 Tensor
- Kernel 需处理 self/other 其中一个为 scalar 的情况
- gpu_kernel 中相关的优化参考 neg.out 中的 *Performance*

Function 接口

代码块

```

1 // schema: add.out(Tensor self, Tensor other, *, Scalar alpha=1, Tensor(a!)
  out) -> Tensor(a!)
2 // schema: add_.Tensor(Tensor(a!) self, Tensor other, *, Scalar alpha=1) ->
  Tensor(a!)
3 // schema: add.Tensor(Tensor self, Tensor other, *, Scalar alpha=1) -> Tensor
4
5 template<typename T>
6 void add_out(
7     const T* self,
8     const int_t self_ndims,
9     const int64_t* self_sizes,
10    const T* other,
11    const int64_t other_ndims,
12    const int64_t* other_sizes,
13    T alpha = 1,
14    T* out,
15    const int64_t out_ndims,
16    const int64_t* out_sizes,
17    const int64_t* element_sizes = nullptr)

```

ISSUE

Dtype 支持

bfloat16 目前缺少对应的指令

Mxformat

`uint8_t, int8_t, int16_t, int, int64_t, float, double, c10::complex<float>, c10::complex<double>, c10::Half, bool, at::BFloat16, c10::complex<c10::Half>`

非 contiguous 的是否能支持（如果支持需要修改 Function 接口）

Memory 是否需要 Align, Align 到多少

目前需要软件 pad 到 tile format

需要优化

目前对于需要 broadcast 的情况, 软件来实现 expand

对于 outer dim broadcast 需要优化

bmm,bmm.out

参考信息

Pytorch 官方文档参考: <https://pytorch.org/docs/stable/generated/torch.bmm>

具体对应的 Op 参考: THPVariable_bmm

Function 接口

代码块

```
1 // schema: bmm(Tensor self, Tensor mat2) -> Tensor
2 // schema: bmm.out(Tensor self, Tensor mat2, *, Tensor(a!) out) -> Tensor(a!)
3
4 // 实现接口参考: https://docs.nvidia.com/cuda/cublas/#cublasgemmstridedbatchedex
```

any.all_out, any, any.dims, any.dim, any.out, any.dimname

对应 pytorch 接口: <https://pytorch.org/docs/stable/generated/torch.any.html>, 具体对应的 Op 参考 THPVariable_any

参考实现:

<https://github.com/pytorch/pytorch/blob/3f069e7679588d5ee4b1d5b2492ca0e20f9320b5/aten/src/ATen/native/ReduceOps.cpp#L1611>

cuda 中 any 和 all 是放在一起实现的

代码块

```
1 // 1. 基本用法 - 全局规约
2 {
3     // any - 任意元素为True则返回True
4     auto input = torch::tensor({1, 0, 1});
5     auto result = torch::any(input); // tensor(True)
6
7     // all - 所有元素为True才返回True
8     input = torch::tensor({1, 1, 1});
9     result = torch::all(input); // tensor(True)
10 }
11
12 // 2. 指定维度规约
13 {
14     // 2x3 tensor
15     auto input = torch::tensor({{1, 0, 1},
16                                 {0, 1, 1}});
17
18     // any 在维度0上规约
19     auto result = torch::any(input, 0); // [1, 1, 1]
20
21     // all 在维度1上规约
22     result = torch::all(input, 1); // [0, 0]
23 }
24
25 // 3. 保持维度 (keepdim=true)
```

```
26 {
27     auto input = torch::tensor({{1, 0, 1},
28                                 {0, 1, 1}});
29
30     // any 在维度0上规约, 保持维度
31     auto result = torch::any(input, 0, true); // [[1, 1, 1]]
32
33     // all 在维度1上规约, 保持维度
34     result = torch::all(input, 1, true); // [[0], [0]]
35 }
36
37 // 4. uint8类型处理
38 {
39     auto input = torch::tensor({1, 0, 1}, torch::kUInt8);
40
41     // 对uint8类型, 输出保持uint8
42     auto result = torch::any(input); // tensor(1, dtype=uint8)
43     result = torch::all(input); // tensor(0, dtype=uint8)
44 }
45
46 // 5. 使用out参数
47 {
48     auto input = torch::tensor({1, 0, 1});
49     auto out = torch::empty({}, torch::kBool);
50
51     // any with out
52     torch::any_out(out, input);
53
54     // all with out
55     torch::all_out(out, input);
56 }
57
58 // 6. 多维度规约
59 {
60     // 2x3x2 tensor
61     auto input = torch::ones({2, 3, 2});
62     std::vector<int64_t> dims = {0, 2};
63
64     // 在维度0和2上规约
65     auto result = torch::any(input, dims); // [[1, 1, 1]]
66     result = torch::all(input, dims); // [1, 1, 1]
67 }
68
69 // 7. 空张量处理
70 {
71     auto input = torch::tensor({});
72 }
```

```

73 // 空张量的any返回false
74 auto result = torch::any(input); // tensor(False)
75
76 // 空张量的all返回true
77 result = torch::all(input); // tensor(True)
78 }

```

any.all_out,any

CUDA 参考实现 any_all_out

代码块

```

1 // schema: any.all_out(Tensor
  self, *, Tensor(a!) out) ->
  Tensor(a!)
2 // schema: all(Tensor self) ->
  Tensor

```

```

tensor([[ 0.6800, -1.9964,  0.2954],
        [-1.5758, -0.9457,  1.0380],
        [-1.4764,  0.5204, -0.0701]], device='cuda:0')
In [32]: a.any()
[call] op=[aten::any], key=[AutogradCUDA]
[redispatch] op=[aten::any], key=[CUDA]
[call] op=[aten::as_strided], key=[CUDA]
Out[32]: [call] op=[aten::reshape], key=[AutogradCUDA]
[call] op=[aten::view], key=[AutogradCUDA]
[redispatch] op=[aten::view], key=[ADInplaceOrView]
[redispatch] op=[aten::view], key=[CUDA]
[call] op=[aten::unbind.int], key=[AutogradCUDA]
[redispatch] op=[aten::unbind.int], key=[ADInplaceOrView]
[redispatch] op=[aten::unbind.int], key=[CUDA]
[call] op=[aten::select.int], key=[CUDA]
[call] op=[aten::as_strided], key=[CUDA]
[call] op=[aten::item], key=[AutogradCUDA]
[call] op=[aten::_local_scalar_dense], key=[AutogradCUDA]
[redispatch] op=[aten::_local_scalar_dense], key=[CUDA]
[call] op=[aten::item], key=[AutogradCUDA]
[call] op=[aten::_local_scalar_dense], key=[AutogradCUDA]
[redispatch] op=[aten::_local_scalar_dense], key=[CUDA]
tensor(True, device='cuda:0')

```

any.dim,any.out

CUDA 参考实现 any_out

代码块

```

1 // schema: any.dim(Tensor self,
  int dim, bool keepdim=False) ->
  Tensor
2 // schema: any.out(Tensor self,
  int dim, bool keepdim=False, *,
  Tensor(a!) out) -> Tensor(a!)

```

```

In [45]: a.shape
Out[45]: torch.Size([2, 3, 4])
In [46]: a.any(dim=0)
[call] op=[aten::any.dim], key=[AutogradCUDA]
[redispatch] op=[aten::any.dim], key=[CUDA]
[call] op=[aten::as_strided], key=[CUDA]
tensor([[True, True, True, True],
        [True, True, True, True],
        [True, True, True, True]], device='cuda:0')

```

all.dimname, all.dimname_out, all.dims, all.dims_out,

官方文档

<https://pytorch.org/docs/stable/generated/torch.all.html>

代码块

```

1 //all.dimname_out(Tensor self, Dimname dim, bool keepdim=False, *, Tensor(a!)
  out) -> Tensor(a!)
2 //all.dimname(Tensor self, Dimname dim, bool keepdim=False) -> Tensor
3 //all.dims(Tensor self, int[]? dim=None, bool keepdim=False) -> Tensor
4 //all.dims_out(Tensor self, int[]? dim=None, bool keepdim=False, *, Tensor(a!)
  out) -> Tensor(a!)
5
6

```

all.dim, all.out

代码块

```

1 //all.dims(Tensor self, int[]? dim=None, bool keepdim=False) -> Tensor
2 //all.dims_out(Tensor self, int[]? dim=None, bool keepdim=False, *, Tensor(a!)
  out) -> Tensor(a!)

```

Function 接口

代码块

```

1 /**
2  * @brief Compute logical any/all reduction over dimensions of input tensor
3  *
4  * @param dev          Device handle for CUDA device management
5  * @param out          Output tensor (pre-allocated, bool or uint8)
6  * @param self         Input tensor
7  * @param dims         Dimensions to reduce over (optional)
8  * @param keepdim     Whether the output tensor has dim retained or not
9  * @param dtype       Data type of input tensor
10 * @param is_all      If true compute logical AND, else compute logical OR
11 *
12 * @note Output Type:
13 *   - For uint8 input: output is uint8
14 *   - For all other inputs: output is bool
15 *
16 * @note Dimension Rules:
17 *   - If dim is empty, reduce over all dimensions
18 *   - Multiple dims are allowed
19 *   - Negative dims are wrapped
20 *
21 * @note Examples:
22 *   any: [1,1,0] -> true
23 *   any: [0,0,0] -> false
24 *   all: [1,1,0] -> false
25 *   all: [1,1,1] -> true

```

```

26  *
27  * @return int 0 on success, error code otherwise
28  */
29  void any_all_fn(
30      SiDeviceHandle dev,           // Device handle
31      void *out,                   // Output tensor (bool or uint8)
32      const void *self,            // Input tensor
33      const int64_t* dims,         // Dimensions to reduce over (optional)
34      bool keepdim,               // Whether to keep reduced dimensions
35      ScalarType dtype,           // Input data type
36      bool is_all                 // true for all(), false for any()
37  );

```

bitwise_and.Tensor, bitwise_and.Tensor_out, bitwise_and_.Tensor, bitwise_and.Scalar_out, bitwise_and.Scalar, bitwise_and.Scalar_Tensor, bitwise_and.Scalar_Tensor_out, bitwise_and_.Scalar

参考信息

只能使用 RVV 实现

Pytorch 官方文档参考: https://pytorch.org/docs/stable/generated/torch.bitwise_and.html

具体对应的 Op 参考: THPVariable_bitwise_and_, THPVariable_bitwise_and

支持 broadcastable 的 tensor 计算

对于 int 类型的 tensor, 运算后的结果是二进制补码对应的十进制数

在线补码网站: <https://www.23bei.com/tool/56.html?>

bitwise_and.Tensor, bitwise_and.Tensor_out, bitwise_and_.Tensor

CUDA 实现搜索: bitwise_and_kernel_cuda

Function 接口

建议在 device 和 kernel 可以联合编译后参考 CUDA 代码来实现

代码块

```

1  // schema: bitwise_and.Tensor(Tensor self, Tensor
   other) -> Tensor
2  // schema: bitwise_and_.Tensor(Tensor(a!) self,
   Tensor other) -> Tensor(a!)
3  // schema: bitwise_and.Tensor_out(Tensor self,
   Tensor other, *, Tensor(a!) out) -> Tensor(a!)
4
5  /**
6  nargs: num of output+input tensors

```

```

In [53]: b = torch.randint(0, 10, (3, 2))
Notes   op=[aten::randint_low], key=[BackendSelect]
[call] op=[aten::randint_low], key=[CPU]
[call] op=[aten::empty_memory_format], key=[BackendSelect]
[redispatch] op=[aten::empty_memory_format], key=[CPU]
[call] op=[aten::random_from], key=[CPU]

In [54]: b & b
[call] op=[aten::_and_.Tensor], key=[AutogradCPU]
[call] op=[aten:bitwise_and.Tensor], key=[AutogradCPU]
[redispatchBoxed] op=[aten:bitwise_and.Tensor], key=[CPU]

```

```

7  args: output+input tensors' ptr
8  dims: output tensor dims
9  sizes: output tensor sizes
10 strides: strides of nargs tensors, 0 means it's a
    broadcasted dim
11 element_sizes: element_size of nargs tensor; if
    element_sizes is nullptr, then the strides will
    be in bytes, otherwise the strides will be in #
    of elements.
12  **/
13 void fn(dev, int nargs, const void* const* args,
    int dims, const int64_t* sizes, const int64_t*
    const* strides, const int64_t*
    element_sizes=nullptr)
14
15 // 可选的形式
16 void fnv2(dev, const void* out, const void* a,
    const void* b, int dims, const int64_t* sizes,
    const int64_t* strides_out, const int64_t*
    strides_a, const int64_t* strides_b, const
    int64_t* element_sizes=nullptr)
17 //updated
18 void fnv2(dev, const void* out, const void* self,
    const void* other, int out_ndims, const int64_t*
    out_sizes, const int64_t* out_strides, const
    int64_t* self_strides, const int64_t*
    other_strides, const int64_t*
    element_sizes=nullptr)

```

ISSUE

Dtype 支持

float32

float16

bfloat16 目前缺少对应的指令

mxformat

非 contiguous 的是否能支持

Memory 是否需要 Align, Align 到多少

bitwise_and.Scalar_out

通过 `at::native::bitwise_and_out` 直接将 scalar 转换为 tensor 的方式, 调用 tensor 的实现版本

```
Tensor& bitwise_and_out(const Tensor& self, const Scalar& other, Tensor& result) {
    return at::bitwise_and_out(result, self, wrapped_scalar_tensor(other));
}
```

bitwise_and.Scalar, bitwise_and.Scalar_Tensor, bitwise_and.Scalar_Tensor_out

通过 `at::native::bitwise_and` 直接将 scalar 转换为 tensor 的方式，调用 tensor 的实现版本

```
Tensor bitwise_and(const Tensor& self, const Scalar& other) {
    return at::bitwise_and(self, wrapped_scalar_tensor(other));
}
```

```
Tensor bitwise_and(const Scalar& self, const Tensor& other) {
    return at::bitwise_and(wrapped_scalar_tensor(self), other);
}
```

bitwise_and_.Scalar

通过 `at::native::bitwise_and_` 直接将函数调用转换为 `tensor.bitwise_and_` 一个 tensor 的方式实现

```
Tensor& bitwise_and_(Tensor& self, const Scalar& other) {
    return self.bitwise_and_(wrapped_scalar_tensor(other));
}
```

bitwise_not, bitwise_not_, bitwise_not.out

参考信息

只能使用 RVV 实现

https://pytorch.org/docs/stable/generated/torch.bitwise_not

具体对应的 Op 参考 `THPVariable_bitwise_not`

对于 int 类型的 tensor，运算后的结果是二进制补码对应的十进制数，每个 tensor 中的数单独按照大小对应的 8bit/16bit/32bit/64bit 来计算，例如对 `a = tensor ([1, 65536])`，`~a = tensor ([-2, -65537])`

在线补码网站：<https://www.23bei.com/tool/56.html?>

Function 接口

建议在 device 和 kernel 可以联合编译后参考 CUDA 代码来实现

代码块

```
1 // schema: bitwise_not(Tensor self) ->
   Tensor
```

```
tensor([[[8, 5],
         [4, 8],
         [6, 6]])
```

```
In [91]: (b, ~b)
[call] op=[aten::bitwise_not], key=[AutogradCPU]
[redispatchBoxed] op=[aten::bitwise_not], key=[CPU]
```

```

2 // schema: bitwise_not_(Tensor(a!) self) ->
  Tensor(a!)
3 // schema: bitwise_not.out(Tensor self, *,
  Tensor(a!) out) -> Tensor(a!)
4
5 /**
6 out: output tensor
7 a: input tensor
8 dims: output tensor dims
9 sizes: output tensor sizes
10 strides: strides of output tensor
11 element_sizes: element_size of nargs
  tensor; if element_sizes is nullptr, then
  the strides will be in bytes, otherwise the
  strides will be in # of elements.
12 */
13 void fn(dev, const void* out, const void*
  a, int dims, const int64_t* sizes, const
  int64_t* strides, const int64_t*
  element_sizes=nullptr)
14 //updated
15 void fn(dev, const void* out, const void*
  self, int out_ndims, const int64_t*
  out_sizes, const int64_t* out_strides,
  const int64_t* element_sizes=nullptr)

```

```

(tensor([[8, 5],
        [4, 8],
        [6, 6]]),
tensor([[ -9, -6],
        [-5, -9],
        [-7, -7]]))

```

ISSUE

Dtype 支持

float32

float16

bfloat16 目前缺少对应的指令

mxformat

非 contiguous 的是否能支持（如果支持需要修改 Function 接口）

Memory 是否需要 Align, Align 到多少

bitwise_or.Tensor, bitwise_or_.Tensor, bitwise_or.Tensor_out

参考信息

只能使用 RVV 实现

https://pytorch.org/docs/stable/generated/torch.bitwise_or

支持 broadcastable 的 tensor 计算

具体对应的 Op 参考 THPVariable_bitwise_orq

代码块

```
1 // schema: bitwise_or.Tensor(Tensor self,
  // Tensor other) -> Tensor
2 // schema: bitwise_or_.Tensor(Tensor(a!)
  // self, Tensor other) -> Tensor(a!)
3 // schema: bitwise_or.Tensor_out(Tensor
  // self, Tensor other, *, Tensor(a!) out) ->
  // Tensor(a!)
```

```
In [149]: (c, d, c | d)
[call] op=[aten::_or_.Tensor], key=[AutogradCPU]
[call] op=[aten::bitwise_or.Tensor], key=[AutogradCPU]
[redispatchBoxed] op=[aten::bitwise_or.Tensor], key=[CPU]
(tensor([ 1, 65536]), tensor([ 22, 65536]), tensor([ 23, 65536]))
```

Function 接口

建议在 device 和 kernel 可以联合编译后参考 CUDA 代码来实现

代码块

```
1 // schema: bitwise_and.Tensor(Tensor self, Tensor
  // other) -> Tensor
2 // schema: bitwise_and_.Tensor(Tensor(a!) self,
  // Tensor other) -> Tensor(a!)
3 // schema: bitwise_and.Tensor_out(Tensor self,
  // Tensor other, *, Tensor(a!) out) -> Tensor(a!)
4
5 /**
6 nargs: num of output+input tensors
7 args: output+input tensors' ptr
8 dims: output tensor dims
9 sizes: output tensor sizes
10 strides: strides of nargs tensors, 0 means it's a
    broadcasted dim
11 element_sizes: element_size of nargs tensor; if
    element_sizes is nullptr, then the strides will
    be in bytes, otherwise the strides will be in #
    of elements.
12 **/
13 void fn(dev, int nargs, const void* const* args,
  int dims, const int64_t* sizes, const int64_t*
  const* strides, const int64_t*
  element_sizes=nullptr)
14
15 // 可选的形式
16 void fnv2(dev, const void* out, const void* a,
  const void* b, int dims, const int64_t* sizes,
  const int64_t* strides_out, const int64_t*
```

```
In [53]: b = torch.randint(0,10,(3,2))
Notes op=[aten::randint.Low], key=[BackendSelect]
[redispatch] op=[aten::randint.Low], key=[CPU]
[call] op=[aten::empty.memory_format], key=[BackendSelect]
[redispatch] op=[aten::empty.memory_format], key=[CPU]
[call] op=[aten::random_from], key=[CPU]
In [54]: b & b
[call] op=[aten::_and_.Tensor], key=[AutogradCPU]
[call] op=[aten::bitwise_and.Tensor], key=[AutogradCPU]
[redispatchBoxed] op=[aten::bitwise_and.Tensor], key=[CPU]
```

```

strides_a, const int64_t* strides_b, const
int64_t* element_sizes=nullptr)
17 //updated
18 void fnv2(dev, const void* out, const void* self,
const void* other, int out_ndims, const int64_t*
out_sizes, const int64_t* out_strides, const
int64_t* self_strides, const int64_t*
other_strides, const int64_t*
element_sizes=nullptr)

```

ISSUE

Dtype 支持

float32

float16

bfloat16 目前缺少对应的指令

mxformat

非 contiguous 的是否能支持

Memory 是否需要 Align, Align 到多少

eq.Tensor,eq.Tensor_out,eq_.Tensor

参考信息

只能使用 RVV 实现

<https://pytorch.org/docs/stable/generated/torch.eq>

支持 broadcastable 的 tensor 计算

具体对应的 Op 参考 THPVariable_eq

代码块

```

1 // schema: eq_.Tensor(Tensor(a!) self,
Tensor other) -> Tensor(a!)
2 // schema: eq.Tensor_out(Tensor self,
Tensor other, *, Tensor(a!) out) ->
Tensor(a!)
3 // schema: eq.Tensor(Tensor self, Tensor
other) -> Tensor

```

```

In [165]: (c, d, c==d)
[call] op=[aten::eq.Tensor], key=[AutogradCPU]
[redispatch] op=[aten::eq.Tensor], key=[CPU]

```

```

(tensor([[ 1, 65536],
          [ 2,    3]]),
tensor([ 22, 65536]),
tensor([[False, True],
        [False, False]]))

```

Function 接口

建议在 device 和 kernel 可以联合编译后参考 CUDA 代码来实现

代码块

```
1 // schema: bitwise_and.Tensor(Tensor self, Tensor
  other) -> Tensor
2 // schema: bitwise_and_.Tensor(Tensor(a!) self,
  Tensor other) -> Tensor(a!)
3 // schema: bitwise_and.Tensor_out(Tensor self,
  Tensor other, *, Tensor(a!) out) -> Tensor(a!)
4
5 /**
6  nargs: num of output+input tensors
7  args: output+input tensors' ptr
8  dims: output tensor dims
9  sizes: output tensor sizes
10 strides: strides of nargs tensors, 0 means it's a
  broadcasted dim
11 element_sizes: element_size of nargs tensor; if
  element_sizes is nullptr, then the strides will
  be in bytes, otherwise the strides will be in #
  of elements.
12 */
13 void fn(dev, int nargs, const void* const* args,
  int dims, const int64_t* sizes, const int64_t*
  const* strides, const int64_t*
  element_sizes=nullptr)
14
15 // 可选的形式
16 void fnv2(dev, const void* out, const void* a,
  const void* b, int dims, const int64_t* sizes,
  const int64_t* strides_out, const int64_t*
  strides_a, const int64_t* strides_b, const
  int64_t* element_sizes=nullptr)
17 // updated
18 void fnv2(dev, const void* out, const void* self,
  const void* other, int out_ndim, const int64_t*
  out_sizes, const int64_t* out_strides, const
  int64_t* self_strides, const int64_t*
  other_strides, const int64_t*
  element_sizes=nullptr)
```

```
In [53]: b = torch.randint(0,10,(3,2))
Notes   op=[aten::randint_low], key=[BackendSelect]
[redispatch] op=[aten::randint_low], key=[CPU]
[call] op=[aten::empty_memory_format], key=[BackendSelect]
[redispatch] op=[aten::empty_memory_format], key=[CPU]
[call] op=[aten::random_from], key=[CPU]

In [54]: b & b
[call] op=[aten::...and...Tensor], key=[AutogradCPU]
[call] op=[aten::bitwise_and.Tensor], key=[AutogradCPU]
[redispatchBoxed] op=[aten::bitwise_and.Tensor], key=[CPU]
```

首个版本开发说明

- (1) 数制支持上, 需要支持类型有: int32_t, bfloat16, float32
- (2) 初版只需要支持输入2为scalar的情况
- (3) 由于硬件限制, 对于int64_t类型, 由框架将数据进行64-32转换

内部接口调用形式:

代码块

```
1  /// @brief Computes element-wise equality
2  /// @tparam T Datatype support [int32_t, sifmt:float32, sifmt::bfloat16]
3  /// @param[in] d_in0 Input tensor 0 with shape [dim0_size, dim1_size]
4  /// @param[in] d_in1 Input tensor 1 with shape [dim0_size, dim1_size] OR scalar
5  /// @param[out] d_out Output tensor with shape [dim0_size, dim1_size]
6  /// @param[in] dim0_size Tensor shape param 0
7  /// @param[in] dim1_size Tensor shape param 1
8  /// @param[in] is_d_in1_scalar Indicates tensor 1 is scalar OR tensor
9  /// @param[in] op_type Indicates OP is EQUAL(0) OR NOT EQUAL(1)
10 template <typename T>
11 void eq(void *d_in0, void *d_in1, void *d_out, int dim0_size, int dim1_size,
        bool is_d_in1_scalar = true, int op_type = 0);
```

ISSUE

Dtype 支持

float32

float16

bfloat16 目前缺少对应的指令

mxformat

非 contiguous 的是否能支持 (如果支持需要修改 Function 接口)

Memory 是否需要 Align, Align 到多少

ge.Scalar,ge.Scalar_out,ge_.Scalar,ge.Tensor,ge.Tensor_out,ge_.Tensor

参考信息

只能使用 RVV 实现

<https://pytorch.org/docs/stable/generated/torch.ge>

支持 broadcastable 的 tensor 计算

具体对应的 Op 参考 THPVariable_ge

ge.Scalar,ge.Scalar_out,ge_.Scalar

```
In [174]: (c, c > 2)
[call] op=[aten::gt.Scalar], key=[AutogradCPU]
[redispatch] op=[aten::gt.Scalar], key=[CPU]
```

代码块

```
1  // schema: ge.Scalar(Tensor self, Scalar
    other) -> Tensor
```

```
(tensor([[ 1, 65536],
          [ 2,   3]]),
 tensor([[False, True],
          [False, True]]))
```

```

2 // schema: ge.Scalar_out(Tensor self,
  Scalar other, *, Tensor(a!) out) ->
  Tensor(a!)
3 // schema: ge_.Scalar(Tensor(a!) self,
  Scalar other) -> Tensor(a!)

```

ge.Tensor,ge.Tensor_out,ge_.Tensor

代码块

```

1 // schema: ge_.Tensor(Tensor(a!) self,
  Tensor other) -> Tensor(a!)
2 // schema: ge.Tensor_out(Tensor self,
  Tensor other, *, Tensor(a!) out) ->
  Tensor(a!)
3 // schema: ge.Tensor(Tensor self, Tensor
  other) -> Tensor

```

```

In [179]: (c, d, c>d)
[call] op=[aten::gt.Tensor], key=[AutogradCUDA]
[redispatch] op=[aten::gt.Tensor], key=[CUDA]

```

```

(tensor([[ 1, 65536],
          [ 2,   3]], device='cuda:0'),
 tensor([ 22, 65536], device='cuda:0'),
 tensor([[False, False],
         [False, False]], device='cuda:0'))

```

Function 接口

建议在 device 和 kernel 可以联合编译后参考 CUDA 代码来实现

代码块

```

1 // schema: bitwise_and.Tensor(Tensor self, Tensor
  other) -> Tensor
2 // schema: bitwise_and_.Tensor(Tensor(a!) self,
  Tensor other) -> Tensor(a!)
3 // schema: bitwise_and.Tensor_out(Tensor self,
  Tensor other, *, Tensor(a!) out) -> Tensor(a!)
4
5 /**
6  nargs: num of output+input tensors
7  args: output+input tensors' ptr
8  dims: output tensor dims
9  sizes: output tensor sizes
10 strides: strides of nargs tensors, 0 means it's a
    broadcasted dim
11 element_sizes: element_size of nargs tensor; if
    element_sizes is nullptr, then the strides will
    be in bytes, otherwise the strides will be in #
    of elements.
12 **/
13 void fn(dev, int nargs, const void* const* args,
    int dims, const int64_t* sizes, const int64_t*

```

```

In [53]: b = torch.randint(0,10,(3,2))
Notes: op=[aten::randint.low], key=[BackendSelect]
[redispatch] op=[aten::randint.low], key=[CPU]
[call] op=[aten::empty.memory_format], key=[BackendSelect]
[redispatch] op=[aten::empty.memory_format], key=[CPU]
[call] op=[aten::random_.from], key=[CPU]
In [54]: b & b
[call] op=[aten::_.and_.Tensor], key=[AutogradCPU]
[call] op=[aten::bitwise_and.Tensor], key=[AutogradCPU]
[redispatchBoxed] op=[aten::bitwise_and.Tensor], key=[CPU]

```

```

14     const* strides, const int64_t*
15     element_sizes=nullptr)
16 // 可选的形式
17 void fnv2(dev, const void* out, const void* a,
18 const void* b, int dims, const int64_t* sizes,
19 const int64_t* strides_out, const int64_t*
20 strides_a, const int64_t* strides_b, const
21 int64_t* element_sizes=nullptr)
22 // updated
23 void fnv2(dev, const void* out, const void* self,
24 const void* other, int out_ndims, const int64_t*
25 out_sizes, const int64_t* out_strides, const
26 int64_t* self_strides, const int64_t*
27 other_strides, const int64_t*
28 element_sizes=nullptr)

```

ISSUE

Dtype 支持

float32

float16

bfloat16 目前缺少对应的指令

mxformat

非 contiguous 的是否能支持

Memory 是否需要 Align, Align 到多少

mul.Tensor, mul_.Tensor, mul.out

参考信息

Pytorch 官方文档参考: <https://pytorch.org/docs/stable/generated/torch.mul.html>

cuda 实现参考:

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/BinaryMulKernel.cu>
- 直接调用封装好的 gpu_kernel, 支持传入一个 Scalar + 一个 Tensor 或者两个 Tensor
- Kernel 需处理 self/other 其中一个为 scalar 的情况
- gpu_kernel 中相关的优化参考 neg.out 中的 *Performance*

Function 接口

```

1 代码块/ schema: mul.out(Tensor self, Tensor other, *, Tensor(a!) out) -> Tensor(a!)
2  // schema: mul.Tensor(Tensor self, Tensor other) -> Tensor
3  // schema: mul_.Tensor(Tensor(a!) self, Tensor other) -> Tensor(a!)
4
5  template<typename T>
6  void mul_out(
7      const T* self,
8      const int64_t self_ndims,
9      const int64_t* self_sizes,
10     const T* other,
11     const int64_t other_ndims,
12     const int64_t* other_sizes,
13     T* out,
14     const int64_t out_ndims,
15     const int64_t* out_sizes,
16     const int64_t* element_sizes = nullptr)

```

ISSUE

Dtype 支持

`uint8_t, int8_t, int16_t, int, int64_t, float, double, c10::complex<float>, c10::complex<double>, c10::Half, bool, at::BFloat16, c10::complex<c10::Half>`

gather, gather.out

参考信息

OP List 链接: <https://siorigin.feishu.cn/record/L6XorIxsxeO1g8c8DdFcqnUtnbd>

Pytorch 官方文档参考: <https://pytorch.org/docs/stable/generated/torch.gather.html>

Output tensor 的 shape 和 index 的 tensor shape 一致

实现参考:

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/ScatterGatherKernel.cu#L219>
- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/ScatterGatherKernel.cu#L117>

Function 接口

代码块

```

1  // schema: gather.out(Tensor self, int dim, Tensor index, *, bool
   sparse_grad=False, Tensor(a!) out) -> Tensor(a!)

```

```

2 // schema: gather(Tensor self, int dim, Tensor index, *, bool
  sparse_grad=False) -> Tensor
3
4 /**
5  * @brief Performs a gather operation on the input tensor `self` using indices
  from the `index` tensor, and stores the result in the output tensor `out`.
6  *
7  * @param dev Device on which the operation is performed.
8  * @param out Pointer to the output tensor, which is modified in
  place.
9  * @param dims Common number of dimensions shared by `out`, `self`,
  and `index`.
10 * @param out_sizes Sizes of each dimension for both the output tensor
  and the index tensor.
11 * @param out_strides Strides of each dimension of the output tensor.
12 * @param self Pointer to the input tensor.
13 * @param self_sizes Sizes of each dimension of the input tensor.
14 * @param self_strides Strides of each dimension of the input tensor.
15 * @param dim Dimension along which elements are gathered from
  `self`.
16 * @param index Pointer to the index tensor, which is a multi-
  dimensional tensor containing indices.
17 * @param index_stride Strides of each dimension of the index tensor.
18 * @param element_sizes Optional pointer to the element sizes of each tensor
  (if nullptr, strides are in bytes).
19 */
20 void fn(dev, void* out, int dims, const int64_t* out_sizes, const int64_t*
  out_strides, const void* self, const int64_t* self_sizes, const int64_t*
  self_strides, int64_t dim, const void* index, const int64_t* index_stride,
  const int64_t* element_sizes=nullptr)
21 //updated
22 void fn(dev, void* out, int ndims, const int64_t* out_sizes, const int64_t*
  out_strides, const void* self, const int64_t* self_sizes, const int64_t*
  self_strides, int64_t dim, const void* index, const int64_t* index_strides,
  const int64_t* element_sizes=nullptr)

```

一些例子

```

代码块
1 # 参考pytorch文档https://pytorch.org/docs/stable/generated/torch.gather.html, 这
  里index=torch.tensor([[0, 0], [1, 0]]), shape为[2,2], 所以ij分别为
2 # 0,0; 0,1; 1,0; 1,1
3 t = torch.tensor([[1, 2], [3, 4]])
4 torch.gather(t, 1, torch.tensor([[0, 0], [1, 0]]))
5

```

```

6 # Out[33]:
7 # tensor([[1, 1],
8 #         [4, 3]])
9
10 # 这里index=torch.tensor([[[1,1]],[[0,1]]]), shape为[2,1,2], 所以ijk分别为
11 # 0,0,0; 0,0,1; 1,0,0; 1,0,1
12 t = torch.arange(8).reshape(2,2,2)
13 torch.gather(t, 1, torch.tensor([[[1,1]],[[0,1]]]))
14
15 # Out[32]:
16 # tensor([[[2, 3]],
17 #         [[4, 7]])
18 #

```

首个版本开发说明

- (1) 支持输入输出2D tensor的数据为float32类型，index内的数据类型为int64。
- (2) dim默认为-1，可支持0/1。
- (3) 维度信息包括input_dim0,input_dim1,index_dim0,index_dim1。output维度与index相同，不需要提供。

内部接口调用形式：

代码块

```

1 template <typename T>
2 void gather(void *d_in, void *d_index, void *d_out, int dim=-1,
3            int in_dim0, int in_dim1,
4            int index_dim0, int
5            index_dim1);

```

ISSUE

Dtype 支持

- AT_DISPATCH_ALL_TYPES_AND_COMPLEX: uint8_t, int8_t, int16_t, int, int64_t, float, double, c10::complex<float>, c10::complex<double>
- c10::Half, bool, at::BFloat16
- MX Format?

index_select

参考信息

OP List 链接: <https://siorigin.feishu.cn/record/XkXcrPomqe6yO0cma8XclHYjnlf>

Pytorch 官方文档参考: https://pytorch.org/docs/stable/generated/torch.index_select.html

实现参考:

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/Indexing.cu#L1350>
- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/Indexing.cu#L1240>
- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/Indexing.cu#L1283>

example

代码块

```
1 >>> x = torch.randn(3, 4)
2 >>> x
3 tensor([[ 1.1880,  0.2477, -1.3225,  1.4690],
4         [ 2.2300, -0.1451, -0.1112, -0.5504],
5         [ 0.1272, -0.8875,  0.2266, -0.3423]])
6 >>> indices = torch.tensor([0, 2])
7 >>> res1 = torch.index_select(x, 0, indices)
8 >>> res1
9 tensor([[ 1.1880,  0.2477, -1.3225,  1.4690],
10        [ 0.1272, -0.8875,  0.2266, -0.3423]])
11 >>> res2 = torch.index_select(x, 1, indices)
12 >>> res2
13 tensor([[ 1.1880, -1.3225],
14         [ 2.2300, -0.1112],
15         [ 0.1272,  0.2266]])
```

Function 接口

代码块

```
1 // schema: index_select(Tensor self, int dim, Tensor index) -> Tensor
2
3 /**
4  * @brief Performs an index_select operation on the input tensor `self` using
5  * the provided `index` tensor, and stores the result in the output tensor `out`.
6  *
7  * @param dev Device on which the operation is performed.
8  * @param out Pointer to the output tensor.
9  * @param dims Number of dimensions in the input/output tensor.
10 * @param out_sizes Sizes of each dimension of the output tensor.
11 * @param out_strides Strides of each dimension of the output tensor.
```

```

11 * @param self Pointer to the input tensor.
12 * @param self_sizes Sizes of each dimension of the input tensor.
13 * @param self_strides Strides of each dimension of the input tensor.
14 * @param dim Dimension along which index_select is performed.
15 * @param index Pointer to the one-dimensional tensor containing
    indices for indexing.
16 * @param index_size Number of elements in the index tensor.
17 * @param index_stride Stride of the index tensor.
18 * @param element_sizes Optional pointer to the element sizes of each tensor
    (if nullptr, strides are in bytes).
19 */
20 void fn(dev, void* out, int dims, const int64_t* out_sizes, const int64_t*
    out_strides, const void* self, const int64_t* self_sizes, const int64_t*
    self_strides, int64_t dim, const void* index, const int64_t index_size, const
    int64_t index_stride, const int64_t* element_sizes=nullptr)
21 //updated
22 void fn(dev, void* out, int ndims, const int64_t* out_sizes, const int64_t*
    out_strides, const void* self, const int64_t* self_sizes, const int64_t*
    self_strides, int64_t dim, const void* index, const int64_t index_sizes, const
    int64_t index_strides, const int64_t* element_sizes=nullptr)

```

ISSUE

Dtype 支持

- AT_ALL_TYPES_AND_COMPLEX: `uint8_t, int8_t, int16_t, int, int64_t, float, double, c10::complex<float>, c10::complex<double>`
- AT_BAREBONES_UNSIGNED_TYPES: `uint16_t, uint32_t, uint64_t`
- `c10::complex<c10::Half>, c10::Half, bool, at::BFloat16`
- MX Format?

silu

unary op, 接口形式同 neg, 输入输出 shape 一致

官方参考文档

- <https://pytorch.org/docs/stable/generated/torch.nn.functional.silu.html>

cuda kernel 参考实现

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/ActivationSiluKernel.cu#L23>

exp 计算需用到 sfu? 会对接口形式有影响吗?

代码块

```

1 Tensor silu(const Tensor & self); // {"schema": "aten::silu(Tensor self) ->
  Tensor", "dispatch": "True", "default": "True"}
2 Tensor & silu_(Tensor & self); // {"schema": "aten::silu_(Tensor(a!) self) ->
  Tensor(a!)", "dispatch": "True", "default": "True"}
3 Tensor & silu_out(const Tensor & self, Tensor & out); // {"schema":
  "aten::silu.out(Tensor self, *, Tensor(a!) out) -> Tensor(a!)", "dispatch":
  "True", "default": "False"}

```

Dtypes

- AT_FLOATING_TYPES: float, double
- AT_COMPLEX_TYPES: c10::complex<float>, c10::complex<double>
- Half, BFloat16: float16, bfloat16

代码块

```

1 AT_DISPATCH_FLOATING_AND_COMPLEX_TYPES_AND2(
2     at::ScalarType::Half,
3     at::ScalarType::BFloat16

```

max, max.dim, max.dim_max

unary op, 接口形式同 neg, 输出 size 为 1

官方参考文档

- 注意 aten::max 需要实现的只是不接收 keep_dim / dim 版本的, 只是单纯的 max_all
- <https://pytorch.org/docs/main/generated/torch.max.html>

cuda kernel 参考实现

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/ReduceMaxValuesKernel.cu#L28>
- cuda 实现会封装一个通用的 gpu_reduce_kernel, 支持多种 reduce 操作, 建议实现类似接口, 其他 op 即可通过统一接口接入
- 不同 reduce op 只需传递不同的 ops (即具体的 reduce 函数, 如 max 为比较函数), ident (reduce 函数的默认值, 如 max 为 scalar_t 类型的最小值)
- 通用接口支持传入要 reduce 的 dim 以及是否 keep_dim

代码块

```

1 template <typename scalar_t, typename out_scalar_t, int vt0=4, typename ops_t,
  typename ident_t=double>

```

```

2 inline void gpu_reduce_kernel(TensorIterator& iter, const ops_t& ops, ident_t
  ident=0,
3                               AccumulationBuffer* acc_buf_ptr=nullptr, int64_t
  base_idx=0)

```

代码块

```

1 Tensor max(const Tensor & self); // {"schema": "aten::max(Tensor self) ->
  Tensor", "dispatch": "True", "default": "False"}
2
3
4 max.dim, max.dim_max用法
5 >>> x = torch.tensor([[1, 2, 3], [4, 5, 6]])
6 >>> values = torch.empty(2, 1, dtype=torch.float32)
7 >>> indices = torch.empty(2, 1, dtype=torch.long)
8 >>> max_values, max_indices = torch.max(x, dim=1, keepdim=True)
9 >>> max_values
10 tensor([[3],
11         [6]])
12 >>> max_indices
13 tensor([[2],
14         [2]])
15
16 >>> torch.max(x, dim=1, keepdim=True, out = (max_indices, max_values))
17 torch.return_types.max_out(
18 values=tensor([[3],
19               [6]]),
20 indices=tensor([[2],
21               [2]]))

```

Dtypes

- AT_ALL_TYPES: uint8_t, int8_t, int16_t, int, int64_t, float, double
- Half, Bool, BFloat16: float16, bool, bfloat16

代码块

```

1 AT_DISPATCH_ALL_TYPES_AND3(
2     kBFloat16, kHalf, kBool

```

maximum, maximum.out

Binary op, element-wise. 两个tensor逐元素比较大小并返回。torch.max(input, other)会调用torch.maximum接口。

官方文档参考

- <https://docs.pytorch.org/docs/stable/generated/torch.maximum.html#torch.maximum>

接口（实现maximum_out即可，maximum继承于maximum_out）

代码块

```
1 Tensor maximum(const Tensor & self, const Tensor & other); // {"schema":  
  "aten::maximum(Tensor self, Tensor other) -> Tensor", "dispatch": "True",  
  "default": "True"}  
2 Tensor & maximum_out(const Tensor & self, const Tensor & other, Tensor & out);  
  // {"schema": "aten::maximum.out(Tensor self, Tensor other, *, Tensor(a!) out)  
  -> Tensor(a!)", "dispatch": "True", "default": "False"}
```

Cuda kernel参考实现

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/MaxMinElementwiseKernel.cu#L14>

如果是Nan值，直接返回，通过a!=a来判断是否Nan值。

Dtypes（不需要支持complex类型，其他都支持）

1. Bool类型
2. AT_DISPATCH_INTEGRAL_TYPES 所有整数类型
3. AT_DISPATCH_FLOATING_TYPES_AND2 所有浮点数类型

代码块

```
1 int32, int8, uint8, int16, int64, float32, bfloat16, half, double, bool
```

mean, mean.dim, mean.out

参考信息

官方参考文档

- <https://pytorch.org/docs/stable/generated/torch.mean.html>

cuda kernel 参考实现

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/ReduceMomentKernel.cu#L44>
- 同样使用通用的 gpu_reduce_kernel 实现
- float16 和 bfloat16 使用 float 作为中间计算类型

- dtype cast:

- cast1: 除 self 为 float16/bfloat16 且 out 为 float 情况下, 其他情况都会将 self 转换为 out dtype
- cast2:
 - self 为 float16/bfloat16 时, 在 reduce 过程中会将其中一个操作数 static_cast 为 float, out dtype 可以为 float 或者 float16/bfloat16, 后者则需要 project 时转换为 float16/bfloat16
 - 因此接口定义的时候 out 的 dtype 是可能和 self 的 dtype 不一样的。

代码块

```
1 Tensor mean(const Tensor & self, OptionalIntArrayRef dim, bool keepdim,
  ::std::optional<ScalarType> dtype); // {"schema": "aten::mean.dim(Tensor self,
  int[1]? dim, bool keepdim=False, *, ScalarType? dtype=None) -> Tensor",
  "dispatch": "True", "default": "True"}
2 Tensor & mean_out(const Tensor & self, OptionalIntArrayRef dim, bool keepdim,
  ::std::optional<ScalarType> dtype, Tensor & out); // {"schema":
  "aten::mean.out(Tensor self, int[1]? dim, bool keepdim=False, *, ScalarType?
  dtype=None, Tensor(a!) out) -> Tensor(a!)", "dispatch": "True", "default":
  "False"}
```

Function 接口

接口定义

- 暂时保证 self 和 out 的 dtype 是一样的。后续需支持 self 为 bfloat16/float16, out 为 float 的情况。

代码块

```
1 void fn(
2     SiDeviceHandle dev, // Device handle
3     void *out, // Output indices (int64_t)
4     const int64_t* out_sizes, // Output sizes [ndims]
5     const int64_t* out_strides, // Output strides [ndims]
6     const void *self, // Input data pointer
7     const int64_t* self_sizes, // Input sizes [ndims]
8     const int64_t* self_strides, // Input strides [ndims]
9     int ndims, // Number of dimensions
10    int* dim, // Reduction dimensions
11    int dim_size, // Number of reduction dimensions
12    bool keepdim = false // Keep reduced dimension
13 );
```

```

代码块 = torch.randn(1, 3)
2 a
3 tensor([[ 0.2294, -0.5481,  1.3288]])
4 torch.mean(a)
5 tensor(0.3367)
6
7 a = torch.randn(4, 4)
8 a
9 tensor([[ -0.3841,  0.6320,  0.4254, -0.7384],
10         [ -0.9644,  1.0131, -0.6549, -1.4279],
11         [ -0.2951, -1.3350, -0.7694,  0.5600],
12         [  1.0842, -0.9580,  0.3623,  0.2343]])
13 torch.mean(a, 1)
14 tensor([ -0.0163, -0.5085, -0.4599,  0.1807])
15 torch.mean(a, 1, True)
16 tensor([[ -0.0163],
17         [ -0.5085],
18         [ -0.4599],
19         [  0.1807]])

```

ISSUE

Dtypes

- AT_ALL_TYPES: uint8_t, int8_t, int16_t, int, int64_t, float, double
- Half, BFloat16: float16, bfloat16

代码块

```

1 AT_DISPATCH_ALL_TYPES_AND_COMPLEX
2 kHalf
3 kBFloat16

```

初版说明

数据类型: float, bfloat16

函数接口: 如果需要使用Tile Core, 要求dim1_size按照1024B对齐, pad value=0

代码块

```

1 /// @brief Returns the mean value of each row of the tensor in the given
dimension.
2 /// @tparam T Datatype support [float32, bfloat16], the datatype of input and
output should be the same, can be different in the future.
3 /// @param d_out Output tensor with shape
4 /// @param d_in Input tensor with shape [dim0_size, dim1_size]

```

```

5  /// @param dim0_size Dimension 0 size of input tensor
6  /// @param dim1_size Dimension 1 size of input tensor, should be 1024B aligned
   for tile core calculation, pad value 0
7  /// @param original_dim1_size Original dimension 1 size of input tensor before
   pad
8  /// @param dim The dimension or dimensions to reduce
9  template <typename T>
10 void mean_dim(void* d_out, void* d_in, int dim0_size, int dim1_size, int
   origianl_dim1_size, int dim);

```

sub.Tensor,sub.out,sub_.Tensor

实际直接调用的是 add_stub，无需 kernel 单独实现，依赖 add.Tensor

代码块

```

1  TORCH_IMPL_FUNC(sub_out) (
2    const Tensor& self, const Tensor& other, const Scalar& alpha, const Tensor&
   result
3  ) {
4    add_stub(device_type(), *this, -alpha);
5    TORCH_INTERNAL_ASSERT(result.scalar_type() == output().dtype());
6  }

```

代码块

```

1  // schema: sub.out(Tensor self, Tensor other, *, Scalar alpha=1, Tensor(a!)
   out) -> Tensor(a!)
2
3  // schema: sub.Tensor(Tensor self, Tensor other, *, Scalar alpha=1) -> Tensor
4
5  // schema: sub_.Tensor(Tensor(a!) self, Tensor other, *, Scalar alpha=1) ->
   Tensor(a!)

```

pow_.Scalar,pow.Tensor_Scalar,pow.Tensor_Scalar_out

Pytorch 文档：

- <https://pytorch.org/docs/stable/generated/torch.pow.html>
- 这三个 api 仅用于 exponent 为 scalar 的情况。

Pytorch 实现：

- meta

- self 与 exponent 均为 IntegralType 且 exponent < 0 时，会直接报错
- 获取 result_type，将 self 转换为 result_type
- Impl
 - <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/PowKernel.cu#L146>
 - 对 exp 为 0.5, -0.5, -1, 2, 3, -2 时做特殊处理，其他直接调用 gpu_kernel + pow functor

接口定义

- 可统一实现一个也兼容 pow_。Tensor, pow.Tensor_Tensor, pow.Tensor_Tensor_out 的 kernel

代码块

```

1  Tensor & pow_out(const Tensor & self, const Scalar & exponent, Tensor & out);
  // {"schema": "aten::pow.Tensor_Scalar_out(Tensor self, Scalar exponent, *,
  Tensor(a!) out) -> Tensor(a!)", "dispatch": "True", "default": "False"}
2  Tensor pow(const Tensor & self, const Scalar & exponent); // {"schema":
  "aten::pow.Tensor_Scalar(Tensor self, Scalar exponent) -> Tensor", "dispatch":
  "True", "default": "True"}
3  Tensor & pow_(Tensor & self, const Scalar & exponent); // {"schema":
  "aten::pow_.Scalar(Tensor(a!) self, Scalar exponent) -> Tensor(a!)",
  "dispatch": "True", "default": "True"}

```

Dtypes

- AT_ALL_TYPES: uint8_t, int8_t, int16_t, int, int64_t, float, double
- Half, BFloat16: float16, bfloat16
- isComplexType: c10::complex<c10::Half>, c10::complex<float>, c10::complex<double>

代码块

```

1  AT_DISPATCH_ALL_TYPES_AND2(kHalf, kBFloat16
2  isComplexType

```

pow_.Tensor, pow.Tensor_Tensor, pow.Tensor_Tensor_out, pow.Scalar, pow.Scalar_out

Pytorch 文档:

- <https://pytorch.org/docs/stable/generated/torch.pow.html>

- 对应 exponent 为 tensor 的情况

Pytorch 实现:

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/PowKernel.cu#L99>

- 支持 self 或者 exponent 为 cpu scalar, 两者不同时为 cpu scalar

- 可考虑跟 pow_。Scalar 的 kernel 实现统一, 类似 add.Tensor 传 bool is_exponent_scalar, bool is_self_scalar 来兼容两种形式

- 其中 pow.Scalar, pow.Scalar_out 实际上调用到 pow.Tensor_Tensor_out 的实现

代码块

```
1 Tensor & pow_out(const Tensor & self, const Tensor & exponent, Tensor & out);  
  // {"schema": "aten::pow.Tensor_Tensor_out(Tensor self, Tensor exponent, *,  
  Tensor(a!) out) -> Tensor(a!)", "dispatch": "True", "default": "False"}  
2 Tensor pow(const Tensor & self, const Tensor & exponent); // {"schema":  
  "aten::pow.Tensor_Tensor(Tensor self, Tensor exponent) -> Tensor", "dispatch":  
  "True", "default": "True"}  
3 Tensor & pow_(Tensor & self, const Tensor & exponent); // {"schema":  
  "aten::pow_.Tensor(Tensor(a!) self, Tensor exponent) -> Tensor(a!)",  
  "dispatch": "True", "default": "True"}  
4 Tensor & pow_out(const Scalar & self, const Tensor & exponent, Tensor & out);  
  // {"schema": "aten::pow.Scalar_out(Scalar self, Tensor exponent, *,  
  Tensor(a!) out) -> Tensor(a!)", "dispatch": "True", "default": "False"}  
5 Tensor pow(const Scalar & self, const Tensor & exponent); // {"schema":  
  "aten::pow.Scalar(Scalar self, Tensor exponent) -> Tensor", "dispatch":  
  "True", "default": "True"}
```

Dtypes

- AT_ALL_TYPES: uint8_t, int8_t, int16_t, int, int64_t, float, double

- AT_COMPLEX_TYPES: c10::complex<float>, c10::complex<double>

- kHalf, kBFloat16

代码块

```
1 AT_DISPATCH_ALL_TYPES_AND_COMPLEX_AND2(  
2     kHalf, kBFloat16
```

triu,triu_,triu.out,tril,tril_,tril.out

Pytorch 文档:

- <https://pytorch.org/docs/stable/generated/torch.triu.html>
- <https://pytorch.org/docs/stable/generated/torch.tril.html>

Pytorch 实现:

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/TriangularOps.cu#L155>
- triu / tril 共用同个 kernel, 通过 bool upper 做标识

代码块

```
1 Tensor & triu_out(const Tensor & self, int64_t diagonal, Tensor & out); //
  {"schema": "aten::triu.out(Tensor self, int diagonal=0, *, Tensor(a!) out) ->
  Tensor(a!)", "dispatch": "True", "default": "False"}
2 Tensor triu(const Tensor & self, int64_t diagonal); // {"schema":
  "aten::triu(Tensor self, int diagonal=0) -> Tensor", "dispatch": "True",
  "default": "True"}
3 Tensor & triu_(Tensor & self, int64_t diagonal); // {"schema":
  "aten::triu_(Tensor(a!) self, int diagonal=0) -> Tensor(a!)", "dispatch":
  "True", "default": "True"}
4 Tensor & tril_out(const Tensor & self, int64_t diagonal, Tensor & out); //
  {"schema": "aten::tril.out(Tensor self, int diagonal=0, *, Tensor(a!) out) ->
  Tensor(a!)", "dispatch": "True", "default": "False"}
5 Tensor tril(const Tensor & self, int64_t diagonal); // {"schema":
  "aten::tril(Tensor self, int diagonal=0) -> Tensor", "dispatch": "True",
  "default": "True"}
6 Tensor & tril_(Tensor & self, int64_t diagonal); // {"schema":
  "aten::tril_(Tensor(a!) self, int diagonal=0) -> Tensor(a!)", "dispatch":
  "True", "default": "True"}
```

注意 dim >= 3 的 case

代码块

```
1 >>> a = torch.randn((2, 3, 5))
2 >>> a
3 tensor([[[[ 1.1182, -0.1174,  0.5041, -0.6076, -1.5491],
4           [ 0.2768,  0.1626, -0.2481, -1.2557,  1.3977],
5           [ 0.0719,  0.3547,  0.5944, -0.4056,  1.6179]],
6
7           [[ 1.1097,  0.2569, -1.0501, -0.8111, -1.1674],
8           [-0.8123,  0.5775, -2.4633,  0.4311, -0.3557],
9           [ 0.4450, -0.0715, -2.1173,  1.5523, -0.7991]]]])
10 >>> torch.triu(a, -1)
11 tensor([[[[ 1.1182, -0.1174,  0.5041, -0.6076, -1.5491],
```

```
12 [ 0.2768, 0.1626, -0.2481, -1.2557, 1.3977],
13 [ 0.0000, 0.3547, 0.5944, -0.4056, 1.6179]],
14
15 [[ 1.1097, 0.2569, -1.0501, -0.8111, -1.1674],
16 [-0.8123, 0.5775, -2.4633, 0.4311, -0.3557],
17 [ 0.0000, -0.0715, -2.1173, 1.5523, -0.7991]]])
```

Dtypes

- AT_ALL_TYPES: `uint8_t, int8_t, int16_t, int, int64_t, float, double`
- AT_COMPLEX_TYPES: `c10::complex<float>, c10::complex<double>`
- Half, Bool, BFloat16, ComplexHalf: `float16, bool, bfloat16, c10::complex<c10::Half>`

代码块

```
1 AT_DISPATCH_ALL_TYPES_AND_COMPLEX_AND4(
2     at::ScalarType::ComplexHalf,
3     at::ScalarType::Half,
4     at::ScalarType::BFloat16,
5     at::ScalarType::Bool
```

masked_fill_.Scalar,masked_fill.Scalar,masked_fill_.Tensor,masked_fill.Tensor

Pytorch 文档

- https://pytorch.org/docs/stable/generated/torch.Tensor.masked_fill_.html#torch.Tensor.masked_fill_
- https://pytorch.ac.cn/docs/stable/generated/torch.Tensor.masked_fill_.html

Pytorch 实现

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/Indexing.cu#L1549>
- 如要求 contiguous 以及 aligned, torch_sipu 会通过 TensorIteratorBridge 来保证 self 和 mask 为 contiguous 且 aligned 且 shape 一致
- value 为 Tensor 时会直接拿到 value 的 scalar 值
- 直接调用 gpu_kernel 实现

代码块

```

1 Tensor & masked_fill_(Tensor & self, const Tensor & mask, const Scalar &
  value); // {"schema": "aten::masked_fill.Scalar(Tensor(a!) self, Tensor mask,
  Scalar value) -> Tensor(a!)", "dispatch": "True", "default": "False"}
2 Tensor masked_fill(const Tensor & self, const Tensor & mask, const Scalar &
  value); // {"schema": "aten::masked_fill.Scalar(Tensor self, Tensor mask,
  Scalar value) -> Tensor", "dispatch": "True", "default": "True"}
3 Tensor & masked_fill_(Tensor & self, const Tensor & mask, const Tensor &
  value); // {"schema": "aten::masked_fill_.Tensor(Tensor(a!) self, Tensor mask,
  Tensor value) -> Tensor(a!)", "dispatch": "True", "default": "False"}
4 Tensor masked_fill(const Tensor & self, const Tensor & mask, const Tensor &
  value); // {"schema": "aten::masked_fill.Tensor(Tensor self, Tensor mask,
  Tensor value) -> Tensor", "dispatch": "True", "default": "True"}
5 Tensor & masked_fill_out(const Tensor & self, const Tensor & mask, const
  Scalar & value, Tensor & out); // {"schema":
  "aten::masked_fill.Scalar_out(Tensor self, Tensor mask, Scalar value, *,
  Tensor(a!) out) -> Tensor(a!)", "dispatch": "True", "default": "True"}
6 Tensor & masked_fill_out(const Tensor & self, const Tensor & mask, const
  Tensor & value, Tensor & out); // {"schema":
  "aten::masked_fill.Tensor_out(Tensor self, Tensor mask, Tensor value, *,
  Tensor(a!) out) -> Tensor(a!)", "dispatch": "True", "default": "True"}

```

Dtypes

- AT_ALL_TYPES: `uint8_t, int8_t, int16_t, int, int64_t, float, double`
- AT_COMPLEX_TYPES: `c10::complex<float>, c10::complex<double>`
- Half, Bool, BFloat16, ComplexHalf: `float16, bool, bfloat16, c10::complex<c10::Half>`

代码块

```

1 AT_DISPATCH_ALL_TYPES_AND_COMPLEX_AND4(
2     kBool, kHalf, kBFloat16, kComplexHalf

```

_local_scalar_dense

无需 kernel 实现，依赖 runtime api `cudaMemcpyAsync`, `cudaStreamSynchronize`

Pytorch 实现

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/CUDAScalar.cu#L45>
- <https://github.com/pytorch/pytorch/blob/v2.5.0/c10/cuda/CUDAFUNCTIONS.h#L77>

代码块

```

1 # 间接调用
2 >>> x = torch.tensor([42])
3 >>> x.item()
4 [call] op=[aten::item], key=[PythonTLSnapshot]
5 [redispatchBoxed] op=[aten::item], key=[AutogradCPU]
6 [call] op=[aten::_local_scalar_dense], key=[PythonTLSnapshot]
7 [redispatchBoxed] op=[aten::_local_scalar_dense], key=[AutogradCPU]
8 [redispatch] op=[aten::_local_scalar_dense], key=[Python]
9 [callBoxed] op=[aten::_local_scalar_dense], key=[Python]
10 [callBoxed] op=[aten::_local_scalar_dense], key=[CPU]
11 42

```

cos,cos_,cos.out

Pytorch 官方文档参考: <https://pytorch.org/docs/stable/generated/torch.cos.html>

参考实现:

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/UnaryGeometricCosKernel.cu>

代码块

```

1 // schema: cos(Tensor self) -> Tensor
2 // schema: cos_(Tensor(a!) self) -> Tensor(a!)
3 // schema: cos.out(Tensor self, *, Tensor(a!) out) -> Tensor(a!)
4
5 /**
6  * @brief Computes the cosine of each element in the input tensor `input` and
7  * stores the result in the output tensor `out`.
8  *
9  * @param dev Device on which the operation is performed.
10 * @param out Pointer to the output tensor where the cosine values
11 * will be stored.
12 * @param input Pointer to the input tensor containing values to
13 * compute cosine.
14 * @param dims Number of dimensions in the input/output tensor.
15 * @param sizes Sizes of each dimension of the input/output tensor.
16 * @param strides_out Strides of each dimension of the output tensor.
17 * @param strides_input Strides of each dimension of the input tensor.
18 * @param element_sizes Optional pointer to the element sizes of each tensor
19 * (if nullptr, strides are in bytes).
20 */
21 void fn(dev, const void* out, const void* input, int dims, const int64_t*
22 sizes, const int64_t* strides_out, const int64_t* strides_input, const
23 int64_t* element_sizes=nullptr)
24 //updated

```

```
19 void fn(dev, const void* out, const void* self, int ndims, const int64_t*
    sizes, const int64_t* out_strides, const int64_t* self_strides, const int64_t*
    element_sizes=nullptr)
```

Dtype 支持

- Floating types + complex types: `float`, `double`, `at::Half`, `at::BFloat16`, `c10::complex<float>`, `c10::complex<double>`, `c10::complex<c10::Half>`

sin, sin_, sin.out

Pytorch 官方文档参考: <https://pytorch.org/docs/stable/generated/torch.sin.html>

参考实现:

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/UnaryGeometricSinkKernel.cu>

Function 接口, Dtype 支持同 `cos`

arange.start_out, arange, arange.out, arange.start, arange.start_step

参考信息

文档:

<https://pytorch.org/docs/stable/generated/torch.arange.html>

参考实现

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/RangeFactories.cu#L213>

代码块

```
1  参考用法
2  # 创建一个预分配的输出张量
3  out = torch.empty(5, dtype=torch.float32)
4
5  # 使用 arange.start_out 填充输出张量
6  torch.arange(0, 5, out=out)
7  print(out) # tensor([0., 1., 2., 3., 4.])
8
9  # 重用同一个输出张量
10 torch.arange(5, 10, out=out)
11 print(out) # tensor([5., 6., 7., 8., 9.])
12
13 #指定步长
14 out = torch.empty(5, dtype=torch.float32)
15 torch.arange(0, 10, step=2, out=out)
```

```
16 print(out) # tensor([0., 2., 4., 6., 8.])
```

c++

```
1 // build/aten/src/ATen/RegisterCompositeExplicitAutograd.cpp
2
3 m.impl("arange",
4 TORCH_FN(wrapper_CompositeExplicitAutograd__arange));
5 m.impl("arange.out",
6 TORCH_FN(wrapper_CompositeExplicitAutograd_out_arange_out));
7 m.impl("arange.start",
8 TORCH_FN(wrapper_CompositeExplicitAutograd_start_arange));
9 m.impl("arange.start_step",
10 TORCH_FN(wrapper_CompositeExplicitAutograd_start_step_arange));
11
12 namespace {
13 at::Tensor wrapper_CompositeExplicitAutograd__arange(const at::Scalar & end,
14 ::std::optional<at::ScalarType> dtype, ::std::optional<at::Layout> layout,
15 ::std::optional<at::Device> device, ::std::optional<bool> pin_memory) {
16     // No device check
17     // DeviceGuard omitted
18     return at::native::arange(end, dtype, layout, device, pin_memory);
19 }
20 } // anonymous namespace
21 namespace {
22 at::Tensor & wrapper_CompositeExplicitAutograd_out_arange_out(const at::Scalar
23 & end, at::Tensor & out) {
24     // No device check
25     // DeviceGuard omitted
26     return at::native::arange_out(end, out);
27 }
28 } // anonymous namespace
29 namespace {
30 at::Tensor wrapper_CompositeExplicitAutograd_start_arange(const at::Scalar &
31 start, const at::Scalar & end, ::std::optional<at::ScalarType> dtype,
32 ::std::optional<at::Layout> layout, ::std::optional<at::Device> device,
33 ::std::optional<bool> pin_memory) {
34     // No device check
35     // DeviceGuard omitted
36     return at::native::arange(start, end, dtype, layout, device, pin_memory);
37 }
38 } // anonymous namespace
39 namespace {
40 at::Tensor wrapper_CompositeExplicitAutograd_start_step_arange(const
41 at::Scalar & start, const at::Scalar & end, const at::Scalar & step,
```

```

::std::optional<at::ScalarType> dtype, ::std::optional<at::Layout> layout,
::std::optional<at::Device> device, ::std::optional<bool> pin_memory) {
35 // No device check
36 // DeviceGuard omitted
37 return at::native::arange(start, end, step, dtype, layout, device,
pin_memory);
38 }
39 } // anonymous namespace
40
41 //aten/src/ATen/native/TensorFactories.cpp
42 // ~~~~~
43
44 Tensor arange(const Scalar& end,
45               std::optional<ScalarType> dtype,
46               std::optional<Layout> layout,
47               std::optional<Device> device,
48               std::optional<bool> pin_memory) {
49     return native::arange(/*start=*/0, end, dtype, layout, device, pin_memory);
50 }
51
52 Tensor arange(const Scalar& start, const Scalar& end,
53               std::optional<ScalarType> dtype,
54               std::optional<Layout> layout,
55               std::optional<Device> device,
56               std::optional<bool> pin_memory) {
57     return native::arange(
58         start, end, /*step=*/1, dtype, layout, device, pin_memory);
59 }
60
61 Tensor arange(
62     const Scalar& start,
63     const Scalar& end,
64     const Scalar& step,
65     std::optional<ScalarType> dtype,
66     std::optional<Layout> layout,
67     std::optional<Device> device,
68     std::optional<bool> pin_memory) {
69     // See [Note: hacky wrapper removal for TensorOptions]
70     TensorOptions options =
TensorOptions().dtype(dtype).layout(layout).device(device).pinned_memory(pin_me
memory);
71
72     bool set_to_integral_dtype = !options.has_dtype() &&
73         // bool inputs are considered integral
74         start.isIntegral(true) &&
75         end.isIntegral(true) &&

```

```

76     step.isIntegral(true);
77
78     Tensor result = set_to_integral_dtype
79         ? at::empty({0}, options.dtype(at::ScalarType::Long))
80         : at::empty({0}, options);
81     return at::arange_out(result, start, end, step);
82 }
83
84 Tensor& arange_out(const Scalar& end, Tensor& result) {
85     return at::arange_out(result, /*start=*/0, end, /*step=*/1);
86 }
87
88
89 最终都dispatch到at::arange_out_cuda的实现在上面链接中

```

Function 接口

代码块

```

1 // schema: arange.start_out(Scalar start, Scalar end, Scalar step=1, *,
  // Tensor(a!) out) -> Tensor(a!)
2 template<typename T, typename T_acc_type>
3 /**
4  * @brief Creates a 1-D tensor of size  $\text{ceil}((\text{end} - \text{start}) / \text{step})$  with values
  * from start to end with step step
5  *
6  * @param dev          Device handle for CUDA/CPU device management
7  * @param start        Starting value of the sequence
8  * @param end          Upper limit of the sequence (exclusive)
9  * @param step         Spacing between values (must be nonzero)
10 * @param dtype        Data type of output tensor
11 * @param out          Optional pre-allocated output tensor (can be non-
  contiguous)
12 *
13 * @note Memory Layout:
14 *     - Fully supports non-contiguous output tensors
15 *     - For non-contiguous tensors, computation is done in a temporary
  contiguous tensor
16 *     and then copied to the output tensor
17 *     - This approach ensures optimal performance while supporting any
  memory layout
18 *
19 * @return Tensor      The output tensor
20 */
21 template<typename T>

```

```

22 void fn_arange(
23     SiDeviceHandle *dev,
24     void* out,
25     const T& start,
26     const T& end,
27     const T& step
28 );
29 支持数据类型:
30 整数类型
31 at::ScalarType::Byte (uint8)
32 at::ScalarType::Char (int8)
33 at::ScalarType::Short (int16)
34 at::ScalarType::Int (int32)
35 at::ScalarType::Long (int64)
36 浮点类型
37 at::ScalarType::Float (float32)
38 at::ScalarType::Double (float64)
39 低精度类型
40 at::ScalarType::Half (float16)
41 at::ScalarType::BFloat16 (bfloat16)

```

非 contiguous 的是否能支持

cuda 实现: 检查输出张量是否连续, 如果不是, 它会创建一个临时的连续张量, 填充数据, 然后将结果复制回原始非连续张量。

可以在传入 kernel 前进行判断, 如果是 contiguous, 可直接传入 kernel, 否则在 pytorch 分配一块连续张量传入 kernel, 然后将填充后的 tensor copy 给原始非 contiguous tensor。

argmax, argmax.out

参考信息

Pytorch 官方文档参考:

<https://pytorch.org/docs/stable/generated/torch.argmax.html#torch.argmax>

参考实现:

<https://github.com/pytorch/pytorch/blob/main/aten/src/ATen/native/cuda/ReduceArgMaxKernel.cu>

<https://github.com/pytorch/pytorch/blob/7bd2e3bca1809a586a2b3072a16c75e871a256b3/aten/src/ATen/native/cuda/Reduce.cuh#L1173>

代码块

```
1 import torch
2
3 # 1. 基本用法 - 1D张量
4 x = torch.tensor([1, 2, 3, 4, 2])
5 indices = torch.argmax(x)
6 print(indices) # tensor(3) - 最大值4的索引位置
7
8 # 2. 指定维度 - 2D张量
9 x = torch.tensor([[1, 2, 3],
10                  [3, 1, 2]])
11 # 按行找最大值(dim=1)
12 indices = torch.argmax(x, dim=1)
13 print(indices) # tensor([2, 0]) - 每行最大值的索引
14
15 # 3. 保持维度(keepdim=True)
16 indices = torch.argmax(x, dim=1, keepdim=True)
17 print(indices) # tensor([[2],
18                  #       [0]])
19
20 # 4. 处理相同的最大值(返回第一个)
21 x = torch.tensor([1, 4, 4, 3])
22 indices = torch.argmax(x)
23 print(indices) # tensor(1) - 返回第一个4的位置
24
25 # 5. 在GPU上使用
26 if torch.cuda.is_available():
27     x_cuda = x.cuda()
28     indices = torch.argmax(x_cuda)
29     print(indices.device) # cuda:0
30
31 # 6. 处理不同数据类型
32 x = torch.tensor([1.1, 2.2, 3.3], dtype=torch.float16)
33 indices = torch.argmax(x)
34 print(indices) # tensor(2)
35
36 # 7. 多维张量
37 x = torch.randn(2, 3, 4)
38 # 在最后一维上找最大值
39 indices = torch.argmax(x, dim=-1)
40 print(indices.shape) # torch.Size([2, 3])
41
42 # 8. 全局最大值
43 x = torch.tensor([[1, 2], [3, 4]])
44 index = torch.argmax(x) # 展平后找最大值
45 print(index) # tensor(3) - 4的位置
46
47 # 9. 带掩码的argmax
```

```

48 mask = torch.tensor([True, False, True, True])
49 x = torch.tensor([1, 2, 3, 4])
50 indices = torch.argmax(x * mask.float())
51 print(indices) # tensor(3)

```

Function 接口

代码块

```

1 // schema: argmax(Tensor self, int? dim=None, bool keepdim=False) -> Tensor
2 // schema: argmax.out(Tensor self, int? dim=None, bool keepdim=False, *,
  Tensor(a!) out) -> Tensor(a!)
3 //at::Tensor & wrapper_CUDA_argmax_out_out(const at::Tensor & self,
  ::std::optional<int64_t> dim, bool keepdim, at::Tensor & out)
4
5 /**
6  * @brief Find indices of maximum values along a dimension
7  *
8  * @param dev Device handle for CUDA device management
9  * @param out Output indices tensor (int64_t)
10 * @param out_sizes Output tensor sizes [ndims]
11 * @param out_strides Output tensor strides [ndims]
12 * @param self Input tensor data pointer
13 * @param self_sizes Input tensor sizes [ndims]
14 * @param self_strides Input tensor strides [ndims]
15 * @param ndim Number of dimensions
16 * @param dim Dimension to reduce along (-1 means flatten)
17 * @param dtype Data type of input tensor
18 * @param keepdim Whether to keep reduced dimension
19 *
20 * @note Memory Layout:
21 * - Output tensor contains indices of maximum values
22 * - For float16/bfloat16 inputs, computation is done in float32
23 * - Indices are 64-bit integers
24 *
25 * @note Performance Optimizations:
26 * - Uses vectorized loads when possible
27 * - Special handling for contiguous tensors
28 * - Different implementations for inner/outer dimensions
29 *
30 * @note Constraints:
31 * - dim must be valid (-ndim <= dim < ndim)
32 * - Memory must be properly aligned
33 * - Maximum dimensions: 32
34 *
35 * @return int 0 on success, error code otherwise
36 */

```

```

37 void fn(
38     SiDeviceHandle dev,           // Device handle
39     void *out,                   // Output indices (int64_t)
40     const int64_t* out_sizes,    // Output sizes [ndims]
41     const int64_t* out_strides,  // Output strides [ndims]
42     const void *self,           // Input data pointer
43     const int64_t* self_sizes,   // Input sizes [ndims]
44     const int64_t* self_strides, // Input strides [ndims]
45     int ndims,                   // Number of dimensions
46     int dim,                     // Reduction dimension
47     int dtype,                   // Input data type
48     bool keepdim = false        // Keep reduced dimension
49 );
50
51 //支持 Half 和 BFloat16 类型, 但是使用 float 作为中间计算类型
52 //其余基本数据类型均使用本身的数据类型作为中间计算类型
53 - uint8_t    // Byte
54 - int8_t     // Char
55 - int16_t    // Short
56 - int32_t    // Int
57 - int64_t    // Long
58 - float      // Float
59 - double     // Double

```

ISSUE

Dtype 支持 //cuda 的实现 支持 *Byte Char Short Int Long Float Double*

float32

float16 使用 float 作为中间计算类型

bfloat16 使用 float 作为中间计算类型

mxformat

非 contiguous 的是否能支持 (如果支持需要修改 Function 接口)

需要额外存储 stride 和 size 信息, 用来进行复杂的索引计算

Memory 是否需要 Align, Align 到多少

isin.Tensor_Tensor, isin.Tensor_Tensor_out

参考信息

<https://pytorch.org/docs/stable/generated/torch.isin>

```
In [188]: (c,d,torch.isin(c, d))
[call] op=[aten::isin.Tensor_Tensor], key=[AutogradCUDA]
[redispatch] op=[aten::isin.Tensor_Tensor], key=[CUDA]
[call] op=[aten::fill_.Scalar], key=[CUDA]
[call] op=[aten::unsqueeze], key=[CUDA]
[call] op=[aten::as_strided], key=[CUDA]
[call] op=[aten::view], key=[CUDA]
[call] op=[aten::eq.Tensor], key=[CUDA]
[call] op=[aten::any.dim], key=[CUDA]
[call] op=[aten::as_strided], key=[CUDA]
[call] op=[aten::copy_], key=[CUDA]
```

```
(tensor([[ 1, 65536],
          [ 2,   3]], device='cuda:0'),
 tensor([ 22, 65536], device='cuda:0'),
 tensor([[False, True],
          [False, False]], device='cuda:0'))
```

参考实现:

<https://github.com/pytorch/pytorch/blob/3f069e7679588d5ee4b1d5b2492ca0e20f9320b5/aten/src/ATen/native/cuda/TensorCompare.cpp#L11>

代码块

```
1 import torch
2
3 # 创建测试数据
4 elements = torch.tensor([1, 2, 3, 4, 5])
5 test_elements = torch.tensor([2, 4, 6])
6
7 # 检查 elements 中的元素是否在 test_elements 中
8 result = torch.isin(elements, test_elements)
9 print(result) # 输出: tensor([False, True, False, True, False])
10
11 # 使用 Tensor 方法
12 result = elements.isin(test_elements)
13 print(result) # 输出: tensor([False, True, False, True, False])
14
15 # 使用 assume_unique 参数
16 result = torch.isin(elements, test_elements, assume_unique=True)
17
18 # 使用 invert 参数获取不在 test_elements 中的元素掩码
19 result = torch.isin(elements, test_elements, invert=True)
20 print(result) # 输出: tensor([True, False, True, False, True])
```

Function 接口

代码块

```
1 - func: isin.Tensor_Tensor_out(Tensor elements, Tensor test_elements, *, bool
  assume_unique=False, bool invert=False, Tensor(a!) out) -> Tensor(a!)
2 variants: function
```

```

3   structured: True
4   dispatch:
5     CPU, CUDA: isin_Tensor_Tensor_out
6     MPS: isin_Tensor_Tensor_out_mps
7
8   //isin.Tensor_Tensor_out(Tensor elements, Tensor test_elements, *, bool
assume_unique=False, bool invert=False, Tensor(a!) out) -> Tensor(a!)
9
10
11  /**
12  * @brief Test each element in 'elements' for membership in 'test_elements'
13  *
14  * @param dev          Device handle for device management
15  * @param elements    Input tensor to check
16  * @param test_elements Elements to test for membership
17  * @param assume_unique If True, assumes both elements and test_elements
contain unique elements
18  * @param invert      If True, return nonzero where element is NOT in
test_elements
19  * @param out         Output boolean tensor (pre-allocated)
20  * @param dtype       Data type of input tensors(optional)
21  *
22  * @note Functionality:
23  *
24  * - Returns boolean tensor of same shape as elements
25  * - True where element is in test_elements, False otherwise
26  * - If invert=True, returns opposite
27  *
28  * @note Performance:
29  *
30  * - For GPU: Materializes cross product for parallel comparison
31  * - assume_unique=True can enable optimizations
32  * - Uses tensor broadcasting internally
33  *
34  * @note Example:
35  * elements = [1, 2, 3, 4]
36  * test_elements = [2, 4, 6]
37  * out = [False, True, False, True]
38  *
39  * @return int 0 on success, error code otherwise
40  */
41  void fn(
42      SiDeviceHandle dev,          // Device handle
43      const void *elements,       // Input elements to check
44      const void *test_elements,  // Test elements to compare against
45      bool assume_unique,        // Whether inputs contain unique elements
46      bool invert,               // Whether to invert the result
47      void *out,                 // Output boolean tensor

```

```
46     int Type(optional) // Data type of input tensors (elements &
    test_elements)
47 );
```

ISSUE

- Dtype 支持 //cuda 的实现 支持 kComplexHalf, kHalf, kBool, kBFloat16
 - float32
 - float16
 - bfloat16
 - mxformat
- 非 contiguous 的是否能支持 (如果支持需要修改 Function 接口)
 - 非连续在具体实现中也是会被转换为 contiguous, 不会涉及到 stride, dim 信息
- Memory 是否需要 Align, Align 到多少

cat.out, cat

参考信息

Pytorch 官方文档参考: <https://pytorch.org/docs/stable/generated/torch.cat.html#torch.cat>

并行优化实现参考:

<https://github.com/pytorch/pytorch/blob/32f585d9346e316e554c8d9bf7548af9f62141fc/aten/src/ATen/native/cuda/Shape.cu#L288C6-L288C18>

串行实现版本参考:

<https://github.com/pytorch/pytorch/blob/32f585d9346e316e554c8d9bf7548af9f62141fc/aten/src/ATen/native/cuda/Shape.cu#L526>

代码块

```
1  int64_t offset = 0; // 跟踪当前复制位置
2  for (const Tensor& t : materialized) {
3      if (cat_should_skip_tensor(t)) continue;
4      // 获取当前张量在连接维度上的大小
5      int64_t dimSize = t.size(dim);
6      // 从结果张量中切出对应大小的视图
7      Tensor nt = at::narrow(result, dim, offset, dimSize);
8      // 将当前张量复制到视图中
9      copy_(nt, t);
10     //更新偏移量
11     offset += dimSize;
12 }
13
```

```

14 //example:
15 // 例如连接两个张量 [2,3] 和 [2,4] 在维度1上:
16 tensor1 = [[1,2,3],      result = [[1,2,3,4,5,6,7],
17         [4,5,6]]         [8,9,10,11,12,13,14]]
18 tensor2 = [[4,5,6,7],
19         [11,12,13,14]]
20
21 // 步骤:
22 1. 创建结果张量 result, 大小为 [2,7]
23 2. narrow 操作切出前3列: nt1 = result[:, 0:3]
24 3. 复制第一个张量: copy_(nt1, tensor1)
25 4. narrow 操作切出后4列: nt2 = result[:, 3:7]
26 5. 复制第二个张量: copy_(nt2, tensor2)

```

Function 接口

代码块

```

1 //schema :
2 //cat.out(Tensor[] tensors, int dim=0, *, Tensor(a!) out) -> Tensor(a!)
3 //cat(Tensor[] tensors, int dim=0) -> Tensor
4 # alias for torch.cat
5 //concat(Tensor[] tensors, int dim=0) -> Tensor
6 //concat.out(Tensor[] tensors, int dim=0, *, Tensor(a!) out) -> Tensor(a!)
7 //concat.names(Tensor[] tensors, Dimname dim) -> Tensor
8 //concat.names_out(Tensor[] tensors, Dimname dim, *, Tensor(a!) out) ->
  Tensor(a!)
9 # alias for torch.cat
10 //concatenate(Tensor[] tensors, int dim=0) -> Tensor
11 //concatenate.out(Tensor[] tensors, int dim=0, *, Tensor(a!) out) -> Tensor(a!)
12 //concatenate.names(Tensor[] tensors, Dimname dim) -> Tensor
13 //concatenate.names_out(Tensor[] tensors, Dimname dim, *, Tensor(a!) out) ->
  Tensor(a!)
14
15 //概括起来拥有两种调用方式
16 //使用整数作为维度参数
17 "concat(TensorList tensors, int64_t dim=0, *, Tensor out=None)"
18 //使用命名维度(Dimname)作为维度参数
19 "concat(TensorList tensors, Dimname dim, *, Tensor out=None)"
20
21 // 原始的实现, 对于每个函数, 有四种情况:
22 case 0: {
23     if (_r.isNone(2)) {
24         // 1. 使用整数维度, 无输出张量
25         return at::concat(tensors, dim);
26     } else {
27         // 2. 使用整数维度, 有输出张量

```

```

28     return at::concat_out(out, tensors, dim);
29 }
30 }
31 case 1: {
32     if (_r.isNone(2)) {
33         // 3. 使用命名维度, 无输出张量
34         return at::concat(tensors, dim);
35     } else {
36         // 4. 使用命名维度, 有输出张量
37         return at::concat_out(out, tensors, dim);
38     }
39 }

```

代码块

```

1 //接口中不包含Dimname, 需要提前处理转换为dim, 默认输出存放在out指针, 不考虑是否return
  output
2 /**
3  * @brief Concatenate multiple tensors along the specified dimension
4  *
5  * @param dev Device handle
6  * @param out Output tensor data pointer (pre-allocated)
7  * @param out_ndims Number of outs' dimensions
8  * @param out_sizes Array of output tensor sizes for each dimension
9  * @param out_strides Array of output tensor strides for each dimension
10 * @param tensors Array of input data pointers
11 * @param num_tensors Number of input tensors
12 * @param tensors_sizes Array of input sizes [num_tensors * tensors_ndims]
13 * @param tensors_strides Array of input strides [num_tensors * tensors_ndims]
14 * @param dim Concatenation dimension
15 * @param element_size element_size of tensor
16 * @param tensors_ndims Number of tensors' dimensions (<= 4)
17 * @param all_contiguous Whether all inputs are memory contiguous
18 * @param all_same_dtype Whether all inputs have same data type
19 * @param all_same_stride Whether all inputs have same stride pattern
20 * @param memory_format Memory format (0: Contiguous, 1: ChannelsLast, 2:
  ChannelsLast3d) (optional)
21 *
22 * @note Memory Layout:
23 *     - tensors_size[i][j]->[(i-1)*tensors_ndim+j]: size of i-th tensor in
  j-th dimension
24 *     - tensors_stride[i][j]->[(i-1)*tensors_ndim+j]: stride of i-th tensor
  in j-th dimension
25 *
26 * @note Performance Optimizations:
27 *     - Fast path when all_contiguous is true

```

```

28 * - Special handling for aligned memory
29 * - Batch processing when all_same_dtype is true
30 *
31 * @note Constraints:
32 * - Maximum dimensions: 4
33 * - All tensors must be 32-bit indexable
34 *
35 * @return int 0 on success, error code otherwise
36 */
37 void fn(
38     SiDeviceHandle dev,           // Device handle
39     void *out,                   // Output data pointer
40     const int64_t out_ndims,     // Outputs' dimensions
41     const int64_t* out_sizes,    // Output sizes [out_ndim]
42     const int64_t* out_strides,  // Output strides [out_ndim]
43     void **tensors,             // Array of input data pointers
44     [num_tensors]
45     int num_tensors,            // Number of input tensors
46     const int64_t tensors_ndims, // tensor dimensions
47     const int64_t* tensors_sizes, // Input sizes [num_tensors* tensors_ndim]
48     const int64_t* tensors_strides, // Input strides [num_tensors*
49     tensors_ndim]
50     int dim,                   // Concatenation dimension
51     int element_size,         // element_size of tensor
52     bool all_contiguous,      // Whether all inputs are contiguous
53     bool all_same_dtype,     // Whether all inputs have same dtype
54     bool all_same_stride,    // Whether all inputs have same stride
55     int memory_format         // Memory format (optional)
56 );
57
58 // input_sizes, input_strides说明, 假设连接两个 2x3 张量
59 int64_t out_size[] = {2, 6}; // 2x6 输出
60 int64_t out_stride[] = {6, 1}; // 连续存储
61
62 int64_t tensors_size[] = {
63     2, 3, // 第一个输入
64     2, 3 // 第二个输入
65 };
66 int64_t tensors_stride[] = {
67     3, 1, // 第一个输入
68     3, 1 // 第二个输入
69 };
70
71 //cuda 的实现 支持 kComplexHalf, kHalf, kBool, kBFloat16

```

Dtype 支持 //cuda 的实现 支持 kComplexHalf, kHalf, kBool, kBFloat16

float32

float16

bfloat16

mxformat

非 contiguous 的是否能支持 (如果支持需要修改 Function 接口)

非连续情况使用较小的批处理大小 (CAT_ARRAY_BATCH_SIZE/2)

需要额外存储 stride 和 size 信息, 用来进行复杂的索引计算

内存格式必须是 Contiguous

Memory 是否需要 Align, Align 到多少

参考 cuda 实现中有 16 字节对齐的检查, 如果对齐可以进行加速

cumsum.out, cumsum, cumsum_

参考信息

Pytorch 官方文档参考: <https://pytorch.org/docs/stable/generated/torch.cumsum.html>

cuda 实现参考:

<https://github.com/pytorch/pytorch/blob/main/aten/src/ATen/native/cuda/CumsumKernel.cu>

<https://github.com/pytorch/pytorch/blob/bb7e8fbd668c7c8931436b4a935b26911cbe0daf/aten/src/ATen/native/cuda/ScanUtils.cuh#L446>

代码块

```
1 import torch
2
3 # 1. 基本用法 - 1D张量
4 x = torch.tensor([1, 2, 3, 4])
5 y = torch.cumsum(x, dim=0)
6 print(y) # tensor([1, 3, 6, 10])
7
8 # 2. 2D张量 - 按行累加
9 x = torch.tensor([[1, 2, 3],
10                  [4, 5, 6],
11                  [7, 8, 9]])
12 y = torch.cumsum(x, dim=0)
13 print(y)
14 # tensor([[ 1,  2,  3],
15          # [ 5,  7,  9],
16          # [12, 15, 18]])
17
```

```

18 # 3. 2D张量 - 按列累加
19 y = torch.cumsum(x, dim=1)
20 print(y)
21 # tensor([[ 1,  3,  6],
22 #         [ 4,  9, 15],
23 #         [ 7, 15, 24]])
24
25 # 4. 指定输出类型
26 y = torch.cumsum(x, dim=0, dtype=torch.float)
27 print(y.dtype) # torch.float32
28
29 # 5. 使用方法形式
30 y = x.cumsum(dim=0)
31
32 # 6. 在GPU上使用
33 if torch.cuda.is_available():
34     x_cuda = x.cuda()
35     y_cuda = torch.cumsum(x_cuda, dim=0)
36     print(y_cuda.device) # cuda:0
37
38 # 7. 处理大维度张量
39 x = torch.randn(2, 3, 4, 5)
40 y = torch.cumsum(x, dim=2) # 在第3维上累加
41 print(y.shape) # torch.Size([2, 3, 4, 5])
42
43 # 8. 负维度索引
44 y = torch.cumsum(x, dim=-1) # 等同于 dim=3

```

Function 接口

代码块

```

1 cumsum.out(Tensor self, int dim, *, ScalarType? dtype=None, Tensor(a!) out) ->
  Tensor(a!)
2
3 void launch_cumsum_cuda_kernel(const TensorBase& result, const TensorBase&
  self, int64_t dim) {
4     AT_DISPATCH_ALL_TYPES_AND_COMPLEX_AND2(
5         ScalarType::Half, ScalarType::BFloat16,
6         self.scalar_type(), "cumsum_cuda",
7         [&]() {
8             scalar_t init = 0;
9             scan_dim<scalar_t>(
10                self,
11                result,
12                dim,
13                init,

```

```

14         std::plus<scalar_t>());
15     });
16 }
17
18 /**
19  * @brief Compute the cumulative sum along a dimension
20  *
21  * @param dev          Device handle for CUDA device management
22  * @param out          Output tensor data pointer (pre-allocated)
23  * @param out_sizes    Output tensor sizes for each dimension [ndim]
24  * @param out_strides  Output tensor strides for each dimension [ndim]
25  * @param self         Input tensor data pointer
26  * @param self_sizes   Input tensor sizes for each dimension [ndim]
27  * @param self_strides Input tensor strides for each dimension [ndim]
28  * @param ndims        Number of dimensions
29  * @param dim          Dimension to compute cumsum along
30  * @param dtype        Data type of tensors
31  * @param all_contiguous Whether all tensors are memory contiguous
32  *
33  * @note Memory Layout:
34  *     - Input and output tensors must have same shape
35  *     - For contiguous tensors, memory is assumed to be packed
36  *     - For non-contiguous tensors, strides specify memory layout
37  *
38  * @note Performance Optimizations:
39  *     - Special fast path when input is contiguous
40  *     - Uses vectorized operations for aligned memory
41  *     - Different implementations for innermost vs outer dimensions
42  *
43  * @note Constraints:
44  *     - dim must be valid (0 <= dim < ndim)
45  *     - Memory must be properly aligned (16 bytes)
46  *     - Maximum dimensions: 32
47  *
48  */
49 void cumsum_fn(
50     SiDeviceHandle dev,          // Device handle
51     const void *out,            // Output data pointer
52     const int64_t* out_sizes,    // Output sizes [ndims]
53     const int64_t* out_strides,  // Output strides [ndims]
54     const void *self,           // Input data pointer
55     const int64_t* self_sizes,   // Input sizes [ndims]
56     const int64_t* self_strides, // Input strides [ndims]
57     int ndims,                  // Number of dimensions
58     int dim,                    // Cumsum dimension
59     int dtype,                  // Data type
60     bool all_contiguous = true  // Whether all tensors are contiguous

```

```
61 );
62
63 support dtype:
64 - uint8_t // Byte
65 - int8_t // Char
66 - int16_t // Short
67 - int32_t // Int
68 - int64_t // Long
69 - float // Float
70 - double // Double
71 - complex<float> // ComplexFloat
72 - complex<double> // ComplexDouble
73 - ScalarType::Half // float16
74 - ScalarType::BFloat16 // bfloat16
```

ISSUE

Dtype 支持

float32

float16

bfloat16

mxformat

非 contiguous 的是否能支持（如果支持需要修改 Function 接口）

For contiguous tensors, memory is assumed to be packed

For non-contiguous tensors, strides specify memory layout

Memory 是否需要 Align, Align 到多少

clamp

Clamp 和 clamp_tensor 在实现逻辑上是一致的，这里设计一个统一接口。clamp_min, clamp_max 是 clamp 实现的特例，也放在一起进行处理

clamp, clamp_, clamp.out, clamp_max, clamp_max_, clamp_max.out, clamp_min, clamp_min_, clamp_min.out

clamp.Tensor, clamp_.Tensor, clamp.Tensor_out, clamp_max.Tensor, clamp_max_.Tensor, clamp_max.Tensor_out, clamp_min.Tensor, clamp_min_.Tensor, clamp_min.Tensor_out

参考信息

pytorch 官方资料: <https://pytorch.org/docs/stable/generated/torch.clamp.html#torch.clamp>

cuda 实现:

<https://github.com/pytorch/pytorch/blob/7ae0ce6360b6e4f944906502d20da24c04debee5/aten/src/ATen/native/cuda/TensorCompare.cu#L61C13-L61C32> (scalar)

<https://github.com/pytorch/pytorch/blob/7ae0ce6360b6e4f944906502d20da24c04debee5/aten/src/ATen/native/cuda/TensorCompare.cu#L44>(tensor)

代码块

```
1 // 1. 基本用法 - clamp(input, min, max)
2 {
3     auto input = torch::tensor({-2, 0, 2, 4, 6});
4
5     // 使用标量限制范围
6     auto output = torch::clamp(input, /*min=*/0, /*max=*/4);
7     // output: [0, 0, 2, 4, 4]
8
9     // 使用 Tensor 限制范围
10    auto min = torch::tensor({-1, 1, 0, 2, -2});
11    auto max = torch::tensor({1, 3, 2, 4, 6});
12    output = torch::clamp(input, min, max);
13    // 每个元素独立限制在对应的范围内
14 }
15
16 // 2. 原地操作 - clamp_(input, min, max)
17 {
18     auto input = torch::tensor({-2, 0, 2, 4, 6});
19
20     // 原地修改
21     input.clamp_(/*min=*/0, /*max=*/4);
22     // input: [0, 0, 2, 4, 4]
23 }
24
25 // 3. 使用 out 参数 - clamp.out(input, min, max, out)
26 {
27     auto input = torch::tensor({-2, 0, 2, 4, 6});
28     auto out = torch::empty_like(input);
29
30     torch::clamp_out(out, input, /*min=*/0, /*max=*/4);
31     // out: [0, 0, 2, 4, 4]
32 }
33
34 // 4. 只限制最小值 - clamp_min
35 {
36     auto input = torch::tensor({-2, 0, 2, 4, 6});
37
38     auto output = torch::clamp_min(input, /*min=*/2);
```

```

39 // output: [2, 2, 2, 4, 6]
40
41 // 原地版本
42 input.clamp_min_(2);
43 }
44
45 // 5. 只限制最大值 - clamp_max
46 {
47     auto input = torch::tensor({-2, 0, 2, 4, 6});
48
49     auto output = torch::clamp_max(input, /*max=*/2);
50     // output: [-2, 0, 2, 2, 2]
51
52     // 原地版本
53     input.clamp_max_(2);
54 }
55
56 // 6. 处理 NaN
57 {
58     auto input = torch::tensor({-2, NaN, 2, 4, 6});
59
60     // NaN 会被保留
61     auto output = torch::clamp(input, /*min=*/0, /*max=*/4);
62     // output: [0, NaN, 2, 4, 4]
63 }
64
65 // 7. 不同数据类型
66 {
67     // 支持所有数值类型
68     auto input_float = torch::tensor({-2.5, 0.5, 2.5}, torch::kFloat32);
69     auto output_float = torch::clamp(input_float, 0.0, 2.0);
70
71     auto input_int = torch::tensor({-2, 0, 2}, torch::kInt32);
72     auto output_int = torch::clamp(input_int, 0, 2);
73 }
74

```

Function 接口

代码块

```

1 //clamp.out(Tensor self, Scalar? min=None, Scalar? max=None, *, Tensor(a!)
  out) -> Tensor(a!)
2 //clamp.Tensor_out(Tensor self, Tensor? min=None, Tensor? max=None, *,
  Tensor(a!) out) -> Tensor(a!)
3
4 /**

```

```

5  * @brief Clamps all elements in input into the range [min, max]
6  *
7  * @param dev      Device handle for device management
8  * @param out      Output data pointer
9  * @param self     Input data pointer
10 * @param min      Lower-bound tensor/scalar value
11 * @param max      Upper-bound tensor/scalar value
12 * @param minmax   Decide trigger clamp/clamp_min/clamp_max
13 * @param is_min_scalar  Decide the min is scalar/tensor
14 * @param is_max_scalar  Decide the max is scalar/tensor
15 *
16 * @note Supported Types:
17 *     - All floating point types (float, double)
18 *     - All integer types (uint8, int8, int16, int32, int64)
19 *     - Half precision (Half, BFloat16)
20 *
21 */
22 template<typename T>
23 void clamp_min_max_fn(
24     SiDeviceHandle dev,           // Device handle
25     void *out,                   // Output data pointer
26     const void *self,           // Input data pointer
27     const void *min,            // Lower bound tensor
28     const void *max,            // Upper bound tensor
29     int minmax,                 // mode (0 for clamp, 1 for clamp_min, 2 for
clamp_max, tmem)
30     bool is_min_scalar,         // Decide the min is scalar/tensor, true-
scalar false-tensor
31     bool is_max_scalar          // Decide the max is scalar/tensor, true-
scalar false-tensor
32 );
33 //该接口参考了void add_bf16(SiDeviceHandle dev, void *d_in0, void *d_in1, void
*d_out, int dim0_size, int dim1_size, sifmt::bfloat16 alpha, bool
is_d_in0_scalar, bool is_d_in1_scalar);
34 //可以处理min/max可能为tensor或scalar的情况

```

ISSUE

Dtype 支持

float32

float16

bfloat16

mxformat

非 contiguous 的是否能支持 (如果支持需要修改 Function 接口)

该算子是 pointwise 的，建议转换为 contiguous 后进行处理

Memory 是否需要 Align, Align 到多少

neg.out, neg, neg_

参考信息

pytorch 官方资料: <https://pytorch.org/docs/stable/generated/torch.neg.html>

cuda 实现参考:

<https://github.com/pytorch/pytorch/blob/94969d0a40e14eb36dd158d9805f48fb2f477cae/aten/src/ATen/native/cuda/UnarySignKernels.cu#L29>

代码块

```
1 import torch
2
3 # 1. 基本用法
4 x = torch.tensor([1, -2, 3])
5 y = torch.neg(x)
6 print(y) # tensor([-1, 2, -3])
7
8 # 2. 使用方法形式
9 x = torch.tensor([1, -2, 3])
10 y = x.neg()
11 print(y) # tensor([-1, 2, -3])
12
13 # 3. 使用负号运算符
14 x = torch.tensor([1, -2, 3])
15 y = -x
16 print(y) # tensor([-1, 2, -3])
17
18 # 4. 原地操作
19 x = torch.tensor([1, -2, 3])
20 x.neg_() # 注意后缀下划线表示原地操作
21 print(x) # tensor([-1, 2, -3])
22
23 # 5. 不同数据类型
24 x = torch.tensor([1.5, -2.5, 3.5], dtype=torch.float32)
25 y = torch.neg(x)
26 print(y) # tensor([-1.5000, 2.5000, -3.5000])
27
28 # 6. 复数类型
29 x = torch.tensor([1+2j, -3-4j])
30 y = torch.neg(x)
31 print(y) # tensor([-1.-2.j, 3.+4.j])
32
```

```

33 # 7. 多维张量
34 x = torch.tensor([[1, -2], [3, -4]])
35 y = torch.neg(x)
36 print(y) # tensor([[ -1,  2], [-3,  4]])
37
38 # 8. 在GPU上使用
39 if torch.cuda.is_available():
40     x_cuda = torch.tensor([1, -2, 3], device='cuda')
41     y_cuda = torch.neg(x_cuda)
42     print(y_cuda) # tensor([-1,  2, -3], device='cuda:0')

```

Function 接口

代码块

```

1  at::Tensor & wrapper_CUDA_neg_out_out(const at::Tensor & self, at::Tensor &
   out)
2
3  /**
4   * @brief Compute element-wise negation of a tensor
5   *
6   * @param dev      Device handle for CUDA device management
7   * @param out      Output tensor data pointer
8   * @param self     Input tensor data pointer
9   * @param sizes    Tensor sizes [ndim]
10  * @param strides  Tensor strides [ndim]
11  * @param ndim     Number of dimensions
12  * @param dtype    Data type of tensors
13  *
14  * @note Memory Layout:
15  *     - Input and output must have same shape
16  *     - Supports non-contiguous tensors
17  *     - Optimized for contiguous case
18  *
19  * @note Performance:
20  *     - Uses vectorized loads for aligned memory
21  *     - JIT compilation for complex types
22  *     - Different implementations for contiguous/strided
23  *
24  */
25  void fn(
26      SiDeviceHandle dev,          // Device handle
27      void *out,                  // Output data pointer
28      const void *self,           // Input data pointer
29      const int64_t* sizes,       // Tensor sizes [ndim]
30      const int64_t* strides,     // Tensor strides [ndim]
31      int ndim,                   // Number of dimensions

```

```

32     int dtype // Data type
33 );
34
35 support Dtype
36 所有基本数值类型 (通过 AT_DISPATCH_ALL_TYPES):
37 uint8_t (Byte)
38 int8_t (Char)
39 int16_t (Short)
40 int32_t (Int)
41 int64_t (Long)
42 float (Float)
43 double (Double)
44 半精度浮点类型 (通过 ScalarType::Half):
45 at::Half (float16)
46 BFloat16 类型 (通过 ScalarType::BFloat16):
47 at::BFloat16 (bfloat16)
48 复数类型 (通过 AT_DISPATCH_COMPLEX_TYPES_AND):
49 c10::complex<float> (ComplexFloat)
50 c10::complex<double> (ComplexDouble)
51 c10::complex<Half> (ComplexHalf)

```

ISSUE

Dtype 支持

float32

float16

bfloat16

mxformat

非 contiguous 的是否能支持 (如果支持需要修改 Function 接口)

Different implementations for contiguous/strided

Memory 是否需要 Align, Align 到多少

aten::as_strided

参考信息

官网文档: https://pytorch.org/docs/stable/generated/torch.as_strided.html#torch.as_strided

参考实现:

<https://github.com/pytorch/pytorch/blob/7ae0ce6360b6e4f944906502d20da24c04debee5/c10/core/TensorImpl.cpp#L846>

用法参考:

代码块 基本用法

```
2 import torch
3
4 # 创建一个连续的张量
5 x = torch.tensor([[1, 2, 3],
6                  [4, 5, 6]])
7 print(x.size())          # torch.Size([2, 3])
8 print(x.stride())       # (3, 1) - 表示移动一行需要步进3个元素, 移动一列需要步进1
   个元素
9
10 # 使用as_strided创建一个新视图
11 # 1. 转置操作
12 y = x.as_strided(size=(3, 2), stride=(1, 3))
13 print(y)
14 # 输出:
15 # tensor([[1, 4],
16 #         [2, 5],
17 #         [3, 6]])
18
19 # 2. 提取对角线元素
20 diag = x.as_strided(size=(2,), stride=(4,))
21 print(diag) # tensor([1, 5])
22
23 # 3. 使用storage_offset参数
24 # 从第二个元素开始, 创建一个1x2的视图
25 z = x.as_strided(size=(1, 2), stride=(3, 1), storage_offset=1)
26 print(z) # tensor([[2, 3]])
27
28 # 4. 创建重叠视图
29 patches = x.as_strided(size=(2, 2, 2), stride=(3, 1, 1))
30 print(patches)
31 # 结果包含2x2大小的重叠窗口
32 # tensor([[[[1, 2],
33 #          [2, 3]],
34 #         [[4, 5],
35 #          [5, 6]]]])
36
37 # 5. 实现卷积的滑动窗口
38 input = torch.arange(16).reshape(4, 4)
39 # tensor([[ 0,  1,  2,  3],
40 #         [ 4,  5,  6,  7],
41 #         [ 8,  9, 10, 11],
42 #         [12, 13, 14, 15]])
43 windows = input.as_strided(size=(2, 2, 3, 3),
44                             stride=(4, 1, 4, 1))
45 print(windows) # 第一个3x3窗口
46 # tensor([[[[ 0,  1,  2],
```

```

47 # [ 4, 5, 6],
48 # [ 8, 9, 10]],
49 # [[ 1, 2, 3],
50 # [ 5, 6, 7],
51 # [ 9, 10, 11]]],
52 # [[[ 4, 5, 6],
53 # [ 8, 9, 10],
54 # [12, 13, 14]],
55 # [[ 5, 6, 7],
56 # [ 9, 10, 11],
57 # [13, 14, 15]]]])

```

Function 接口

代码块

```

1  /**
2   * @brief Creates a view of the input tensor with specified size, stride, and
   optional storage offset
3   *
4   * @param dev Device handle for CUDA/CPU device management
5   * @param self Input tensor data pointer
6   * @param self_sizes Input tensor sizes [ndims]
7   * @param self_strides Input tensor strides [ndims]
8   * @param self_ndims Number of dimensions for input tensor
9   * @param out Output tensor data pointer (same storage as input)
10  * @param size Desired output sizes [out_ndims]
11  * @param stride Desired output strides [out_ndims]
12  * @param out_ndims Number of dimensions for output tensor
13  * @param storage_offset Storage offset for output tensor (in elements)
14  * @param dtype Data type of tensor
15  *
16  * @note Memory Layout:
17  * - Same underlying storage is shared between input and output
18  * - No memory allocation occurs
19  * - Output tensor accesses input tensor's storage according to
   size/stride/offset
20  *
21  * @note Safety:
22  * - Caller is responsible for ensuring valid size/stride/offset
   combinations
23  * - Invalid parameters may cause out-of-bounds memory access
24  *
25  * @note Examples:
26  * 1. Transpose a 2x3 tensor to 3x2:
27  * - Input: sizes=[2,3], strides=[3,1]
28  * - Output: sizes=[3,2], strides=[1,3]

```

```

29  *
30  *   2. Get diagonal elements of a square matrix:
31  *     - Input: sizes=[NxN]
32  *     - Output: sizes=[N], strides=[N+1]
33  *
34  *   3. Extract a subtensor with offset:
35  *     - Input: sizes=[H,W]
36  *     - Output: sizes=[h,w], storage_offset=(i*W+j)
37  *
38  */
39  void as_strided(
40      SiDeviceHandle *dev,           // Device handle
41      const void *self,              // Input tensor data pointer
42      const int64_t* self_sizes,     // Input tensor sizes
43      const int64_t* self_strides,   // Input tensor strides
44      int self_ndims,                // Number of input dimensions
45      void *out,                     // Output tensor data pointer (same as input)
46      const int64_t* size,           // Output tensor sizes
47      const int64_t* stride,         // Output tensor strides
48      int out_ndims,                 // Number of output dimensions
49      int64_t storage_offset = 0,    // Storage offset (in elements)
50      int dtype = kFloat              // Data type
51  );

```

ISSUE

Dtype 支持

view 操作，不涉及到计算，应该支持所有的 type

非 contiguous 的是否能支持（如果支持需要修改 Function 接口）

Memory 是否需要 Align，Align 到多少

eq.Scalar_out, eq.Scalar, eq_.Scalar, eq.Tensor_out, eq.Tensor,
eq_.Tensor, nq.Scalar_out, nq.Scalar, nq_.Scalar, nq.Tensor_out,
nq.Tensor, nq_.Tensor [TODO]

cuda 中 eq 和 nq 是通过同一个 kernel 实现的，逻辑上反转即可，scalar 和 tensor 可以通过 flag 来控制，实现统一算子接口

参考信息

<https://pytorch.org/docs/stable/generated/torch.eq.html>

<https://pytorch.org/docs/stable/generated/torch.ne.html#torch.ne>

实现参考

<https://github.com/pytorch/pytorch/blob/7ae0ce6360b6e4f944906502d20da24c04debee5/aten/>

代码块

```

1  eq
2  # 1. 张量与张量比较
3  x = torch.tensor([1, 2, 3])
4  y = torch.tensor([1, 1, 3])
5  result = torch.eq(x, y) # 或 x.eq(y) 或 x == y
6  # result: tensor([True, False, True])
7
8  # 2. 张量与标量比较
9  x = torch.tensor([1, 2, 3])
10 result = torch.eq(x, 2) # 或 x.eq(2) 或 x == 2
11 # result: tensor([False, True, False])
12
13 ne
14 # 1. 张量与张量比较
15 x = torch.tensor([1, 2, 3])
16 y = torch.tensor([1, 1, 3])
17 result = torch.ne(x, y) # 或 x.ne(y) 或 x != y
18 # result: tensor([False, True, False])
19
20 # 2. 张量与标量比较
21 x = torch.tensor([1, 2, 3])
22 result = torch.ne(x, 2) # 或 x.ne(2) 或 x != 2
23 # result: tensor([True, False, True])

```

Function 接口

代码块

```

1  /**
2   * @brief Performs element-wise equality/inequality comparison
3   *
4   * @param dev          Device handle for CUDA/CPU device management
5   * @param out          Output tensor (bool type) to store the result
6   * @param self         First input tensor
7   * @param other        Second input tensor or scalar
8   * @param out_sizes    output sizes
9   * @param out_strides  output strides (can be non-contiguous)
10 * @param self_sizes    First input sizes
11 * @param self_strides  First input strides (can be non-contiguous)
12 * @param other_sizes   Second input sizes (nullptr if scalar)
13 * @param other_strides Second input strides (can be non-contiguous)
14 * @param self_ndims    Number of self dimensions
15 * @param other_ndims   Number of other dimensions

```

```

16 * @param dtype          Data type of inputs
17 * @param is_other_scalar Whether other is a scalar
18 * @param op_type        0 for eq, 1 for ne
19 *
20 * @note Memory Layout:
21 *     - Fully supports non-contiguous tensors
22 *     - Handles arbitrary strides efficiently
23 *     - Optimized for both contiguous and non-contiguous cases
24 */
25 void compare_eq_ne(
26     SiDeviceHandle *dev,          // Device handle
27     void *out,                   // Output tensor data pointer (bool type)
28     const void *self,            // First input tensor data pointer
29     const void *other,           // Second input tensor/scalar data pointer
30     const int64_t* out_sizes,     // output sizes
31     const int64_t* sout_strides, // output strides
32     const int64_t* self_sizes,   // First input sizes
33     const int64_t* self_strides, // First input strides
34     const int64_t* other_sizes,  // Second input sizes (nullptr if scalar)
35     const int64_t* other_strides, // Second input strides (nullptr if scalar)
36     int self_ndims,              // Number of self dimensions
37     int other_ndims,             // Number of other dimensions
38     int dtype,                   // Data type of inputs (optional)
39     bool is_other_scalar,        // Whether other is a scalar
40     int op_type                  // 0 for eq, 1 for ne
41 );

```

支持的数据类型

基本数值类型

整型

46 kByte / kUInt8 (uint8)

47 kChar / kInt8 (int8)

48 kShort / kInt16 (int16)

49 kInt / kInt32 (int32)

50 kLong / kInt64 (int64)

浮点型

52 kFloat (float32)

53 kDouble (float64)

复数类型

55 kComplexFloat (complex64)

56 kComplexDouble (complex128)

57 kComplexHalf (complex32)

低精度类型

59 kHalf (float16)

60 kBFloat16 (brain floating point)

特殊类型

62 kBool (布尔)

```
63 AT_FLOAT8_TYPES (Float8类型, 包括 kFloat8_e4m3fn 和 kFloat8_e5m2)
64 高级无符号类型
65 AT_BAREBONES_UNSIGNED_TYPES
66 kUInt16 (uint16)
67 kUInt32 (uint32)
68 kUInt64 (uint64)
```

ISSUE

Dtype 支持

AT_EXPAND(AT_ALL_TYPES_AND_COMPLEX), kComplexHalf, kHalf, kBFloat16, kBool, AT_EXPAND(AT_FLOAT8_TYPES), AT_EXPAND(AT_BAREBONES_UNSIGNED_TYPES));

非 contiguous 的是否能支持 (如果支持需要修改 Function 接口)

pytorch 中使用了 TensorIteratorBase 传递 self, other, out, 需考虑非 contiguous 实现

考虑到 broadcast 的情况, out 和 self 的 size/stride 存在不一致的情况, 需要传入

Memory 是否需要 Align, Align 到多少

fill_.Scalar, fill_.Scalar_out, fill_.Tensor, fill_.Tensor_out

参考信息

文档

https://pytorch.org/docs/stable/generated/torch.Tensor.fill_.html

实现:

<https://github.com/pytorch/pytorch/blob/2e0c98ff058eb6661f991f05557a673187ead84/aten/src/ATen/native/cuda/FillKernel.cu#L22>

代码块

```
1 fill_.Scalar(Tensor(a!) self, Scalar value) -> Tensor(a!)
2 fill_.Tensor(Tensor(a!) self, Tensor value) -> Tensor(a!)
3
4 fill_.Tensor虽然支持输入tensor, 但是只能是标量张量, 有限制value.dim() == 0, 最终调用的
  还是标量的kernel
5
6 参考用法:
7 # 创建张量
8 x = torch.zeros(2, 3)
9 print(x)
10 # tensor([[0., 0., 0.],
11 #         [0., 0., 0.]])
```

```

12
13 # 填充标量值
14 x.fill_(5)
15 print(x)
16 # tensor([[5., 5., 5.],
17 #         [5., 5., 5.]])
18
19 # 填充不同类型的值
20 y = torch.zeros(2, 3, dtype=torch.int)
21 y.fill_(3.14) # 会被截断为整数
22 print(y)
23 # tensor([[3, 3, 3],
24 #         [3, 3, 3]])
25
26 # 创建张量
27 x = torch.zeros(2, 3)
28
29 # 创建填充后的新张量
30 y = torch.fill(x, 7)
31 print(x) # 原张量不变
32 # tensor([[0., 0., 0.],
33 #         [0., 0., 0.]])
34 print(y) # 新张量被填充
35 # tensor([[7., 7., 7.],
36 #         [7., 7., 7.]])
37
38 # 创建标量张量
39 scalar_tensor = torch.tensor(9.5)
40
41 # 使用标量张量填充
42 x = torch.zeros(2, 3)
43 x.fill_(scalar_tensor)
44 print(x)
45 # tensor([[9.5000, 9.5000, 9.5000],
46 #         [9.5000, 9.5000, 9.5000]])

```

Function 接口

代码块

```

1 /**
2  * @brief Fills a tensor with a scalar value
3  *
4  * @param dev Device handle for CUDA/CPU device management
5  * @param self Tensor data pointer to be filled
6  * @param self_sizes Tensor sizes [ndims]

```

```

7  * @param self_strides Tensor strides [ndims] (supports non-contiguous
   tensors)
8  * @param ndims        Number of dimensions
9  * @param value        Scalar value to fill the tensor with
10 * @param dtype        Data type of the tensor (optional)
11 *
12 * @note Supported Types:
13 *     - All numeric types (float, int, etc.)
14 *     - Complex types
15 *     - Boolean type
16 *     - Half precision (float16, bfloat16)
17 *     - Float8 types
18 *     - Unsigned types
19 *
20 * @note Memory Layout:
21 *     - Supports both contiguous and non-contiguous tensors
22 *     - Handles arbitrary strides efficiently
23 *
24 * @note Special Cases:
25 *     - For single-element CPU tensors, a specialized fast path is used
26 *
27 * @return int 0 on success, error code otherwise
28 */
29 int fill(
30     SiDeviceHandle *dev,           // Device handle
31     void *self,                   // Tensor data pointer
32     const int64_t* self_sizes,     // Tensor sizes
33     const int64_t* self_strides,   // Tensor strides
34     int ndims,                    // Number of dimensions
35     const void* value,             // Scalar value pointer
36     int dtype                      // Data type of tensor (optional)
37 );

```

ISSUE

- Dtype 支持
- 非 contiguous 的是否能支持 (如果支持需要修改 Function 接口)
 - pytorch 中使用了 TensorIteratorBase 传递 self, 需考虑非 contiguous 实现
- Memory 是否需要 Align, Align 到多少
- 对于 tensor/scalar
 - 如果使用 fill_stub, 会通过 pytorch 的原始实现来对 tensor 进行处理, 可以直接获取 (TensorIterator& iter, const Scalar& value) 作为参数

- 如果不使用 fill_stub，需要自己手动处理不同的参数，对 tensor 进行判断，提取标量后和 scalar 进行相同处理

div.out, div_.Tensor, div.Tensor, div.Scalar, div_.Scalar

参考信息

参考实现:

<https://github.com/pytorch/pytorch/blob/9cd52da45c3030601bdd27e8a948f23be12b21ee/aten/src/ATen/native/cuda/BinaryDivTrueKernel.cu#L20>

资料: <https://pytorch.org/docs/stable/generated/torch.div.html#torch.div>

cuda 中只注册了 div.out, div_.Tensor, div.Tensor, cuda 通过 div_true_stub 注册 div, cuda 的 kernel 中有针对第二个输入是 scalar 的情况优化, 会将 scalar 取倒数后按照乘法计算

代码块

```
1 # 张量与张量除法
2 x = torch.tensor([10.0, 20.0, 30.0])
3 y = torch.tensor([2.0, 5.0, 10.0])
4
5 # 三种等价方式
6 result1 = torch.div(x, y)
7 result2 = x / y
8 result3 = x.div(y)
9
10 print(result1) # tensor([5.0000, 4.0000, 3.0000])
11
12 a = torch.tensor([10.0, 20.0, 30.0])
13 c = torch.div(a, 2.0)
14 print(c) # tensor([5.0000, 10.0000, 15.0000])
15
16 # 使用运算符
17 c = a / 2.0
18 print(c) # tensor([5.0000, 10.0000, 15.0000])
```

Function 接口

代码块

```
1 div.Tensor(Tensor self, Tensor other) -> Tensor
2 div_.Tensor(Tensor(a!) self, Tensor other) -> Tensor(a!)
3 div.out(Tensor self, Tensor other, *, Tensor(a!) out) -> Tensor(a!)
4
5 /**
```

```

6 * @brief Performs element-wise division of tensor by another tensor or scalar
7 *
8 * @param dev Device handle for CUDA/CPU device management
9 * @param self Dividend tensor (numerator)
10 * @param other Divisor tensor or scalar (denominator)
11 * @param self_sizes Tensor sizes of self [ndims]
12 * @param self_strides Tensor strides of self [ndims]
13 * @param other_sizes Tensor sizes of other [ndims] (nullptr if scalar)
14 * @param other_strides Tensor strides of other [ndims] (nullptr if scalar)
15 * @param ndims Number of dimensions
16 * @param out Output tensor (can be same as self for in-place)
17 * @param out_sizes Output tensor sizes [ndims]
18 * @param out_strides Output tensor strides [ndims]
19 * @param is_self_scalar
20 *
21 * @note Supported Types:
22 *     - Floating point types (float, double)
23 *     - Complex types (complex<float>, complex<double>)
24 *     - Half precision (float16, bfloat16)
25 *
26 * @note Performance Optimizations:
27 *     - Special fast path for scalar divisors
28 *     - For scalar division, computes  $a * (1/b)$  instead of  $a / b$ 
29 *     - Handles non-contiguous tensors efficiently
30 *
31 * @note Broadcast Support:
32 *     - Automatically handles broadcasting between tensors
33 *     - Caller must ensure output tensor has correct broadcast shape
34 *
35 * @return int 0 on success, error code otherwise
36 */
37 int div(
38     SiDeviceHandle dev,
39     const void* self,
40     const void* other,
41     const int64_t* self_sizes,
42     const int64_t* self_strides,
43     const int64_t* other_sizes,
44     const int64_t* other_strides,
45     int self_ndims,
46     int other_ndims,
47     void* out,
48     const int64_t* out_sizes,
49     const int64_t* out_strides
50     bool is_self_scalar
51 );
52

```

- 53 支持数据类型
- 54 浮点类型:
- 55 float32 (torch.float)
- 56 float64 (torch.double)
- 57 float16 (torch.half)
- 58 bfloat16 (torch.bfloat16)
- 59 复数类型:
- 60 complex64 (torch.complex64)
- 61 complex128 (torch.complex128)
- 62 CUDA 实现不直接支持整数类型 (如 int8, int16, int32, int64), 因为在 PyTorch 中, 整数除法会返回浮点结果 (真除法行为)。

首个版本开发说明

- (1) 数制支持上, 初版需要支持float, bf16类型
- (2) 除数 (输入2) 需要支持tensor与scalar两种形式
- (3) 输出与被除数 (输入1) 的shape一致
- (4) 如果除数为tensor, 那么除数与被除数shape应相同 (如不同, 则由框架处理进行broadcast), 做elementwise除法
- (5) 对齐限制: 硬件目前要求dim0_size * sizeof(T)为1024B对齐
- (6) 输入输出参数, 均由框架进行预处理 (reshape为一维并padding至符合对齐限制)

内部接口调用形式:

代码块

```
1  template <typename T>
2  void div(void *d_in0, void *d_in1, void *d_out, int dim0_size, bool
    is_d_in1_scalar);
```

ISSUE

- 非 contiguous 的是否能支持 (如果支持需要修改 Function 接口)
 - 需要考虑到 broadcast 和第二个输入为 scalar 的情况。需要 stride, size, dim 信息

exponential_

参考信息

参考实现:

<https://github.com/pytorch/pytorch/blob/main/aten/src/ATen/native/cuda/DistributionExponentialKernel.cu>

<https://github.com/pytorch/pytorch/blob/6055a4f612782ca944f2e0465f7497b7f18de4e9/aten/src/ATen/native/cuda/DistributionTemplates.h#L558>

资料: https://pytorch.org/docs/stable/generated/torch.Tensor.exponential_.html

代码块

```
1 import torch
2
3 # 创建一个张量
4 x = torch.tensor([1.0, 2.0, 3.0])
5
6 # 原地进行指数变换
7 x.exponential_()
8
9 print(x) # 输出: tensor([ 2.7183,  7.3891, 20.0855])
```

Function 接口

代码块

```
1 (Tensor(a!)) self, float lambda=1, *, Generator? generator=None) -> Tensor(a!)
2
3 /**
4  * @brief Fills the input tensor with samples from an exponential distribution.
5  *
6  * @param self The output tensor to fill with samples.
7  * @param lambda The rate parameter (lambda) of the exponential
8  * distribution.
9  * @param gen Optional random number generator for reproducibility.
10 *
11 * @return void* ptr to the output tensor filled with samples.
12 */
13 void* fn(
14     SiDeviceHandle *dev,
15     void* self,
16     double lambda,
17     std::optional<Generator> gen = std::nullopt
18     const int64_t* self_sizes,
19     const int64_t* self_strides,
20     int ndims,
21 );
22
23 (update)
24
25
```

- 26 支持的数据类型:
- 27 浮点类型: float32, float64, float16, bfloat16
- 28 复数类型: complex64, complex128

ISSUE

非 contiguous 的是否能支持 (如果支持需要修改 Function 接口)

通过 exponential_stub 注册, 支持非 contiguous 数据, 需要提供 size, stride, dim 信息

`addmm.out, addmm, mm.out, mm, _addmm_activation.out, _addmm_activation`

参考信息

文档: <https://pytorch.org/docs/stable/generated/torch.mm.html>

<https://pytorch.org/docs/stable/generated/torch.addmm.html>

Addmm:

$$\text{out} = \beta \text{input} + \alpha (\text{mat1}_i @ \text{mat2}_i)$$

参考实现:

<https://github.com/pytorch/pytorch/blob/919d54b7b1bcba4dfd425675f7c29438e1435469/aten/src/ATen/native/cuda/Blas.cpp#L318>

代码块

```
1
2 TORCH_IMPL_FUNC(addmm_out_cuda)(const Tensor& self, const Tensor& mat1, const
  Tensor& mat2, const Scalar& beta, const Scalar& alpha, const Tensor& result) {
3   addmm_out_cuda_impl(const_cast<Tensor&>(result), self, mat1, mat2, beta,
  alpha);
4 }
5
6 TORCH_IMPL_FUNC(addmm_activation_out_cuda)(const Tensor& self, const Tensor&
  mat1, const Tensor& mat2, const Scalar& beta, const Scalar& alpha, bool
  use_gelu, const Tensor& result) {
7   addmm_out_cuda_impl(const_cast<Tensor&>(result), self, mat1, mat2, beta,
  alpha, use_gelu ? Activation::GELU : Activation::RELU);
8 }
9
10 TORCH_IMPL_FUNC(mm_out_cuda)(const Tensor& self, const Tensor& mat2, const
  Tensor& result) {
11   addmm_out_cuda_impl(const_cast<Tensor&>(result), result, self, mat2, 0, 1);
```

```
12 }
13
14 三个op通过同一个kernel实现
15
16 >>> A = torch.tensor([[1, 2], [3, 4]])
17 >>> B = torch.tensor([[5, 6], [7, 8]])
18 >>> input = torch.tensor([[1, 1], [1, 1]])
19
20 >>> C = torch.mm(A,B)
21 >>> C
22 tensor([[19, 22],
23         [43, 50]])
24
25 >>> C = torch.addmm(input, A, B)
26 >>> C
27 tensor([[20, 23],
28         [44, 51]])
29
30 >>> C = torch.addmm(input, A, B, beta=1, alpha=2.0)
31 >>> C
32 tensor([[ 39,  45],
33         [ 87, 101]])
34
35 torch._addmm_activation(input, mat1, mat2, beta=0.5, alpha=2.0,use_gelu=True)
36 等价于
37 gelu(beta * input + alpha * torch.mm(mat1, mat2))
38
39 torch._addmm_activation(input, mat1, mat2, beta=0.5, alpha=2.0,use_gelu=False)
40 等价于
41 relu(beta * input + alpha * torch.mm(mat1, mat2))
42
43 >>> result = torch._addmm_activation(input, mat1, mat2, use_gelu=True)
44 >>> equivalent = gelu(input + torch.mm(mat1, mat2))
45 >>> result
46 tensor([[ 2.7273e+00,  1.2938e+00,  4.7260e-02, -2.6163e-03],
47         [ 1.1547e+00, -1.3293e-01, -9.3054e-02,  1.9921e-01],
48         [-1.1569e-04,  2.0479e+00,  6.1474e-01,  1.1607e+00]])
49 >>> equivalent
50 tensor([[ 2.7273e+00,  1.2938e+00,  4.7260e-02, -2.6163e-03],
51         [ 1.1547e+00, -1.3293e-01, -9.3054e-02,  1.9921e-01],
52         [-1.1569e-04,  2.0479e+00,  6.1474e-01,  1.1607e+00]])
53
54 >>> result = torch._addmm_activation(input, mat1, mat2, beta=0.5,
55 alpha=2.0,use_gelu=True)
56 >>> equivalent = gelu(0.5 * input + 2.0 * torch.mm(mat1, mat2))
57 >>> result
58 tensor([[ 3.1910e+00, -1.6997e-01, -6.1870e-02, -1.6937e-02],
```

```

58     [ 2.7823e+00, -2.6890e-02, -3.9319e-03, -1.6966e-01],
59     [-6.1636e-06, 4.7033e+00, 2.5606e+00, 3.8821e+00]])
60 >>> equivalent
61 tensor([[ 3.1910e+00, -1.6997e-01, -6.1870e-02, -1.6937e-02],
62         [ 2.7823e+00, -2.6890e-02, -3.9319e-03, -1.6966e-01],
63         [-6.1636e-06, 4.7033e+00, 2.5606e+00, 3.8821e+00]])

```

Function 接口

代码块

```

1  /**
2   * @brief 执行矩阵乘法并加偏置, 可选激活函数:  $out = beta * bias + alpha * (A @ B)$ 
3   *
4   * @param dev 设备句柄, 用于 CUDA 设备管理
5   * @param a_ptr 第一个矩阵的数据指针 (A)
6   * @param b_ptr 第二个矩阵的数据指针 (B)
7   * @param bias_ptr 偏置张量的数据指针, 用于加到矩阵乘法结果
8   * @param a_sizes A 矩阵的维度 [2], 形状为 (m, k)
9   * @param a_strides A 矩阵的步长 [2]
10  * @param b_sizes B 矩阵的维度 [2], 形状为 (k, n)
11  * @param b_strides B 矩阵的步长 [2]
12  * @param bias_sizes 偏置张量的维度, 可以是 [2] 或 [1]
13  * @param bias_strides 偏置张量的步长
14  * @param bias_ndims 偏置张量的维数 (1 或 2)
15  * @param out_ptr 输出张量的数据指针
16  * @param out_sizes 输出张量的维度 [2], 形状为 (m, n)
17  * @param out_strides 输出张量的步长 [2]
18  * @param dtype 所有张量的数据类型 ID
19  * @param alpha A @ B 的缩放因子
20  * @param beta 偏置的缩放因子
21  * @param activation 激活函数类型: 0=None, 1=ReLU, 2=GELU
22  *
23  * @note 支持的数据类型:
24  * - 浮点类型 (float, double)
25  * - 半精度类型 (half, bfloat16)
26  * - 8位浮点类型 (float8_e4m3fn, float8_e5m2)
27  *
28  * @note 内存布局:
29  * - 支持各种内存布局 (行优先、列优先、非连续)
30  * - 自动处理转置和共轭转置
31  * - 如果 A 或 B 是非连续的, 可能会产生临时副本
32  *
33  * @note 张量广播:
34  * - 偏置可以是向量 (1D) 或矩阵 (2D)
35  * - 如果偏置是向量, 会被广播到结果形状
36  *

```

```

37  */
38  void addmm_fn(
39      SiDeviceHandle dev,
40      const void* a_ptr,
41      const void* b_ptr,
42      const void* bias_ptr,
43      const int64_t* a_sizes,
44      const int64_t* a_strides,
45      const int64_t* b_sizes,
46      const int64_t* b_strides,
47      const int64_t* bias_sizes,
48      const int64_t* bias_strides,
49      int bias_ndims,
50      void* out_ptr,
51      const int64_t* out_sizes,
52      const int64_t* out_strides,
53      int dtype,
54      float alpha,
55      float beta,
56      int activation
57  );
58
59  默认a,b,out的ndims均为2
60
61  支持的输入类型
62  1. 浮点类型
63  float32 (torch.float, at::kFloat) - 单精度浮点
64  float64 (torch.double, at::kDouble) - 双精度浮点
65  2. 半精度类型
66  float16 (torch.half, at::kHalf) - 半精度浮点
67  bfloat16 (torch.bfloat16, at::kBFloat16) - Brain 浮点格式
68  3. 低精度类型 (特定 GPU 架构支持)
69  float8_e4m3fn (torch.float8_e4m3fn, at::kFloat8_e4m3fnuz) - 8位浮点 (4位指数, 3位尾数)
70  float8_e5m2 (torch.float8_e5m2, at::kFloat8_e5m2fnuz) - 8位浮点 (5位指数, 2位尾数)

```

ISSUE

非 contiguous 的是否能支持 (如果支持需要修改 Function 接口)

addmm_out_cuda_impl 使用 prepare_matrix_for_cublas 方法来处理非连续输入

index.Tensor,index.Tensor_out, index_put_impl, _index_put_impl,
_index_put_impl.out, index_put_, index_put.out

参考信息

文档:

https://pytorch.org/docs/stable/generated/torch.Tensor.index_put_.html#torch.Tensor.index_put_

参考实现:

<https://github.com/pytorch/pytorch/blob/981807cfc2dc8036f0e83f038455b49ba4b833/aten/src/ATen/native/cuda/IndexKernel.cu#L196>

<https://github.com/pytorch/pytorch/blob/981807cfc2dc8036f0e83f038455b49ba4b833/aten/src/ATen/native/cuda/IndexKernel.cu#L248>

代码块

```
1  >>> x = torch.tensor([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
2  >>> row_indices = torch.tensor([0, 2]) # 选择第0行和第2行
3  >>> result = x[row_indices]
4  >>> result
5  tensor([[1, 2, 3],
6         [7, 8, 9]])
7  >>> row_indices = torch.tensor([0, 1, 2])
8  >>> col_indices = torch.tensor([1, 2, 0])
9  >>> result = x[row_indices, col_indices]
10 >>> result
11 tensor([2, 6, 7])
12 >>>
13
14 >>> x = torch.zeros(3, 3) # 默认是 float 类型
15 >>> indices = (torch.tensor([0, 2]), torch.tensor([1, 1]))
16 >>> values = torch.tensor([10, 20], dtype=torch.float) # 显式指定为浮点类型
17 >>> result = x.index_put(indices, values)
18 >>>
19 >>> result
20 tensor([[ 0., 10.,  0.],
21         [ 0.,  0.,  0.],
22         [ 0., 20.,  0.]])
23
24 >>> x = torch.ones(3, 3) # 浮点类型
25 >>> x
26 tensor([[1., 1., 1.],
27         [1., 1., 1.],
28         [1., 1., 1.]])
29 >>> indices = (torch.tensor([0, 0]), torch.tensor([1, 1])) # 索引位置 (0,1) 重复
30 >>> values = torch.tensor([10.0, 20.0]) # 浮点类型匹配 x
31 >>> result = x.index_put(indices, values, accumulate=True)
32 >>> result
33 tensor([[ 1., 31.,  1.],
```

```

34     [ 1., 1., 1.],
35     [ 1., 1., 1.]]
36 >>>
37
38 cuda实现中共用一个入口
39 template <typename scalar_t>
40 void index_kernel_impl(TensorIteratorBase& iter, const IntArrayRef index_size,
41     const IntArrayRef index_stride) {
42     gpu_index_kernel(iter, index_size, index_stride, []C10_DEVICE(char* const
43     out_data, const char* const in_data, const int64_t offset) {
44         *reinterpret_cast<scalar_t*>(out_data) = *reinterpret_cast<const scalar_t*>
45         (in_data + offset);
46     });
47 }
48
49 template <typename scalar_t>
50 void index_put_kernel_impl(TensorIterator& iter, const IntArrayRef index_size,
51     const IntArrayRef index_stride) {
52     gpu_index_kernel(iter, index_size, index_stride, []C10_DEVICE(char* const
53     out_data, const char* const in_data, const int64_t offset) {
54         *reinterpret_cast<scalar_t*>(out_data + offset) = *reinterpret_cast<const
55         scalar_t*>(in_data);
56     });
57 }
58

```

Function 接口

代码块

```

1  /**
2   * @brief 执行高级索引操作 (index 和 index_put)
3   *
4   * @param dev          设备句柄, 管理 CUDA 上下文
5   * @param out          结果张量 (对于 index 是输出, 对于 index_put 是被索引的目
6   * @param indices      索引张量数组, 每个维度一个索引, 使用 nullptr 表示不索引的维
7   * @param values
8   * @param self        源张量 (对于 index 是被索引的源, 对于 index_put 是输入值)
9   * @param out_ndims
10  * @param out_sizes
11  * @param out_strides
12  * @param indices_ndims
13  * @param indices_sizes 每个索引张量的尺寸数组, 表示有效索引范围
14  * @param indices_strides 每个索引张量的步长数组, 配合索引使用
15  * @param num_indices 索引张量的数量

```

```

16 * @param values_ndims
17 * @param values_sizes
18 * @param values_strides
19 * @param self_ndims
20 * @param self_sizes
21 * @param self_strides
22 * @param is_put true=index_put 操作, false=index 操作
23 * @param accumulate 当 is_put=true 时是否累加值而不是替换 (仅对 index_put 有
效)
24 *
25 * @note 支持的索引类型:
26 * - 整数索引张量: 选择特定索引位置的数据
27 * - nullptr: 表示该维度不进行索引, 相当于整个维度都被选择
28 * - 布尔掩码: 通过布尔值为 true 的位置进行索引
29 *
30 * @note 边界处理:
31 * - 索引从 0 到 (dimension_size-1)
32 * - 负索引在内部被修正: -1 表示最后一个元素
33 * - 索引超出范围时会抛出异常
34 *
35 * @note 内存布局:
36 * - 支持非连续输入和输出张量
37 * - 内部高效处理各种跨步和偏移计算
38 *
39 * @return int 0表示成功, 非零表示错误码
40 */
41 void fn(
42     SiDeviceHandle dev,
43     const void* self,
44     const void** indices,
45     const void* values,
46     void* out,
47     const int64_t self_ndims,
48     const int64_t* self_sizes,
49     const int64_t* self_strides,
50     const int64_t* indices_ndims,
51     const int64_t* indices_sizes,
52     const int64_t* indices_strides,
53     int64_t num_indices,
54     const int64_t values_ndims,
55     const int64_t* values_sizes,
56     const int64_t* values_strides,
57     const int64_t out_ndims,
58     const int64_t* out_sizes,
59     const int64_t* out_strides,
60     bool is_put,
61     bool accumulate

```

```
62 );
63
64 index 操作支持:
65 所有基本类型: int8, int16, int32, int64, uint8, float, double
66 所有复数类型: complex64, complex128
67 半精度类型: float16 (kHalf)
68 布尔类型: bool (kBool)
69 BFloat16: bfloat16 (kBFloat16)
70 复数半精度: complex32 (kComplexHalf)
71
72 index_put 操作支持所有 index 支持的类型, 并额外支持:
73 8位浮点类型: float8_e4m3fn, float8_e5m2 (AT_FLOAT8_TYPES)
```

ISSUE

- 非 contiguous 的是否能支持 (如果支持需要修改 Function 接口)
- CUDA 实现的 index_put 不支持 accumulate=true 参数

gt.Tensor_out gt.Tensor gt_.Tensor gt.Scalar_out gt.Scalar gt_.Scalar,
lt.Tensor_out lt.Tensor lt_.Tensor lt.Scalar_out lt.Scalar lt_.Scalar

代码块

```
1 //Cuda中 gt lt le使用了同一个接口, 通过设定一个flag来确认具体的method
2
3 void gt_kernel_cuda(TensorIteratorBase& iter) {
4     compare_kernel_with_scalars(iter, OpType::GT);
5 }
6
7 void le_kernel_cuda(TensorIteratorBase& iter) {
8     compare_kernel_with_scalars(iter, OpType::LE);
9 }
10
11 void lt_kernel_cuda(TensorIteratorBase& iter) {
12     compare_kernel_with_scalars(iter, OpType::LT);
13 }
```

gt.Tensor_out gt.Tensor gt_.Tensor gt.Scalar_out gt.Scalar gt_.Scalar

参考信息

参考实现:

<https://github.com/pytorch/pytorch/blob/9cd52da45c3030601bdd27e8a948f23be12b21ee/aten/src/ATen/native/cuda/CompareKernels.cu#L86>

资料:

<https://pytorch.org/docs/stable/generated/torch.gt.html>

代码块

```
1  gt.Scalar_out(Tensor self, Scalar other, *, Tensor(a!)) out) -> Tensor(a!)
2  gt.Tensor_out(Tensor self, Tensor other, *, Tensor(a!)) out) -> Tensor(a!)
3
4  # 张量与张量比较
5  x = torch.tensor([1, 2, 3, 4])
6  y = torch.tensor([2, 2, 2, 2])
7
8  # 三种等价方式
9  result1 = torch.gt(x, y)
10 result2 = x > y
11 result3 = x.gt(y)
12
13 print(result1) # tensor([False, False, True, True])
14
15
16 x = torch.tensor([1, 2, 3, 4])
17 # 比较张量中的每个元素是否大于2
18 result = x > 2
19 print(result) # tensor([False, False, True, True])
20
21
22 # 不同形状的张量比较 (使用广播)
23 x = torch.tensor([[1, 2], [3, 4]])
24 y = torch.tensor([2, 1])
25
26 result = x > y
27 print(result)
28 # tensor([[False, True],
29 #         [True, True]])
```

lt.Tensor_out lt.Tensor lt_.Tensor lt.Scalar_out lt.Scalar lt_.Scalar

参考信息

参考实现:

<https://github.com/pytorch/pytorch/blob/9cd52da45c3030601bdd27e8a948f23be12b21ee/aten/src/ATen/native/cuda/CompareKernels.cu#L94>

文档: <https://pytorch.org/docs/stable/generated/torch.lt.html>

代码块

```

1  import torch
2
3  a = torch.tensor([1, 2, 3, 4])
4  b = torch.tensor([2, 2, 2, 2])
5
6  # 返回布尔张量 运算符形式
7  result = a < b
8  print(result) # tensor([True, False, False, False])
9
10
11 # 张量与张量比较
12 a = torch.tensor([1, 2, 3, 4])
13 b = torch.tensor([2, 2, 2, 2])
14 result = torch.lt(a, b)
15 print(result) # tensor([True, False, False, False])
16
17 # 张量与标量比较
18 result = torch.lt(a, 2)
19 print(result) # tensor([True, False, False, False])
20
21 # 批量比较
22 batch = torch.tensor([[1, 2, 3], [4, 5, 6]])
23 threshold = torch.tensor([2, 4, 3])
24
25 result = batch < threshold
26 print(result)
27 # tensor([[ True,  True, False],
28 #         [False, False, False]])

```

le.Scalar_out, le.Scalar, le.Tensor_out, le.Tensor

文档: <https://pytorch.org/docs/stable/generated/torch.le.html>

代码块

```

1  import torch
2
3  # 创建两个张量
4  a = torch.tensor([1, 2, 3, 4])
5  b = torch.tensor([2, 2, 2, 2])
6
7  # 使用 torch.le 比较
8  result = torch.le(a, b)
9  print("a <= b:", result)
10 # a <= b: tensor([True, True, False, False])
11
12 # 使用张量方法比较

```

```

13 result_method = a.le(b)
14 print("a.le(b):", result_method)
15 # a.le(b): tensor([ True,  True, False, False])
16
17 # 创建一个张量
18 x = torch.tensor([[1, 2, 3], [4, 5, 6]])
19
20 # 与标量比较
21 result = torch.le(x, 3)
22 print(result)
23 # tensor([[ True,  True,  True],
24 #         [False, False, False]])
25
26 # 不同形状的张量比较
27 a = torch.tensor([[1, 2], [3, 4]])
28 b = torch.tensor([2, 1]) # 将被广播为 [[2, 1], [2, 1]]
29
30 result = torch.le(a, b)
31 print(result)
32 # tensor([[ True, False],
33 #         [False, False]])

```

Function 接口

代码块

```

1 /**
2  * @brief Performs element-wise gt/lt/le comparison
3  *
4  * @param dev Device handle for CUDA/CPU device management
5  * @param out Output tensor (bool type) to store the result
6  * @param self First input tensor
7  * @param other Second input tensor or scalar
8  * @param out_sizes output sizes
9  * @param out_strides output strides (can be non-contiguous)
10 * @param self_sizes First input sizes
11 * @param self_strides First input strides (can be non-contiguous)
12 * @param other_sizes Second input sizes (nullptr if scalar)
13 * @param other_strides Second input strides (can be non-contiguous)
14 * @param self_ndims Number of self dimensions
15 * @param other_ndims Number of other dimensions
16 * @param dtype Data type of inputs
17 * @param is_other_scalar Whether other is a scalar
18 * @param OpType 0 for gt, 1 for lt, 2 for le
19 *
20 * @note Memory Layout:
21 * - Fully supports non-contiguous tensors

```

```

22  *      - Handles arbitrary strides efficiently
23  *      - Optimized for both contiguous and non-contiguous cases
24  */
25  void compare_gt(
26      SiDeviceHandle *dev,           // Device handle
27      void *out,                     // Output tensor data pointer (bool type)
28      const void *self,              // First input tensor data pointer
29      const void *other,             // Second input tensor/scalar data pointer
30      const int64_t* out_sizes,      // output sizes
31      const int64_t* sout_strides,   // output strides
32      const int64_t* self_sizes,     // First input sizes
33      const int64_t* self_strides,   // First input strides
34      const int64_t* other_sizes,    // Second input sizes (nullptr if scalar)
35      const int64_t* other_strides,  // Second input strides (nullptr if scalar)
36      int self_ndims,                // Number of self dimensions
37      int other_ndims,               // Number of other dimensions
38      int dtype,                     // Data type of inputs (optional)
39      bool is_other_scalar,          // Whether other is a scalar
40      int OpType                      // 0 for gt, 1 for lt, 2 for le
41  );

```

ISSUE

Dtype 支持

AT_EXPAND(AT_ALL_TYPES_AND_COMPLEX), kComplexHalf, kHalf, kBFloat16, kBool, AT_EXPAND(AT_FLOAT8_TYPES), AT_EXPAND(AT_BAREBONES_UNSIGNED_TYPES));

非 contiguous 的是否能支持（如果支持需要修改 Function 接口）

pytorch 中使用了 TensorIteratorBase 传递 self, other, out, 需考虑非 contiguous 实现

考虑到 broadcast 的情况, out 和 self 的 size/stride 存在不一致的情况, 需要传入

Memory 是否需要 Align, Align 到多少

resize_, resize_output, resize.out, resize

参考信息

文档: https://pytorch.org/docs/stable/generated/torch.Tensor.resize_.html

参考实现:

<https://github.com/pytorch/pytorch/blob/32f585d9346e316e554c8d9bf7548af9f62141fc/aten/src/ATen/native/cuda/Resize.cpp#L17>

可以观察到在 cuda 的实现实际上就是 `cudaMemcpyAsync`

代码块

```
1 import torch
2
3 # 创建张量
4 x = torch.randn(2, 3)
5 print(f"原始张量: {x.shape}")
6 # 原始张量: torch.Size([2, 3])
7
8 # 调整大小为 3x2
9 x.resize_(3, 2)
10 print(f"调整后张量: {x.shape}")
11 # 调整后张量: torch.Size([3, 2])
12
13 # 创建一个张量
14 x = torch.tensor([1, 2, 3, 4, 5, 6]).view(2, 3)
15 print(f"原始张量:\n{x}")
16 # 原始张量:
17 # tensor([[1, 2, 3],
18 #         [4, 5, 6]])
19
20
21 # 增加大小 - 新元素值不可预测!
22 x.resize_(3, 3)
23 print(f"增加大小后(新值不确定):\n{x}")
24 # 增加大小后(新值不确定):
25 # tensor([[1, 2, 3],
26 #         [4, 5, 6],
27 #         [?, ?, ?]]) # 新元素的值是未定义的
28
29
30 # 减少大小 - 截断元素
31 y = torch.tensor([1, 2, 3, 4, 5, 6]).view(2, 3)
32 y.resize_(1, 3)
33 print(f"减少大小后:\n{y}")
34 # 减少大小后:
35 # tensor([[1, 2, 3]])
36
37 tensor = torch.tensor([1, 2, 3, 4, 5, 6])
38
39 # 1. resize_: 实际改变底层存储, 可增加或减少元素
40 tensor_resized = tensor.clone()
41 tensor_resized.resize_(2, 4) # 增加元素, 新元素未初始化
42
43 # 2. view: 仅改变视图, 总元素数必须相同
```

```

44  tensor_view = tensor.view(2, 3) # 必须是相同元素数
45
46  # 3. reshape: 类似view, 但可处理非连续张量
47  tensor_reshaped = tensor.reshape(3, 2) # 可能创建副本
48
49  # 4. contiguous+view: 确保连续后改变视图
50  tensor_cont = tensor.contiguous().view(2, 3)
51
52  # 大多数情况下, 推荐使用 view, reshape 或 permute 来调整张量形状, 这些方法更安全且保留
    梯度关系。resize_ 主要用于需要实际改变内存分配的特殊场景。

```

cuda 的处理流程 (大部分工作在 torch 侧)

1. 入口点: resize_cuda_函数

代码块

```

1  const Tensor& resize_cuda_(const Tensor& self, IntArrayRef size,
    std::optional<MemoryFormat> optional_memory_format) {
2      // 处理具名张量
3      if (self.has_names()) {
4          return resize_named_tensor_(self, size, optional_memory_format);
5      }
6
7      // 记录旧存储大小 (用于确定性算法)
8      auto* self_ = self.unsafeGetTensorImpl();
9      int64_t old_storage_nbytes = self_->unsafe_storage() ? self_-
    >unsafe_storage().nbytes() : 0;
10
11     // 调整实际大小
12     resize_impl_cuda_(self_, size, /*strides=*/std::nullopt);
13
14     // 处理内存格式调整
15     if (optional_memory_format.has_value()) {
16         // ...重新设置步长
17         self_->empty_tensor_restride(memory_format);
18     }
19
20     // 确定性内存初始化 (如果需要)
21     if (C10_UNLIKELY(at::globalContext().deterministicAlgorithms() && ...)) {
22         at::native::fill_resize_deterministic_(self, old_storage_nbytes);
23     }
24
25     return self;
26 }

```

2. 核心实现: `resize_impl_cuda_`函数

代码块

```
1  inline TensorImpl* resize_impl_cuda_(TensorImpl* self, IntArrayRef size,
2  at::OptionalIntArrayRef stride) {
3      // 如果大小和步长没变, 直接返回
4      if (self->sizes() == size && (!stride || self->strides() == stride)) {
5          return self;
6      }
7      // 计算元素大小和存储偏移
8      const auto itemsize = self->dtype().itemsize();
9      const auto storage_offset = self->storage_offset();
10
11     size_t storage_size = 1;
12     if (stride) {
13         // 使用给定步长设置大小
14         self->set_sizes_and_strides(size, *stride);
15         // 计算需要的存储字节数
16         storage_size = at::detail::computeStorageNbytes(size, *stride, itemsize,
17 storage_offset);
18     } else {
19         // 设置连续内存布局
20         self->set_sizes_contiguous(size);
21         // 计算连续内存所需字节数
22         storage_size = at::detail::computeStorageNbytesContiguous(size, itemsize,
23 storage_offset);
24     }
25     // 根据需要调整存储大小
26     maybe_resize_storage_cuda(self, storage_size);
27     return self;
28 }
```

3. 存储扩容检查: `maybe_resize_storage_cuda` 函数

代码块

```
1  static inline void maybe_resize_storage_cuda(TensorImpl* self, size_t
2  new_size_bytes) {
3      // 空张量不需要调整存储
4      if (self->numel() == 0) {
5          return;
6      }
```

```

7   const Storage &storage = self->unsafe_storage();
8   TORCH_CHECK(storage, "Tensor: invalid null storage");
9
10  // 仅当需要更大存储时才调整
11  if (new_size_bytes > storage.nbytes()) {
12      resize_bytes_cuda(storage.unsafeGetStorageImpl(), new_size_bytes);
13  }
14  }

```

4. 实际内存调整: `resize_bytes_cuda` 函数

代码块

```

1  void resize_bytes_cuda(StorageImpl* storage, size_t size_bytes) {
2      // 检查存储是否可调整大小
3      TORCH_CHECK(storage->resizable(), "Trying to resize storage that is not
4      resizable");
5      auto allocator = storage->allocator();
6      TORCH_CHECK(allocator != nullptr, "Trying to resize storage without an
7      allocator");
8
9      c10::Device device = storage->device();
10
11     // 特殊情况: 调整为0大小
12     if (size_bytes == 0) {
13         storage->set_data_ptr_noswap(at::DataPtr(nullptr, device));
14         storage->set_nbytes(0);
15         return;
16     }
17
18     // 设置正确的CUDA设备上下文
19     c10::cuda::CUDAGuard guard(device.index());
20
21     // 分配新内存
22     at::DataPtr data = allocator->allocate(size_bytes);
23
24     // 如果存在旧数据, 复制到新内存
25     if (storage->data_ptr()) {
26         at::globalContext().lazyInitCUDA();
27
28         // 异步复制数据 - 高效的CUDA内存管理
29         C10_CUDA_CHECK(
30             cudaMemcpyAsync(
31                 data.get(), // 目标: 新内存
32                 storage->data(), // 源: 旧内存
33                 std::min(storage->nbytes(), size_bytes), // 复制量: 较小的值
34                 cudaMemcpyDeviceToDevice, // 复制类型: 设备到设备

```

```

33         c10::cuda::getCurrentCUDAStream()); // 使用当前CUDA流
34     }
35
36     // 用新的数据指针替换旧的
37     storage->set_data_ptr_noswap(std::move(data));
38     storage->set_nbytes(size_bytes);
39 }

```

rsub.Tensor, rsub.Tensor_out, rsub.Scalar, rsub.Scalar_out

redispatch 到 sub

代码块

```

1  aten/src/ATen/native/BinaryOps.cpp
2
3  // rsub.Scalar -> rsub.Tensor
4  Tensor rsub(const Tensor& self, const Scalar& other, const Scalar& alpha) {
5      return native::rsub(self, wrapped_scalar_tensor(other), alpha);
6  }
7  // rsub.Tensor -> sub.Tensor
8  Tensor rsub(const Tensor& self, const Tensor& other, const Scalar& alpha) {
9      return at::sub(other, self, alpha); // redispatch!
10 }

```

floor_divide, floor_divide.out, floor_divide_.Tensor, floor_divide.Scalar, floor_divide.Scalar_out, floor_divide_.Scalar

参考信息

文档: https://pytorch.org/docs/stable/generated/torch.floor_divide.html

参考实现:

<https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/BinaryDivFloorKernel.cu>

floor_divide.Scalar, floor_divide.Scalar_out, floor_divide_.Scalar 在 build/aten/src/ATen/RegisterCompositeExplicitAutograd.cpp 中注册

```

build/aten/src/ATen/RegisterCompositeExplicitAutograd.cpp
1  m.impl("floor_divide.Scalar",
2  TORCH_FN(wrapper_CompositeExplicitAutograd_Scalar_floor_divide));
3  m.impl("floor_divide.Scalar_out",
4  TORCH_FN(wrapper_CompositeExplicitAutograd_Scalar_out_floor_divide_out));

```

```
5 m.impl("floor_divide_.Scalar",
6 TORCH_FN(wrapper_CompositeExplicitAutograd_Scalar_floor_divide_));
```

在 aten/src/ATen/native/BinaryOps.cpp 中 `redispatch` 到 `floor_divide`

代码块

```
1 Tensor floor_divide(const Tensor& self, const Scalar& other) {
2     return at::floor_divide(self, wrapped_scalar_tensor(other));
3 }
4
5 Tensor& floor_divide_(Tensor& self, const Scalar& other) {
6     return at::floor_divide_out(self, self, wrapped_scalar_tensor(other));
7 }
```

通过最终 `floor_divide` 注册到 `div_floor_stub`

代码块

```
1 import torch
2
3 a = torch.tensor([5, 6, 7, 8])
4 b = torch.tensor([2, 2, 3, 3])
5
6 # 张量与张量整除
7 # 使用函数
8 result = torch.floor_divide(a, b)
9 print(result) # tensor([2, 3, 2, 2])
10
11 # 使用运算符
12 result = a // b
13 print(result) # tensor([2, 3, 2, 2])
14
15 # 张量方法
16 result = a.floor_divide(b)
17 print(result) # tensor([2, 3, 2, 2])
18
19 a = torch.tensor([5, 6, 7, 8])
20
21 # 张量与标量整除
22 # 使用函数
23 result = torch.floor_divide(a, 2)
24 print(result) # tensor([2, 3, 3, 4])
25
26 # 使用运算符
27 result = a // 2
```

```

28 print(result) # tensor([2, 3, 3, 4])
29
30 # 标量在左边
31 result = torch.floor_divide(2, a)
32 print(result) # tensor([0, 0, 0, 0]) # 因为 2/5, 2/6 等都小于1, 向下取整为0
33
34 # 广播支持
35 a = torch.tensor([[5, 6], [7, 8]])
36 b = torch.tensor([2, 3])
37
38 result = a // b
39 print(result)
40 # tensor([[2, 2],
41 #         [3, 2]])
42
43 #当操作数中有浮点数时, 结果为浮点数:
44 a = torch.tensor([5.0, 6.0, 7.0])
45 b = torch.tensor([2.0, 2.0, 2.0])
46
47 # 普通除法
48 div_result = a / b
49 print(div_result) # tensor([2.5000, 3.0000, 3.5000])
50
51 # 整除
52 floor_div_result = a // b
53 print(floor_div_result) # tensor([2., 3., 3.])
54
55 # 等价于
56 floor_div_alternative = torch.floor(a / b)
57 print(floor_div_alternative) # tensor([2., 3., 3.])

```

Function 接口

代码块

```

1  /**
2   * @brief Performs element-wise floor division (a // b)
3   *
4   * @param dev Device handle for CUDA/CPU device management
5   * @param out Output tensor to store the result
6   * @param self Dividend tensor (left operand)
7   * @param other Divisor tensor or scalar (right operand)
8   * @param out_sizes Output tensor sizes
9   * @param out_strides Output tensor strides (can be non-contiguous)
10  * @param self_sizes Dividend tensor sizes
11  * @param self_strides Dividend tensor strides (can be non-contiguous)
12  * @param other_sizes Divisor tensor sizes (nullptr if scalar)

```

```

13 * @param other_strides Divisor tensor strides (can be non-contiguous)
14 * @param out_ndims Number of output tensor dimensions
15 * @param self_ndims Number of dividend tensor dimensions
16 * @param other_ndims Number of divisor tensor dimensions
17 * @param self_dtype Data type of dividend tensor
18 * @param other_dtype Data type of divisor tensor
19 * @param is_other_scalar Whether other is a scalar
20 *
21 * @note Memory Layout:
22 * - Fully supports non-contiguous tensors
23 * - Handles arbitrary strides efficiently
24 * - Optimized for both contiguous and non-contiguous cases
25 *
26 * @note Behavior:
27 * - Performs division followed by floor operation
28 * - Handles type promotion according to PyTorch rules
29 * - Properly handles division by zero according to IEEE standards
30 * - For integer inputs, returns integer floor division (equivalent to
Python //)
31 */
32 void floor_divide_kernel(
33     SiDeviceHandle *dev, // Device handle
34     void *out, // Output tensor data pointer
35     const void *self, // Dividend tensor data pointer
36     const void *other, // Divisor tensor/scalar data pointer
37     const int64_t* out_sizes, // Output sizes
38     const int64_t* out_strides, // Output strides
39     const int64_t* self_sizes, // Dividend tensor sizes
40     const int64_t* self_strides, // Dividend tensor strides
41     const int64_t* other_sizes, // Divisor tensor sizes (nullptr if scalar)
42     const int64_t* other_strides, // Divisor tensor strides (nullptr if
scalar)
43     int out_ndims, // Number of output dimensions
44     int self_ndims, // Number of dividend dimensions
45     int other_ndims, // Number of divisor dimensions
46     int self_dtype, // Data type of dividend
47     int other_dtype, // Data type of divisor
48     bool is_other_scalar // Whether divisor is a scalar
49 );
50
51 支持的数据类型
52 整数类型: uint8, int8, int16, int32, int64
53 浮点类型: float16 (half), bfloat16, float32, float64

```

(1) 数制支持上，初版需要支持int32类型，如果输入数据类型为int64，由框架进行int64->int32转换

(2) 除数（输入2） 需要支持tensor与scalar两种形式

(3) 输出与被除数（输入1） 的shape一致

(4) 如果除数为tensor，那么除数与被除数shape应相同（如不同，则由框架处理进行broadcast），做elementwise除法

(5) 对齐限制：全流程使用RVV Core，无对齐限制

(6) 输入输出参数，均由框架进行预处理（reshape为一维）

内部接口调用形式：

代码块

```
1  /// @brief Tensor floor division
2  /// @tparam T Datatype support [int32_t]
3  /// @param[in] d_in0 Dividend tensor with shape [dim0_size]
4  /// @param[in] d_in1 Divisor tensor with shape [dim0_size] or scalar if
   is_d_in1_scalar is true
5  /// @param[out] d_out Output tensor with shape [dim0_size]
6  /// @param[in] dim0_size Tensor shape param
7  /// @param[in] is_d_in1_scalar Indicates d_in1 is scalar or tensor
8  template <typename T>
9  void floor_div(void *d_in0, void *d_in1, void *d_out, int dim0_size, bool
   is_d_in1_scalar);
```

ISSUE

非 contiguous 的是否能支持（如果支持需要修改 Function 接口）

pytorch 中使用了 TensorIteratorBase 传递 self, other, out, 需考虑非 contiguous 实现

考虑到 broadcast 的情况，out 和 self 的 size/stride 存在不一致的情况，需要传入

Memory 是否需要 Align, Align 到多少

repeat_interleave.self_Tensor, repeat_interleave.self_int

参考信息

参考文档：https://pytorch.org/docs/stable/generated/torch.repeat_interleave.html

repeat_interleave 是由以下几个 op 组合完成：

is_conj, conj, is_neg, _neg_view, flatten, reshape, expand_symint, index_select, unsqueeze, clone

代码块 所有元素重复相同次数

```
2 >>> x = torch.tensor([1, 2, 3])
3 >>> repeated = torch.repeat_interleave(x, repeats=3)
4 >>> repeated
5 tensor([1, 1, 1, 2, 2, 2, 3, 3, 3])
6
7 # 指定维度重复
8 >>> x = torch.tensor([[1, 2], [3, 4]])
9 >>> repeated_rows = torch.repeat_interleave(x, repeats=2, dim=0)
10 >>> repeated_rows
11 tensor([[1, 2],
12         [1, 2],
13         [3, 4],
14         [3, 4]])
15 >>> repeated_cols = torch.repeat_interleave(x, repeats=3, dim=1)
16 >>> repeated_cols
17 tensor([[1, 1, 1, 2, 2, 2],
18         [3, 3, 3, 4, 4, 4]])
19
20 # 每个元素重复不同次数
21 >>> x = torch.tensor([1, 2, 3])
22 >>> repeats = torch.tensor([1, 2, 3])
23 >>> repeated = torch.repeat_interleave(x, repeats)
24 >>> repeated
25 tensor([1, 2, 2, 3, 3, 3])
26
27 # 在多维张量上使用不同的重复次数
28 >>> x = torch.tensor([[1, 2, 3], [4, 5, 6]])
29 >>> row_repeats = torch.tensor([1, 3]) # 第一行重复1次, 第二行重复3次
30 >>> repeated_rows = torch.repeat_interleave(x, row_repeats, dim=0)
31 >>> repeated_rows
32 tensor([[1, 2, 3],
33         [4, 5, 6],
34         [4, 5, 6],
35         [4, 5, 6]])
```

代码块

```
1 // repeat_interleave.self_Tensor
2 Tensor repeat_interleave_symint(
3     const Tensor& self,
4     const Tensor& repeats,
5     std::optional<int64_t> dim,
6     std::optional<SymInt> output_size) {
7     Tensor input = self;
8 }
```

```

9 // Store conj and neg bits
10 const auto conj = input.is_conj();
11 if (conj) {
12     input = input.conj();
13 }
14 const auto neg = input.is_neg();
15 if (neg) {
16     input = input._neg_view();
17 }
18
19 if (!dim) {
20     input = input.flatten();
21     dim = 0;
22 }
23
24 Tensor repeats_ = repeats;
25 if (repeats.dim() == 0 || (repeats.dim() == 1 && repeats.sym_size(0) == 1)) {
26     repeats_ =
repeats.reshape({1}).expand_symint({input.sym_size(dim.value())});
27 } else if (repeats.dim() == 1) {
28     TORCH_CHECK(
29         repeats.sym_size(0) == input.sym_size(dim.value()),
30         "repeats must have the same size as input along dim, but got
repeats.size(0) = ",
31         repeats.sym_size(0), " and input.size(", dim.value(), ") = ",
input.sym_size(dim.value())
32     );
33 } else {
34     AT_ERROR("repeats must be 0-dim or 1-dim tensor");
35 }
36
37 auto ret = input.index_select(
38     dim.value(), at::repeat_interleave_symint(repeats_,
std::move(output_size)));
39 // Restore conj and neg bits
40 if (conj) {
41     ret = ret.conj();
42 }
43 if (neg) {
44     ret = ret._neg_view();
45 }
46 return ret;
47 }
48
49 // repeat_interleave.self_int
50 Tensor repeat_interleave_symint(
51     const Tensor& self,

```

```

52     c10::SymInt repeats,
53     std::optional<int64_t> dim_opt,
54     std::optional<SymInt> output_size) {
55     Tensor input = dim_opt ? self : self.flatten();
56     int64_t dim = c10::maybe_wrap_dim(dim_opt.value_or(0), self.dim());
57     TORCH_CHECK(repeats >= 0, "Repeats must be non-negative");
58
59     input = input.unsqueeze(dim + 1);
60     auto expand_shape = input.sym_sizes().vec();
61     expand_shape[dim + 1] = repeats;
62     input = input.expand_symint(expand_shape);
63
64     // This argument doesn't really make sense for the scalar overload, but
65     // exists
66     // for consistency with the tensor overload
67     if (output_size) {
68         auto calculated_size = (repeats * expand_shape[dim]).guard_int(__FILE__,
69         __LINE__);
70         TORCH_CHECK(*output_size == calculated_size, "repeat_interleave: Invalid
71         output_size, expected ",
72         calculated_size, " but got ", *output_size);
73     }
74     return input.clone(at::MemoryFormat::Contiguous).flatten(dim, dim + 1);
75 }

```

addmv, addmv.out, addmv_

参考信息

pytorch官方文档: <https://pytorch.org/docs/stable/generated/torch.addmv.html>

Addmv 通过cuda cublas的gemv实现, 参考cublas文档

<https://docs.nvidia.com/cuda/cublas/#cublas-t-gemv>

调用位置:

<https://github.com/pytorch/pytorch/blob/164d2c887b45ec2673b9bbc08e457bfab0a1eb3c/aten/src/ATen/native/cuda/Blas.cpp#L839>

代码块

```

1  >>> input = torch.tensor([1, 2, 3], dtype=torch.float)
2  >>> mat = torch.tensor([[2, 3, 4], [5, 6, 7], [8, 9, 10]], dtype=torch.float)
3  >>> vec = torch.tensor([1, 2, 3], dtype=torch.float)

```

```

4 >>> result = torch.addmv(input, mat, vec)
5 >>> result
6 tensor([21., 40., 59.])
7 >>> result = torch.addmv(input, mat, vec, beta=2, alpha=3)
8 >>> result
9 tensor([ 62., 118., 174.])
10 >>> manual = input + mat @ vec
11 >>> manual
12 tensor([21., 40., 59.])
13 >>> manual2 = 2 * input + 3 * (mat @ vec)
14 >>> manual2
15 tensor([ 62., 118., 174.])

```

Function 接口

代码块

```

1  /**
2   * @brief Performs matrix-vector multiplication with scaling and addition
   (addmv operation)
3   *
4   * @param dev Device handle for CUDA device management
5   * @param out Output tensor to store the result
6   * @param self Input vector to be scaled by beta
7   * @param mat Input matrix for matrix-vector multiplication
8   * @param vec Input vector for matrix-vector multiplication
9   * @param beta Scalar multiplier for self
10  * @param alpha Scalar multiplier for the matrix-vector product
11  * @param out_sizes Output tensor sizes
12  * @param out_strides Output tensor strides (can be non-contiguous)
13  * @param self_sizes Input vector sizes
14  * @param self_strides Input vector strides (can be non-contiguous)
15  * @param mat_sizes Matrix tensor sizes
16  * @param mat_strides Matrix tensor strides (can be non-contiguous)
17  * @param vec_sizes Vector tensor sizes
18  * @param vec_strides Vector tensor strides (can be non-contiguous)
19  * @param out_ndims Number of output tensor dimensions
20  * @param self_ndims Number of input vector dimensions
21  * @param mat_ndims Number of matrix tensor dimensions
22  * @param vec_ndims Number of vector tensor dimensions
23  * @param self_dtype Data type of input vector
24  * @param mat_dtype Data type of matrix tensor
25  * @param vec_dtype Data type of vector tensor
26  *
27  * @note Memory Layout:
28  * - Fully supports non-contiguous tensors
29  * - Optimized for different matrix layouts:

```

```

30 *      - Row-major: Uses direct gemv with 'n' mode
31 *      - Column-major: Uses transposed gemv with 't' mode
32 *      - Handles arbitrary strides efficiently
33 *
34 * @note Behavior:
35 *      - Computes: out = beta * self + alpha * (mat @ vec)
36 *      - If mat is empty, either zeroes out or computes beta * self
37 *      - Handles type promotion according to PyTorch rules
38 *      - For vectors with stride 0, creates contiguous copies
39 */
40 template <typename T>
41 void addmv_cuda_kernel(
42     SiDeviceHandle *dev,           // Device handle
43     void *out,                     // Output tensor data pointer
44     const void *self,              // Input vector data pointer
45     const void *mat,               // Matrix tensor data pointer
46     const void *vec,               // Vector tensor data pointer
47     const T& beta,                 // Scalar multiplier for self
48     const T& alpha,                // Scalar multiplier for mat @ vec
49     const int64_t* out_sizes,      // Output sizes
50     const int64_t* out_strides,    // Output strides
51     const int64_t* self_sizes,     // Input vector sizes
52     const int64_t* self_strides,   // Input vector strides
53     const int64_t* mat_sizes,      // Matrix tensor sizes
54     const int64_t* mat_strides,    // Matrix tensor strides
55     const int64_t* vec_sizes,      // Vector tensor sizes
56     const int64_t* vec_strides,    // Vector tensor strides
57     int out_ndims,                 // Number of output dimensions
58     int self_ndims,                // Number of input vector dimensions
59     int mat_ndims,                 // Number of matrix dimensions
60     int vec_ndims,                 // Number of vector dimensions
61     int self_dtype,                // Data type of input vector
62     int mat_dtype,                 // Data type of matrix
63     int vec_dtype                   // Data type of vector
64 );

```

支持的数据类型:

浮点类型: `float`, `double`

复数类型: `complex<float>`, `complex<double>`

特殊类型: `Half`, `BFloat16`

ISSUE

cuda的实现支持非contiguous输入，也支持行优先/列优先模式

conj_physical.out

参考信息

pytorch官方文档: https://pytorch.org/docs/stable/generated/torch.conj_physical.html

参考实现:

<https://github.com/pytorch/pytorch/blob/836955bdbdeb299e6937065299564fb44ec422c2/aten/src/ATen/native/cuda/UnaryComplexKernels.cu#L66>

代码块

```
1 conj_physical 是 PyTorch 中用于处理复数张量的函数，它返回一个新的张量，其中每个复数元素
  都被替换为其共轭 (conjugate)
2 对于复数类型 (kComplexFloat、kComplexDouble、kComplexHalf)，执行真正的共轭操作
3 对于所有其他类型 (整数、浮点数、布尔值)，执行直接复制操作 (no-op)
4
5 >>> x = torch.tensor([1+2j, 3-4j, 5+0j])
6 >>> y = torch.conj_physical(x)
7 >>> x
8 tensor([1.+2.j, 3.-4.j, 5.+0.j])
9 >>> y
10 tensor([1.-2.j, 3.+4.j, 5.-0.j])
```

Function 接口

代码块

```
1 /**
2  * @brief Computes the complex conjugate of each element in the input tensor
3  *
4  * @param dev Device handle for CUDA device management
5  * @param out Output tensor to store the result
6  * @param self Input tensor
7  * @param out_sizes Output tensor sizes
8  * @param out_strides Output tensor strides (can be non-contiguous)
9  * @param self_sizes Input tensor sizes
10 * @param self_strides Input tensor strides (can be non-contiguous)
11 * @param out_ndims Number of output tensor dimensions
12 * @param self_ndims Number of input tensor dimensions
13 * @param out_dtype Data type of output tensor
14 * @param self_dtype Data type of input tensor
15 *
16 * @note Memory Layout:
17 * - Fully supports non-contiguous tensors
18 * - Handles arbitrary strides efficiently
19 * - Optimized for both contiguous and non-contiguous cases
```

```

20 *
21 * @note Behavior:
22 *   - For complex types: returns the complex conjugate ( $a-bi$  for  $a+bi$ )
23 *   - For non-complex types: performs a direct copy (no-op)
24 *   - Preserves NaN values
25 *   - Supports complex64, complex128, and complex half types
26 *   - For boolean and real types, performs a direct copy
27 *
28 * @note Performance:
29 *   - Uses JIT compilation when available for optimal performance
30 *   - Falls back to explicit kernel implementation when JIT is unavailable
31 *   - Optimized for parallel execution on CUDA devices
32 */

```

```

33 template <typename T>
34 void conj_kernel_cuda(
35     SiDeviceHandle *dev,           // Device handle
36     void *out,                    // Output tensor data pointer
37     const void *self,              // Input tensor data pointer
38     const int64_t* out_sizes,      // Output sizes
39     const int64_t* out_strides,    // Output strides
40     const int64_t* self_sizes,     // Input tensor sizes
41     const int64_t* self_strides,   // Input tensor strides
42     int out_ndims,                 // Number of output dimensions
43     int self_ndims,                // Number of input dimensions
44     int out_dtype,                 // Data type of output
45     int self_dtype                 // Data type of input
46 );

```

支持的数据类型

```

49 kByte (uint8_t)
50 kChar (int8_t)
51 kShort (int16_t)
52 kInt (int32_t)
53 kLong (int64_t)

```

所有浮点类型:

```

55 kFloat (float)
56 kDouble (double)
57 kHalf (at::Half)
58 kBFloat16 (at::BFloat16)

```

布尔类型:

```

60 kBool (bool)

```

复数类型

```

62 kComplexFloat (c10::complex<float>)
63 32位复数类型, 由两个单精度浮点数组成
64 kComplexDouble (c10::complex<double>)
65 64位复数类型, 由两个双精度浮点数组成
66 kComplexHalf (c10::complex<at::Half>)

```

67 16位复数类型，由两个半精度浮点数组成

68

ISSUE

cuda中通过TensorIteratorBase来传递输入和输出，支持非contiguous

dot

参考信息

pytorch资料: <https://pytorch.org/docs/stable/generated/torch.dot.html>

cuda的实现使用的是cublas, 参考资料: <https://docs.nvidia.com/cuda/cublas/#cublas-t-dot>

Computes the dot product of two 1D tensors.

Function 接口

代码块

```
1
2 /**
3  * @brief Computes the dot product of two 1-D vectors
4  *
5  * @param dev Device handle for CUDA device management
6  * @param out Output tensor to store the result (scalar)
7  * @param self First input vector
8  * @param other Second input vector
9  * @param self_sizes First vector sizes
10 * @param other_sizes Second vector sizes
11 * @param self_strides First vector strides (can be non-contiguous)
12 * @param other_strides Second vector strides (can be non-contiguous)
13 * @param self_ndims Number of dimensions in first vector (should be 1)
14 * @param other_ndims Number of dimensions in second vector (should be 1)
15 * @param self_dtype Data type of first vector
16 * @param other_dtype Data type of second vector
17 *
18 * @note Memory Layout:
19 * - Supports non-contiguous vectors with arbitrary strides
20 * - For vectors of size 1, strides are treated as 1 regardless of
   actual value
21 * - Optimized for both contiguous and non-contiguous cases
22 *
23 * @note Behavior:
24 * - Computes  $\text{sum}(\text{self}[i] * \text{other}[i])$  for all  $i$ 
```

```

25 * - For complex types, performs standard complex multiplication (not
conjugated)
26 * - Returns a scalar tensor of the same type as inputs
27 * - If either input is a zero tensor, returns an efficient zero tensor
28 * - For complex inputs with conjugation flags, special handling is
applied:
29 * - If self is conjugated but other is not: equivalent to
vdot(self.conj(), other)
30 * - If other is conjugated but self is not: equivalent to
vdot(other.conj(), self)
31 * - If both are conjugated: equivalent to dot(self.conj(),
other.conj()).conj()
32 *
33 * @note Requirements:
34 * - Both vectors must have the same number of elements
35 * - Both vectors must have the same data type
36 * - Number of elements must be  $\leq$  INT_MAX
37 * - Strides must be  $\leq$  INT_MAX
38 *
39 * @note Performance:
40 * - Uses cuBLAS dot/dotc functions for optimal performance
41 * - Sets cuBLAS pointer mode to DEVICE for direct output to tensor
42 */
43 template <typename T>
44 void dot_cuda_kernel(
45     SiDeviceHandle *dev, // Device handle
46     void *out, // Output tensor data pointer (scalar)
47     const void *self, // First vector data pointer
48     const void *other, // Second vector data pointer
49     const int64_t self_sizes, // First vector sizes
50     const int64_t other_sizes, // Second vector sizes
51     const int64_t self_strides, // First vector strides
52     const int64_t other_strides, // Second vector strides
53     int self_dtype, // Data type of first vector
54     int other_dtype, // Data type of second vector
55     bool self_is_conj, // Whether first vector is conjugated
56     bool other_is_conj // Whether second vector is conjugated
57 );

```

支持的数据类型

浮点类型:

float (单精度浮点数)

double (双精度浮点数)

半精度类型:

Half (16位半精度浮点数)

BFloat16 (16位脑浮点数)

复数类型:

```
67 complex<float> (复数, 实部和虚部都是单精度浮点数)
68 complex<double> (复数, 实部和虚部都是双精度浮点数)
```

ISSUE

支持非连续向量, 需要明确处理步长

replication_pad1d, replication_pad1d.out

参考信息

pytorch文档: <https://pytorch.org/docs/stable/generated/torch.nn.ReplicationPad1d.html>

参考实现

<https://github.com/pytorch/pytorch/blob/2e7c9d33e7f933ac3b723cb3bb05b9c88432c25c/aten/src/ATen/native/cuda/ReplicationPadding.cu#L419>

代码块

```
1 import torch
2 import torch.nn as nn
3
4 # 创建一个输入张量
5 input = torch.tensor([[[[1, 2, 3, 4, 5]]], dtype=torch.float)
6
7 # 创建一个 ReplicationPad1d 层
8 pad_layer = nn.ReplicationPad1d(2) #左右均pad 2 个
9
10 # 应用填充
11 output = pad_layer(input)
12 print(output)
13 # 输出: tensor([[[[1, 1, 1, 2, 3, 4, 5, 5, 5]]]])
14
15 pad_layer = nn.ReplicationPad1d(1,2) # 左侧pad 1 右侧 pad 2
16 # 应用填充
17 output = pad_layer(input)
18 print(output)
19 # 输出: tensor([[[[1., 1., 2., 3., 4., 5., 5., 5.]]]])
20
21
22 m = nn.ReplicationPad1d(2)
23 input = torch.arange(8, dtype=torch.float).reshape(1, 2, 4)
24 # tensor([[[[0., 1., 2., 3.], [4., 5., 6., 7.]]]])
25 m(input)
26 # tensor([[[[0., 0., 0., 1., 2., 3., 3., 3.],
27 # [4., 4., 4., 5., 6., 7., 7., 7.]]]])
28
29 # using different paddings for different sides
```

```

30 m = nn.ReplicationPad1d((3, 1))
31 m(input)
32 # tensor([[[[0., 0., 0., 0., 1., 2., 3.],
33 #          [4., 4., 4., 4., 5., 6., 7., 7.]]]])

```

Function 接口

代码块

```

1  /**
2   * @brief Performs 1D replication padding on input tensor
3   *
4   * @param dev Device handle for CUDA device management
5   * @param out Output tensor to store the result
6   * @param input Input tensor to be padded
7   * @param padL Number of padding elements on the left side
8   * @param padR Number of padding elements on the right side
9   * @param out_sizes Output tensor sizes
10  * @param out_strides Output tensor strides (can be non-contiguous)
11  * @param input_sizes Input tensor sizes
12  * @param input_strides Input tensor strides (can be non-contiguous)
13  * @param out_ndims Number of dimensions in output tensor
14  * @param input_ndims Number of dimensions in input tensor
15  * @param input_dtype Data type of input tensor
16  *
17  * @note Memory Layout:
18  *   - Supports non-contiguous tensors with arbitrary strides
19  *   - Handles both 2D (C, W) and 3D (N, C, W) input formats
20  *   - For 2D inputs, the operation is applied to each channel
    independently
21  *
22  * @note Behavior:
23  *   - Pads by replicating the edge values of the input tensor
24  *   - For example, with input [1, 2, 3] and padding (2, 1):
25  *     Output will be [1, 1, 1, 2, 3, 3]
26  *   - Supports all data types including floating point, integer, and
    complex types
27  *
28  * @note Performance:
29  *   - Uses CUDA kernel with block size optimization for different tensor
    sizes
30  *   - Processes tensors in chunks to handle large batch sizes efficiently
31  *   - Automatically handles tensor dimensions > 65535 by splitting into
    multiple kernel launches
32  */
33 template <typename T>
34 void replication_pad1d_out_kernel(

```

```

35     SiDeviceHandle *dev,           // Device handle
36     void *out,                    // Output tensor data pointer
37     const void *input,           // Input tensor data pointer
38     int64_t padL,                 // Left padding size
39     int64_t padR,                 // Right padding size
40     const int64_t* out_sizes,     // Output sizes
41     const int64_t* out_strides,   // Output strides
42     const int64_t* input_sizes,   // Input tensor sizes
43     const int64_t* input_strides, // Input tensor strides
44     int out_ndims,                // Number of output dimensions
45     int input_ndims,              // Number of input dimensions
46     int input_dtype               // Data type of input
47 );

```

新接口:

```

50 template <typename T>
51 void replication_pad_out_kernel(
52     const skernel::Tensor<T>& self,
53     skernel::Tensor<T>& out,
54     const void* padding,
55     int padding_arr_size,
56 );

```

支持数据类型:

```

59 整数类型: uint8_t, int8_t, int16_t, int32_t, int64_t
60 浮点类型: float, double
61 半精度类型: Half (kHalf), BFloat16 (kBFloat16)
62 复数类型: complex<float>, complex<double>

```

ISSUE

- 支持非contiguous, 同时支持2D,3D输入
- cuda的实现通过分块处理超过 CUDA 网格限制 (65535) 的大型张量

replication_pad2d, replication_pad2d.out

参考信息

pytorch文档: <https://pytorch.org/docs/stable/generated/torch.nn.ReplicationPad2d.html>

参考实现:

<https://github.com/pytorch/pytorch/blob/2e7c9d33e7f933ac3b723cb3bb05b9c88432c25c/aten/src/ATen/native/cuda/ReplicationPadding.cu#L534>

代码块

```
1 >>> m = nn.ReplicationPad2d(2)
```

```

2 >>> input = torch.arange(9, dtype=torch.float).reshape(1, 1, 3, 3)
3 >>> input
4 tensor([[[[0., 1., 2.],
5          [3., 4., 5.],
6          [6., 7., 8.]]]]])
7 >>> m(input)
8 tensor([[[[0., 0., 0., 1., 2., 2., 2.],
9          [0., 0., 0., 1., 2., 2., 2.],
10         [0., 0., 0., 1., 2., 2., 2.],
11         [3., 3., 3., 4., 5., 5., 5.],
12         [6., 6., 6., 7., 8., 8., 8.],
13         [6., 6., 6., 7., 8., 8., 8.],
14         [6., 6., 6., 7., 8., 8., 8.]]]]])
15
16 >>> m = nn.ReplicationPad2d((1, 1, 2, 0)) # (padding_left, padding_right,
padding_top, padding_bottom)
17 >>> m(input)
18 tensor([[[[0., 0., 1., 2., 2.],
19          [0., 0., 1., 2., 2.],
20          [0., 0., 1., 2., 2.],
21          [3., 3., 4., 5., 5.],
22          [6., 6., 7., 8., 8.]]]]])
23 >>>

```

Function 接口

代码块

```

1  /**
2   * @brief Performs 2D replication padding on input tensor
3   *
4   * @param dev Device handle for CUDA device management
5   * @param out Output tensor to store the result
6   * @param input Input tensor to be padded
7   * @param padL Number of padding elements on the left side
8   * @param padR Number of padding elements on the right side
9   * @param padT Number of padding elements on the top side
10  * @param padB Number of padding elements on the bottom side
11  * @param out_sizes Output tensor sizes
12  * @param out_strides Output tensor strides (can be non-contiguous)
13  * @param input_sizes Input tensor sizes
14  * @param input_strides Input tensor strides (can be non-contiguous)
15  * @param out_ndims Number of dimensions in output tensor
16  * @param input_ndims Number of dimensions in input tensor
17  * @param input_dtype Data type of input tensor
18  *
19  * @note Memory Layout:

```

```

20 * - Supports non-contiguous tensors with arbitrary strides
21 * - Handles both 3D (C, H, W) and 4D (N, C, H, W) input formats
22 * - For 3D inputs, the operation is applied to each channel
independently
23 *
24 * @note Behavior:
25 * - Pads by replicating the edge values of the input tensor
26 * - For example, with a 3x3 input and padding (1,1,1,1), the corner
values
27 * are replicated diagonally and edge values are replicated outward
28 * - Supports all data types including floating point, integer, and
complex types
29 *
30 * @note Performance:
31 * - Uses CUDA kernel with block size optimization for different tensor
sizes
32 * - Processes tensors in chunks to handle large batch sizes efficiently
33 * - Automatically handles tensor dimensions > 65535 by splitting into
multiple kernel launches
34 * - Optimized for both small and large padding sizes
35 *
36 * @note Requirements:
37 * - Input tensor must fit into 32-bit index math
38 * - Output tensor must have correct size: [N, C, H+padT+padB,
W+padL+padR]
39 */
40 和replication_pad1d接口一样
41 template <typename T>
42 void replication_pad_out_kernel(
43     const skernel::Tensor<T>& self,
44     skernel::Tensor<T>& out,
45     const void* padding,
46     int padding_arr_size,
47 );

```

ISSUE

- 支持非contiguous，同时支持3D,4D输入
- cuda的实现通过分块处理超过 CUDA 网格限制（65535）的大型张量
- cuda的实现没有接收padR，padB。而是通过padL和padT计算出原始数据的起始位置进行填充，然后再根据output的大小，在padding部分直接复制最近的边界值

replication_pad3d, replication_pad3d.out

参考信息

参考信息

pytorch文档: <https://pytorch.org/docs/stable/generated/torch.nn.ReplicationPad3d.html>

参考实现:

<https://github.com/pytorch/pytorch/blob/2e7c9d33e7f933ac3b723cb3bb05b9c88432c25c/aten/src/ATen/native/cuda/ReplicationPadding.cu#L602>

代码块

```
1 >>> m = nn.ReplicationPad3d(1)
2 >>> input = torch.randn(2, 2, 2, 2, 2)
3 >>> output = m(input)
4 >>> input
5 tensor([[[[[[ 0.1190, -2.7971],
6             [ 0.3847, -0.9118]],
7
8             [[ 0.7726, -0.1665],
9             [-0.2704, -0.9113]]],
10
11            [[[-0.1435, -2.3199],
12             [-0.3886,  0.3390]],
13
14            [[-1.7399,  0.3719],
15             [ 1.6875,  0.3726]]]]],
16
17           [[[[[ 0.1996, -0.6908],
18             [-0.3071, -1.2073]],
19
20            [[ 1.0142, -0.7493],
21             [ 1.9937, -0.4326]]]]],
22
23           [[[[[ 0.4510, -0.8870],
24             [ 0.7344, -1.1570]],
25
26            [[-0.7517,  0.7564],
27             [-0.1884, -0.0280]]]]]])
28
29 >>> output
30 tensor([[[[[[ 0.1190,  0.1190, -2.7971, -2.7971],
31             [ 0.1190,  0.1190, -2.7971, -2.7971],
32             [ 0.3847,  0.3847, -0.9118, -0.9118],
33             [ 0.3847,  0.3847, -0.9118, -0.9118]],
34
35            [[ 0.1190,  0.1190, -2.7971, -2.7971],
36             [ 0.1190,  0.1190, -2.7971, -2.7971],
37             [ 0.3847,  0.3847, -0.9118, -0.9118],
38             [ 0.3847,  0.3847, -0.9118, -0.9118]]]]]])
```

39 [0.1190, **0.1190**, -2.7971, -2.7971],
40 [0.3847, **0.3847**, -0.9118, -0.9118],
41 [0.3847, 0.3847, -0.9118, -0.9118]],
42
43 [[0.7726, 0.7726, -0.1665, -0.1665],
44 [0.7726, **0.7726**, -0.1665, -0.1665],
45 [-0.2704, -0.2704, -0.9113, -0.9113],
46 [-0.2704, -0.2704, -0.9113, -0.9113]],
47
48 [[0.7726, 0.7726, -0.1665, -0.1665],
49 [0.7726, **0.7726**, -0.1665, -0.1665],
50 [-0.2704, -0.2704, -0.9113, -0.9113],
51 [-0.2704, -0.2704, -0.9113, -0.9113]]],
52
53
54 [[[-0.1435, -0.1435, -2.3199, -2.3199],
55 [-0.1435, -0.1435, -2.3199, -2.3199],
56 [-0.3886, -0.3886, 0.3390, 0.3390],
57 [-0.3886, -0.3886, 0.3390, 0.3390]],
58
59 [[-0.1435, -0.1435, -2.3199, -2.3199],
60 [-0.1435, -0.1435, -2.3199, -2.3199],
61 [-0.3886, -0.3886, 0.3390, 0.3390],
62 [-0.3886, -0.3886, 0.3390, 0.3390]],
63
64 [[-1.7399, -1.7399, 0.3719, 0.3719],
65 [-1.7399, -1.7399, 0.3719, 0.3719],
66 [1.6875, 1.6875, 0.3726, 0.3726],
67 [1.6875, 1.6875, 0.3726, 0.3726]],
68
69 [[-1.7399, -1.7399, 0.3719, 0.3719],
70 [-1.7399, -1.7399, 0.3719, 0.3719],
71 [1.6875, 1.6875, 0.3726, 0.3726],
72 [1.6875, 1.6875, 0.3726, 0.3726]]]],
73
74
75
76 [[[[0.1996, 0.1996, -0.6908, -0.6908],
77 [0.1996, 0.1996, -0.6908, -0.6908],
78 [-0.3071, -0.3071, -1.2073, -1.2073],
79 [-0.3071, -0.3071, -1.2073, -1.2073]],
80
81 [[0.1996, 0.1996, -0.6908, -0.6908],
82 [0.1996, 0.1996, -0.6908, -0.6908],
83 [-0.3071, -0.3071, -1.2073, -1.2073],
84 [-0.3071, -0.3071, -1.2073, -1.2073]],
85

```

86     [[ 1.0142, 1.0142, -0.7493, -0.7493],
87     [ 1.0142, 1.0142, -0.7493, -0.7493],
88     [ 1.9937, 1.9937, -0.4326, -0.4326],
89     [ 1.9937, 1.9937, -0.4326, -0.4326]],
90
91     [[ 1.0142, 1.0142, -0.7493, -0.7493],
92     [ 1.0142, 1.0142, -0.7493, -0.7493],
93     [ 1.9937, 1.9937, -0.4326, -0.4326],
94     [ 1.9937, 1.9937, -0.4326, -0.4326]]],
95
96
97     [[[ 0.4510, 0.4510, -0.8870, -0.8870],
98     [ 0.4510, 0.4510, -0.8870, -0.8870],
99     [ 0.7344, 0.7344, -1.1570, -1.1570],
100    [ 0.7344, 0.7344, -1.1570, -1.1570]],
101
102    [[ 0.4510, 0.4510, -0.8870, -0.8870],
103    [ 0.4510, 0.4510, -0.8870, -0.8870],
104    [ 0.7344, 0.7344, -1.1570, -1.1570],
105    [ 0.7344, 0.7344, -1.1570, -1.1570]],
106
107    [[-0.7517, -0.7517, 0.7564, 0.7564],
108    [-0.7517, -0.7517, 0.7564, 0.7564],
109    [-0.1884, -0.1884, -0.0280, -0.0280],
110    [-0.1884, -0.1884, -0.0280, -0.0280]],
111
112    [[-0.7517, -0.7517, 0.7564, 0.7564],
113    [-0.7517, -0.7517, 0.7564, 0.7564],
114    [-0.1884, -0.1884, -0.0280, -0.0280],
115    [-0.1884, -0.1884, -0.0280, -0.0280]]]]])
116    >>>

```

Function 接口

代码块

```

1  /**
2   * @brief Performs 3D replication padding on input tensor
3   *
4   * @param dev          Device handle for CUDA device management
5   * @param out          Output tensor to store the result
6   * @param input        Input tensor to be padded
7   * @param pfront       Number of padding elements on the front side
8   * @param pback        Number of padding elements on the back side
9   * @param ptop         Number of padding elements on the top side
10  * @param pbottom      Number of padding elements on the bottom side
11  * @param pleft        Number of padding elements on the left side

```

```

12 * @param pright      Number of padding elements on the right side
13 * @param out_sizes   Output tensor sizes
14 * @param out_strides Output tensor strides (can be non-contiguous)
15 * @param input_sizes Input tensor sizes
16 * @param input_strides Input tensor strides (can be non-contiguous)
17 * @param out_ndims   Number of dimensions in output tensor
18 * @param input_ndims Number of dimensions in input tensor
19 * @param input_dtype Data type of input tensor
20 *
21 * @note Memory Layout:
22 *     - Supports non-contiguous tensors with arbitrary strides
23 *     - Handles both 4D (C, D, H, W) and 5D (N, C, D, H, W) input formats
24 *     - For 4D inputs, the operation is applied to each channel
independently
25 *
26 * @note Behavior:
27 *     - Pads by replicating the edge values of the input tensor in all
three dimensions
28 *     - Edge values are replicated outward, and corner values are
replicated diagonally
29 *     - Supports all data types including floating point, integer, and
complex types
30 *
31 * @note Performance:
32 *     - Uses CUDA kernel with block size optimization for different tensor
sizes
33 *     - Processes tensors in chunks to handle large batch sizes efficiently
34 *     - Automatically handles tensor dimensions > 65535 by splitting into
multiple kernel launches
35 *     - Optimized for both small and large padding sizes
36 *
37 * @note Requirements:
38 *     - Input tensor must fit into 32-bit index math
39 *     - Output tensor must have correct size: [N, C, D+pfloor+pback,
H+ptop+pbottom, W+pleft+pright]
40 *     - At least one of the output dimensions (D, H, W) must be >= 1
41 */

```

和replication_pad1d接口一样

```

42 template <typename T>
43 void replication_pad_out_kernel(
44     const skernel::Tensor<T>& self,
45     skernel::Tensor<T>& out,
46     const void* padding,
47     int padding_arr_size,
48 );

```

ISSUE

- 支持非contiguous，同时支持4D，5D输入
- cuda的实现通过分块处理超过 CUDA 网格限制（65535）的大型张量
- cuda的实现没有接收 *pright*, *pbottom*, *pback*。而是通过 *pleft*, *ptop* 和 *pfront* 计算出原始数据的起始位置进行填充，然后再根据 *output* 的大小，在 *padding* 部分直接复制最近的边界值

reflection_pad1d, reflection_pad1d.out

Pytorch 文档

- <https://pytorch.org/docs/stable/generated/torch.nn.ReflectionPad1d.html>

Pytorch 实现

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/ReflectionPad.cu#L430>
- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/ReflectionPad.cu#L96>
- tensor dim 数量必须为 2 或者 3:
<https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/PadNd.cpp#L205>
- pad_l/pad_r 约束: Padding size should be less than the corresponding input dimension

接口: 两个 tensor + pad_l + pad_r

代码块

```
1 Tensor & reflection_pad1d_out(const Tensor & self, c10::SymIntArrayRef
padding, Tensor & out); // {"schema": "aten::reflection_pad1d.out(Tensor self,
SymInt[2] padding, *, Tensor(a!) out) -> Tensor(a!)", "dispatch": "True",
"default": "False"}
2 Tensor reflection_pad1d(const Tensor & self, c10::SymIntArrayRef padding); //
{"schema": "aten::reflection_pad1d(Tensor self, SymInt[2] padding) -> Tensor",
"default": "True", "dispatch": "True"}
```

代码块

```
1 m = torch.nn.ReflectionPad1d(5)
2 input = torch.arange(8, dtype=torch.float).reshape(1, 4, 2)
3 print(input)
4 print(m(input))
5 => RuntimeError: Argument #4: Padding size should be less than the
corresponding input dimension, but got: padding (5, 5) at dimension 2 of input
```

[1, 4, 2]

Function 接口

代码块

```
1 template <typename T>
2 void reflection_pad_out(const skernel::Tensor<T>& self, skernel::Tensor<T>&
  out, const void* padding, int padding_arr_size);
```

reflection_pad2d, reflection_pad2d.out

Pytorch 文档

- <https://pytorch.org/docs/stable/generated/torch.nn.ReflectionPad3d.html>

Pytorch 实现

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/ReflectionPad.cu#L271>
- Tensor dim 数量必须为 3 或者 4:
<https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/PadNd.cpp#L212>
- 约束: Padding size should be less than the corresponding input dimension

接口: 两个 tensor + padding_left, padding_right, padding_top, padding_bottom

对应关系:

Shape:

- Input: (N, C, H_{in}, W_{in}) or (C, H_{in}, W_{in}) .
- Output: (N, C, H_{out}, W_{out}) or (C, H_{out}, W_{out}) where
 $H_{out} = H_{in} + padding_top + padding_bottom$
 $W_{out} = W_{in} + padding_left + padding_right$

代码块

```
1 Tensor & reflection_pad2d_out(const Tensor & self, c10::SymIntArrayRef
  padding, Tensor & out); // {"schema": "aten::reflection_pad2d.out(Tensor self,
  SymInt[4] padding, *, Tensor(a!)) out) -> Tensor(a!)", "dispatch": "True",
  "default": "False"}
2 Tensor reflection_pad2d(const Tensor & self, c10::SymIntArrayRef padding); //
  {"schema": "aten::reflection_pad2d(Tensor self, SymInt[4] padding) -> Tensor",
  "dispatch": "True", "default": "False"}
```

```

1 m = torch.nn.ReflectionPad2d(3)
2 input = torch.arange(8, dtype=torch.float).reshape(1, 2, 4)
3 print(input)
4 print(m(input))
5 => RuntimeError: Argument #6: Padding size should be less than the
corresponding input dimension, but got: padding (3, 3) at dimension 1 of input
[1, 2, 4]

```

Fuction接口

同reflection_pad1d, 接口复用。

reflection_pad3d, reflection_pad3d.out

Pytorch 文档

- <https://pytorch.org/docs/stable/generated/torch.nn.ReflectionPad3d.html>

Pytorch 实现

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/ReflectionPad.cu#L570>
- Tensor dim 数量必须为 4或者5:
<https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/PadNd.cpp#L219>
- 约束: Padding size should be less than the corresponding input dimension

接口: 两个 tensor + padding_left, padding_right, padding_top, padding_bottom, padding_front, padding_back

对应关系:

Shape:

- Input: $(N, C, D_{in}, H_{in}, W_{in})$ or $(C, D_{in}, H_{in}, W_{in})$.
- Output: $(N, C, D_{out}, H_{out}, W_{out})$ or $(C, D_{out}, H_{out}, W_{out})$, where
 $D_{out} = D_{in} + \text{padding_front} + \text{padding_back}$
 $H_{out} = H_{in} + \text{padding_top} + \text{padding_bottom}$
 $W_{out} = W_{in} + \text{padding_left} + \text{padding_right}$

代码块

```

1 Tensor & reflection_pad3d_out(const Tensor & self, c10::SymIntArrayRef
padding, Tensor & out); // {"schema": "aten::reflection_pad3d.out(Tensor self,
SymInt[6] padding, *, Tensor(a!)) out) -> Tensor(a!)", "dispatch": "True",
"default": "False"}

```

```
2 Tensor reflection_pad3d(const Tensor & self, c10::SymIntArrayRef padding); //  
{"schema": "aten::reflection_pad3d(Tensor self, SymInt[6] padding) -> Tensor",  
"dispatch": "True", "default": "True"}
```

代码块

```
1 m = torch.nn.ReflectionPad3d(3)  
2 input = torch.arange(32, dtype=torch.float).reshape(1, 1, 2, 4, 4)  
3 print(input)  
4 print(m(input))  
5  
6 RuntimeError: Argument #8: Padding size should be less than the corresponding  
input dimension, but got: padding (3, 3) at dimension 2 of input [1, 1, 2, 4,  
4]
```

Fuction接口

同reflection_pad1d, 接口复用。

scatter_.src, scatter.src, scatter.src_out

Pytorch 文档

- https://pytorch.org/docs/stable/generated/torch.Tensor.scatter_.html#torch.Tensor.scatter_ 的上半部分

Pytorch 实现

- 非 deterministic
- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/ScatterGatherKernel.cu#L477>
- index dim 数量与 src 相同, 但 sizes 不要求和 src 相同, 详见 pytorch 文档

代码块

```
1 Tensor scatter(const Tensor & self, int64_t dim, const Tensor & index, const  
Tensor & src); // {"schema": "aten::scatter.src(Tensor self, int dim, Tensor  
index, Tensor src) -> Tensor", "dispatch": "True", "default": "True"}  
2 Tensor & scatter_(Tensor & self, int64_t dim, const Tensor & index, const  
Tensor & src); // {"schema": "aten::scatter_.src(Tensor(a!) self, int dim,  
Tensor index, Tensor src) -> Tensor(a!)", "dispatch": "True", "default":  
"True"}  
3 Tensor & scatter_out(const Tensor & self, int64_t dim, const Tensor & index,  
const Tensor & src, Tensor & out); // {"schema": "aten::scatter.src_out(Tensor
```

```
self, int dim, Tensor index, Tensor src, *, Tensor(a!) out) -> Tensor(a!)",  
"dispatch": "True", "default": "False"}
```

和 gather 为互逆操作

代码块

```
1 import torch  
2 self = torch.zeros(3, 3, dtype=torch.int64)  
3 src = torch.tensor([  
4     [1, 2, 3],  
5     [4, 5, 6],  
6     [7, 8, 9]  
7 ])  
8 index = torch.tensor([  
9     [0, 2, 1],  
10    [1, 0, 2],  
11    [2, 1, 0]  
12 ])  
13 self.scatter_(dim=1, index=index, src=src)  
14 new_src = self.gather(dim=1, index=index)  
15 # src == new_src
```

Dtypes

- AT_ALL_TYPES: uint8_t, int8_t, int16_t, int, int64_t, float, double
- AT_COMPLEX_TYPES: c10::complex<float>, c10::complex<double>
- Half, Bool, BFloat16: float16, bool, bfloat16

代码块

```
1 AT_DISPATCH_ALL_TYPES_AND_COMPLEX_AND3(  
2     at::ScalarType::Half, at::ScalarType::Bool, at::ScalarType::BFloat16
```

scatter_.value, scatter.value, scatter.value_out

Pytorch 文档

- https://pytorch.org/docs/2.5/generated/torch.Tensor.scatter_.html#torch.Tensor.scatter_ 的下半部分

Pytorch 实现

- deterministic

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/ScatterGatherKernel.cu#L485>

代码块

```
1 Tensor scatter(const Tensor & self, int64_t dim, const Tensor & index, const
  Scalar & value); // {"schema": "aten::scatter.value(Tensor self, int dim,
  Tensor index, Scalar value) -> Tensor", "dispatch": "True", "default": "True"}
2 Tensor & scatter_(Tensor & self, int64_t dim, const Tensor & index, const
  Scalar & value); // {"schema": "aten::scatter_.value(Tensor(a!) self, int dim,
  Tensor index, Scalar value) -> Tensor(a!)", "dispatch": "True", "default":
  "True"}
3 Tensor & scatter_out(const Tensor & self, int64_t dim, const Tensor & index,
  const Scalar & value, Tensor & out); // {"schema":
  "aten::scatter.value_out(Tensor self, int dim, Tensor index, Scalar value, *,
  Tensor(a!) out) -> Tensor(a!)", "dispatch": "True", "default": "False"}
```

Dtypes

- AT_ALL_TYPES: `uint8_t, int8_t, int16_t, int, int64_t, float, double`
- AT_COMPLEX_TYPES: `c10::complex<float>, c10::complex<double>`
- Half, Bool, BFloat16: `float16, bool, bfloat16`

代码块

```
1 AT_DISPATCH_ALL_TYPES_AND_COMPLEX_AND3(
2     at::ScalarType::Half, at::ScalarType::Bool, at::ScalarType::BFloat16
```

sum, sum.dim_IntList, sum.IntList_out

参考信息

文档: <https://pytorch.org/docs/stable/generated/torch.sum.html>

参考实现:

<https://github.com/pytorch/pytorch/blob/32f585d9346e316e554c8d9bf7548af9f62141fc/aten/src/ATen/native/cuda/ReduceSumProdKernel.cu#L173>

注册到 `sum_stub`

代码块

```
1 import torch
2
3 # 基本用法 - torch.sum
4 # 创建一个示例张量
```

```

5 x = torch.tensor([[1, 2, 3], [4, 5, 6]], dtype=torch.float)
6
7 # 对整个张量求和
8 total_sum = torch.sum(x)
9 print(f"整个张量的和: {total_sum}") # tensor(21.)
10
11 # 沿着指定维度求和
12 row_sum = torch.sum(x, dim=0) # 沿行方向求和 (合并行)
13 print(f"行求和结果: {row_sum}") # tensor([5., 7., 9.])
14
15 col_sum = torch.sum(x, dim=1) # 沿列方向求和 (合并列)
16 print(f"列求和结果: {col_sum}") # tensor([6., 15.])
17
18 # 使用keepdim=True保持维度
19 row_sum_keep = torch.sum(x, dim=0, keepdim=True)
20 print(f"保持维度的行求和: {row_sum_keep}") # tensor([[5., 7., 9.]])
21 print(f"形状: {row_sum_keep.shape}") # torch.Size([1, 3])
22
23 # 多维度求和 - sum.dim_IntList
24 # 创建一个3D张量
25 y = torch.tensor([[[[1, 2], [3, 4]],
26                    [[5, 6], [7, 8]]], dtype=torch.float)
27 print(f"原始形状: {y.shape}") # torch.Size([2, 2, 2])
28
29 # 沿多个维度求和 (dim=[0, 2])
30 multi_dim_sum = torch.sum(y, dim=[0, 2])
31 print(f"沿维度[0,2]的求和: {multi_dim_sum}") # tensor([14., 22.])
32 print(f"结果形状: {multi_dim_sum.shape}") # torch.Size([2])
33
34 # 保持维度
35 multi_dim_sum_keep = torch.sum(y, dim=[0, 2], keepdim=True)
36 print(f"保持维度的多维求和: {multi_dim_sum_keep}")
37 # tensor([[[[14.],
38           [22.]]]])
39 print(f"结果形状: {multi_dim_sum_keep.shape}") # torch.Size([1, 2, 1])
40
41

```

Function 接口

代码块

```

1 /**
2  * @brief Performs tensor summation along specified dimensions
3  *
4  * @param dev Device handle for CUDA/CPU device management
5  * @param out Output tensor to store the sum result

```

```

6 * @param self Input tensor to be summed
7 * @param out_sizes Output tensor sizes
8 * @param out_strides Output tensor strides (can be non-contiguous)
9 * @param self_sizes Input tensor sizes
10 * @param self_strides Input tensor strides (can be non-contiguous)
11 * @param dim Array of dimensions to reduce along
12 * @param out_ndims Number of dimensions to reduce (length of dims
array)
13 * @param self_ndims Number of dimensions in the input tensor
14 * @param dim_ndims Number of dimensions in the input tensor
15 * @param out_dtype Data type of output tensor
16 * @param input_dtype Data type of input tensor
17 * @param keepdim Whether to keep the reduced dimensions with size 1
18 *
19 * @note Memory Layout:
20 * - Fully supports non-contiguous tensors
21 * - Handles arbitrary strides efficiently
22 * - Optimized for both contiguous and non-contiguous cases
23 *
24 * @note Behavior:
25 * - Accumulates sum in higher precision for floating point types
26 * - Performs efficient reduction using appropriate algorithms for each
device
27 * - Handles empty tensors and edge cases correctly
28 * - Supports type promotion from input to output if needed
29 */

```

```

30 void sum_kernel(
31     SiDeviceHandle *dev, // Device handle
32     void *out, // Output tensor data pointer
33     const void *self, // Input tensor data pointer
34     const int64_t* out_sizes, // Output sizes
35     const int64_t* out_strides, // Output strides
36     const int64_t* self_sizes, // Input tensor sizes
37     const int64_t* self_strides, // Input tensor strides
38     const int64_t* dim, // Dimensions to sum over
39     int out_ndims, // Number of dimensions in input tensor
40     int self_ndims, // Number of dimensions in input tensor
41     int dim_ndims, // Number of dimensions to reduce
42     int out_dtype, // Data type of output tensor
43     int input_dtype, // Data type of input tensor
44     bool keepdim // Whether to keep reduced dimensions as
size 1
45 );

```

46 支持的数据类型

47 整数类型: int8, uint8, int16, int32, int64

48 浮点类型: float16, bfloat16, float32, float64

49 布尔类型: bool

```
50 复数类型: complex64, complex128
51 复数半精度: complex32
```

首个版本开发说明

- (1) 初版只需要支持dim=-1, keepdim=true的参数配置
- (2) keepdim相当于unsqueeze操作, 框架可以处理, 算子不需要做相关适配
- (3) 数制支持上, 初版只需要支持float类型
- (4) 初版可以只考虑tensor为continuous的情况
- (5) 执行流程为: torch将输入tensor (shape为[A, B, C]) 进行shape变换, 转为[A * B, C]传入skernel, skernel对最后一维进行累加计算, 输出shape为[A * B]的tensor, 然后torch将其还原到[A, B, 1]的shape形式

内部接口调用形式:

代码块

```
1  template <typename T>
2  void sum(void *d_in, void *d_out, int dim0_size, int dim1_size, int dim = -1);
```

备注:

- (1) 初版dim1_size参数, 可实现为1024字节对齐 (256个float), 配合tile reg的load
实际开发情况, 可将对齐要求缩减至32字节对齐 (8个float)

ISSUE

- Cuda reduce sum kernel 通过并行处理和多级归约来实现, 如果使用 sum_stub 来获取参数, 此时只能获得 input 以及 output 的 shape/stride。上面接口提供的是完整参数列表, 根据具体 sum sipu 实现, 可以筛选其中一部分的参数作为接口使用
- 接受非 contiguous

zero_

Pytorch 文档: https://pytorch.org/docs/stable/generated/torch.Tensor.zero_.html

代码块

```
1  Tensor& zero_(Tensor &self) {
2      int64_t nelements = c10::multiply_integers(self.sizes());
3      if (self.device() == at::kCPU &&
4          self.is_non_overlapping_and_dense() &&
5          nelements < internal::GRAIN_SIZE) {
6          return zero_cpu_(self, nelements);
```

```

7     }
8     return self.fill_(0);
9 }
10
11 非cpu backend会直接调用fill_
12
13
14 import torch
15
16 # 创建一个张量
17 x = torch.ones(3, 4)
18 print(x)
19 tensor([[1., 1., 1., 1.],
20         [1., 1., 1., 1.],
21         [1., 1., 1., 1.]])
22
23 # 将所有元素设置为零
24 x.zero_()
25 print(x)
26 tensor([[0., 0., 0., 0.],
27         [0., 0., 0., 0.],
28         [0., 0., 0., 0.]])

```

代码块

```

1 void zero_(SiDeviceHandle dev, void const* self, int dims, const int64_t*
  self_sizes, const int64_t* self_strides, dtype)
2

```

topk, topk.values

Pytorch 文档

- <https://pytorch.org/docs/stable/generated/torch.topk.html>

Pytorch 实现

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/TensorTopK.cpp#L47>
- 特定情况下直接 sort+narrow，其他情况下实现看起来是基于基数排序实现的（具体实现细节未细看）。依赖 cub 相关 apis。
- indices 默认 dtype 为 int64

```

1  ...::std::tuple<Tensor &,Tensor &> topk_out(const Tensor & self, c10::SymInt k,
    int64_t dim, bool largest, bool sorted, Tensor & values, Tensor & indices); //
    {"schema": "aten::topk.values(Tensor self, SymInt k, int dim=-1, bool
    largest=True, bool sorted=True, *, Tensor(a!) values, Tensor(b!) indices) ->
    (Tensor(a!) values, Tensor(b!) indices)", "dispatch": "True", "default":
    "False"}
2  ...::std::tuple<Tensor,Tensor> topk(const Tensor & self, c10::SymInt k, int64_t
    dim, bool largest, bool sorted); // {"schema": "aten::topk(Tensor self, SymInt
    k, int dim=-1, bool largest=True, bool sorted=True) -> (Tensor values, Tensor
    indices)", "dispatch": "True", "default": "True"}

```

Dtypes

- AT_ALL_TYPES: `uint8_t, int8_t, int16_t, int, int64_t, float, double`
- Half, BFloat16: `float16, bfloat16`

代码块

```
1  AT_DISPATCH_ALL_TYPES_AND2(at::ScalarType::Half, at::ScalarType::BFloat16
```

接口

- 在调用 kernel 前保证 contiguous, 因此不传 strides

代码块

```

1  /**
2   * @brief Compute the Top-K elements (values and indices) along a specified
    dimension of a tensor.
3   *
4   * @details
5   * Extracts the top K largest/smallest elements from the input tensor along
    the given dimension,
6   * writing results to separate value and index buffers. Supports both sorted
    and unsorted output.
7   *
8   * @param[in] dev          Device handle (e.g., GPU/accelerator) for
    computation
9   * @param[out] values      Output buffer for top K values (shape: original
    tensor with target dim set to K)
10  * @param[out] indices     Output buffer for original indices of top K values
    (int64_t data type)
11  * @param[in] self        Input tensor data pointer (device memory)
12  * @param[in] dims        Number of dimensions in the input tensor
13  * @param[in] self_sizes  Array of dimension sizes (length = dims)

```

```

14 * @param[in] k Number of elements to select (0 < k <= self_sizes[dim])
15 * @param[in] dim Target dimension index (0-based, 0 <= dim < dims)
16 * @param[in] largest Select largest elements if true, smallest otherwise
17 * @param[in] sorted Sort results in descending order if true, preserve
    original order otherwise
18 *
19 */
20 void topk(SiDeviceHandle dev, void* values, void* indices, const void* self,
    int dims, const int64_t* self_sizes, int k, int dim, bool largest, bool sorted)
21
22 void topk(void* values, void* indices, const void* self, const int64_t
    batch_size, const int64_t size, int k, bool largest, bool sorted)

```

sigmoid.out, sigmoid, sigmoid_

pytorch 文档

- <https://pytorch.org/docs/stable/generated/torch.nn.Sigmoid.html>

pytorch 实现

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/UnarySpecialOpKernel.cu#L124>
- 使用 gpu_kernel + element-wise lambda function 实现

Unary op, 接口定义同 neg

代码块

```

1 Tensor sigmoid(const Tensor & self); // {"schema": "aten::sigmoid(Tensor self)
    -> Tensor", "dispatch": "True", "default": "True"}
2 Tensor & sigmoid_(Tensor & self); // {"schema": "aten::sigmoid_(Tensor(a!)
    self) -> Tensor(a!)", "dispatch": "True", "default": "True"}
3 Tensor & sigmoid_out(const Tensor & self, Tensor & out); // {"schema":
    "aten::sigmoid.out(Tensor self, *, Tensor(a!) out) -> Tensor(a!)", "dispatch":
    "True", "default": "False"}

```

Dtypes

- AT_FLOATING_TYPES: float, double
- Half, BFloat16: float16, bfloat16
- isComplexType: c10::complex<c10::Half>, c10::complex<float>, c10::complex<double>

```
1 isComplexType(common_dtype)
2 AT_DISPATCH_FLOATING_TYPES_AND2(at::ScalarType::Half, at::ScalarType::BFloat16
```

rsqrt,rsqrt_,rsqrt.out

Pytorch 文档

- <https://pytorch.org/docs/stable/generated/torch.rsqrt.html>

Pytorch 实现

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/UnaryOpsKernel.cu#L96>
- 同样是 gpu_kernel + elementwise functor

Unary op, 接口定义同 neg

代码块

```
1 Tensor rsqrt(const Tensor & self); // {"schema": "aten::rsqrt(Tensor self) ->
  Tensor", "dispatch": "True", "default": "True"}
2 Tensor & rsqrt_(Tensor & self); // {"schema": "aten::rsqrt_(Tensor(a!) self) -
  > Tensor(a!)", "dispatch": "True", "default": "True"}
3 Tensor & rsqrt_out(const Tensor & self, Tensor & out); // {"schema":
  "aten::rsqrt.out(Tensor self, *, Tensor(a!) out) -> Tensor(a!)", "dispatch":
  "True", "default": "False"}
```

Dtypes

- AT_FLOATING_TYPES: float, double
- Half, BFloat16: float16, bfloat16
- isComplexType: c10::complex<c10::Half>, c10::complex<float>, c10::complex<double>

代码块

```
1 AT_DISPATCH_FLOATING_TYPES_AND2(
2     ScalarType::BFloat16, ScalarType::Half,
3     inline bool isComplexType(ScalarType t) {
4         return (
5             t == ScalarType::ComplexHalf || t == ScalarType::ComplexFloat ||
6             t == ScalarType::ComplexDouble);
7     }
```

首个版本开发说明

(1) 数据类型支持float32, bfloat16, float16

代码块

```
1 //calculate the reciprocal of the square-root of the input tensor with size
  dim0*dim1.
2 //total size should be 1024B aligned
3 template <typename T>
4 //only sifmt::bfloat16, sifmt::float16, //sifmt::float32 are supported
5 void rsqrt(void *d_in, void *d_out, int dim0_size, int dim1_size );
```

sort,sort.stable,sort.values_stable

Pytorch 文档

- <https://pytorch.org/docs/stable/generated/torch.sort.html>

Pytorch 实现

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/SortStable.cu#L218>
- 基于基数排序实现, cub 相关 apis 具体实现在 <https://github.com/NVIDIA/cccl/blob/main/cub/cub/device/>

代码块

```
1 // 在 RegisterCUDA.cpp 实现了 sort/sort_out 接口, 调用的同样是
  structured_sort_stable_out
2 ::std::tuple<Tensor, Tensor> sort(const Tensor & self, ::std::optional<bool>
  stable, int64_t dim, bool descending); // {"schema": "aten::sort.stable(Tensor
  self, *, bool? stable, int dim=-1, bool descending=False) -> (Tensor values,
  Tensor indices)", "dispatch": "True", "default": "True"}
3 ::std::tuple<Tensor &, Tensor &> sort_out(const Tensor & self,
  ::std::optional<bool> stable, int64_t dim, bool descending, Tensor & values,
  Tensor & indices); // {"schema": "aten::sort.values_stable(Tensor self, *,
  bool? stable, int dim=-1, bool descending=False, Tensor(a!) values, Tensor(b!)
  indices) -> (Tensor(a!) values, Tensor(b!) indices)", "dispatch": "True",
  "default": "False"}
```

Dtypes

- AT_ALL_TYPES: uint8_t, int8_t, int16_t, int, int64_t, float, double
- Half, Bool, BFloat16: float16, bool, bfloat16

代码块

```
1 AT_DISPATCH_ALL_TYPES_AND3(  
2     kBool, kHalf, kBFloat16
```

reciprocal,reciprocal_,reciprocal.out

Pytorch 文档

- <https://pytorch.org/docs/stable/generated/torch.reciprocal.html>
- 输入为 integral types 时, 会将 input 转换为 float types, 具体的处理逻辑在 TensorIteratror 的 UNARY_FLOAT_OP_CONFIG () 里的 promote_inputs_to_common_dtype 以及 compute_types 里。因此需要处理的 dtypes 为 float/complex types

Pytorch 实现

- <https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/UnaryFractionKernels.cu#L99>
- 实现为 gpu_kernel + unary functor

Unary op, 接口定义同 neg

代码块

```
1 Tensor reciprocal(const Tensor & self); // {"schema": "aten::reciprocal(Tensor  
self) -> Tensor", "dispatch": "True", "default": "True"}  
2 Tensor & reciprocal_(Tensor & self); // {"schema":  
"aten::reciprocal_(Tensor(a!) self) -> Tensor(a!)", "dispatch": "True",  
"default": "True"}  
3 Tensor & reciprocal_out(const Tensor & self, Tensor & out); // {"schema":  
"aten::reciprocal.out(Tensor self, *, Tensor(a!) out) -> Tensor(a!)",  
"dispatch": "True", "default": "False"}
```

dtypes

- AT_FLOATING_TYPES: float, double
- AT_COMPLEX_TYPES: c10::complex<float>, c10::complex<double>
- Half, BFloat16: float16, bfloat16

代码块

```
1 AT_DISPATCH_FLOATING_AND_COMPLEX_TYPES_AND2(  
2     ScalarType::Half, ScalarType::BFloat16
```

where.self,where.self_out,where,where.ScalarSelf,where.ScalarOther,where.Scalar

具体选择哪个搜索 `static PyObject* THPVariable_where (PyObject* self_, PyObject* args, PyObject* kwargs)`

参考信息

参考实现:

<https://github.com/pytorch/pytorch/blob/9cd52da45c3030601bdd27e8a948f23be12b21ee/aten/src/ATen/native/TensorCompare.cpp#L631>

资料:

<https://pytorch.org/docs/stable/generated/torch.where.html>

其中, `where.self` 通过 `where.self_out` 实现, 可以搜索 `Tensor& where_self_out (const Tensor& condition, const Tensor& self, const Tensor& other, Tensor& out)`

Function 接口

`self` 和 `other` 的 `tensor size` 和 `output tensor` 一致

说明: `condition` 判断条件, `pytorch` 会把 `condition` 处理为一个 `bool tensor`, 该 `tensor` 的 `shape` 和 `out` 完全一致。该算子实现, 对于 `out[i]`, 根据 `condition tensor`, 如果该位置为 `true`, 就赋值 `self[i]`, 否则赋值 `other[i]`。不区分 `scalar` 还是 `tensor`, 框架侧 `padding` 广播好数据。

代码块

```
1 // schema: where.self_out(Tensor condition, Tensor self, Tensor other, *,
2 // Tensor(a!) out) -> Tensor(a!)
3 /**
4 out_ndims: output tensor dims
5 out_sizes: output tensor sizes
6 out_strides: strides of output tensor, 0 means it's a broadcasted dim
7 self_strides: strides of self tensor, 0 means it's a broadcasted dim
8 other_strides: strides of other tensor, 0 means it's a broadcasted dim
9 element_sizes: element_size of nargs tensor; if element_sizes is nullptr, then
10 the strides will be in bytes, otherwise the strides will be in # of elements.
11 */
12 void fn(dev, const void* out, const void* self, const void* other, int
13 out_ndims, const int64_t* out_sizes, const int64_t* out_strides, const
14 int64_t* self_strides, const int64_t* other_strides, const int64_t*
15 element_sizes=nullptr)
16
17 接口修改如下:
18 template<typename T>
19 void where(const void* out, const void* condition, const void* self, const
20 void* other, int out_ndims, const int64_t* out_shapes, const int64_t*
21 out_strides, const int64_t* self_strides, const int64_t* other_strides)
```

数据类型

- kComplexHalf, kHalf, kBFloat16, kBool
- AT_EXPAND(AT_ALL_TYPES_AND_COMPLEX): c10::kByte, c10::kChar, c10::kInt, c10::kLong, c10::kShort, c10::kDouble, c10::kFloat
- AT_EXPAND(AT_FLOAT8_TYPES))

代码例子

代码块

```

1  import torch
2
3  x = torch.randn(3, 2)
4  print(x)
5  y = torch.where(x > 0, -1, torch.tensor(1))
6  print(y)
7
8  #####输出结果#####
9  x=tensor([[ -0.9950, -0.7007],
10           [ 0.1992, -0.4350],
11           [ 0.6579, -0.3728]])
12  self_converted.sizes: [], self_converted.strides: [] # scalar tensor的sizes和
    strides为空
13  condition tensor:  0  0
14     1  0
15     1  0
16  [ CPUBoolType{3,2} ]
17  out shapes: 6, # iter.shape()
18  iter tensor1 strides: [2, 1], iter tensor1 sizes: [3, 2]
19  iter tensor2 strides: [], iter tensor2 sizes: [] # scalar tensor经过
    tensoriterator处理, size和strides仍为空
20  iter tensor3 strides: [], iter tensor3 sizes: []
21  iter strides(1): [1] # 广播后会条件展平为一维处理, 压缩为线性。
    coalesce_dimensions处理逻辑。
22  iter strides(2): [0]
23  iter strides(3): [0]
24  y=tensor([[ 1,  1],
25            [-1,  1],
26            [-1,  1]])

```

注册在 build/aten/src/ATen/RegisterCompositeExplicitAutograd.cpp

使用到的 op: zeros_like, empty_strided, copy_, empty_symint

view, view_as, reshape, _reshape_alias

参考资料: <https://pytorch.org/docs/stable/generated/torch.Tensor.view.html>

<https://pytorch.org/docs/stable/generated/torch.reshape.html>

https://pytorch.org/docs/stable/generated/torch.Tensor.view_as.html

view, view_as 和 reshape 底层调用的都是 **alias_with_sizes_and_strides** 中的 **make_tensor**, 不需要算子实现

代码块

```
1  template <template <typename...> typename Container>
2  Tensor alias_with_sizes_and_strides(
3      const Tensor& self,
4      const Container<c10::SymInt>& sizes,
5      const Container<c10::SymInt>& strides) {
6      //caller should make sure that sizes and strides are valid for self
7      //(storage is sufficient, strides are non-negative, strides and sizes array
8      size is the same)
9      Tensor self_;
10     if (self.is_quantized()) {
11         self_ = at::detail::make_tensor<QTensorImpl>(
12             c10::TensorImpl::VIEW, Storage(self.storage()), self.key_set(),
13             self.dtype(), get_qtensorimpl(self)->quantizer());
14         self_.unsafeGetTensorImpl()->set_sizes_and_strides(sizes, strides,
15             self.sym_storage_offset());
16     } else {
17         self_ = at::detail::make_tensor<TensorImpl>(
18             c10::TensorImpl::VIEW, Storage(self.storage()), self.key_set(),
19             self.dtype());
20         self_.unsafeGetTensorImpl()->set_sizes_and_strides(sizes, strides,
21             self.sym_storage_offset());
22     }
23     namedinference::propagate_names(self_, self);
24     return self_;
25 }
```

view 操作要求能够找到一组有效的步长, 使得在不重新排列内存的情况下可以实现形状变换。这在以下情况下可能失败:

1. 非连续张量: 如果张量经过了转置等操作, 导致内存布局不连续
2. 跨越不连续子空间的维度: 当新形状要求将不连续的内存区域视为连续时

例如，对于一个形状为 [2, 3, 4] 的张量，如果经过 permute (2, 0, 1) 变为 [4, 2, 3]，然后尝试 view (8, 3)，这将失败，因为无法找到合适的步长。

与 view 一样，view_as 也要求张量在内存中是连续的

reshape: 内部会先检查是否可以直接使用 view; 如果不行，会先调用 contiguous () 创建一个连续的副本，再对这个副本使用 view

代码块

```
1 x = torch.randn(2, 3, 4)
2
3 # 将张量展平为一维
4 flat_view = x.view(-1)
5 print("展平后:", flat_view.shape) # torch.Size([24])
6
7 # 改变为 2x12 的形状, 第二维自动计算
8 reshaped = x.reshape(2, -1)
9 print("自动推断维度:", reshaped.shape) # torch.Size([2, 12])
10
11
12 # 添加一个维度 (批次维度)
13 batched = x.view(1, 2, 3, 4)
14 print("添加批次维度:", batched.shape) # torch.Size([1, 2, 3, 4])
15
16 # 移除一个维度 (当该维度大小为1时)
17 squeezed = batched.view(2, 3, 4)
18 print("移除批次维度:", squeezed.shape) # torch.Size([2, 3, 4])
19
20
21 # 创建一个非连续张量 (通过转置)
22 y = torch.randn(2, 3, 4)
23 y_transposed = y.transpose(0, 2) # 形状变为 [4, 3, 2]
24 print("转置后是否连续:", y_transposed.is_contiguous()) # False
25
26 # 使用 view - 会失败
27 try:
28     y_view = y_transposed.view(4, 6)
29 except RuntimeError as e:
30     print("view 错误:", e)
31
32 # 使用 reshape - 会成功
33 y_reshape = y_transposed.reshape(4, 6)
34 print("reshape 成功:", y_reshape.shape) # torch.Size([4, 6])
35
36 # 先使用 contiguous() 再使用 view
```

```

37 y_view = y_transposed.contiguous().view(4, 6)
38 print("contiguous().view() 成功:", y_view.shape) # torch.Size([4, 6])
39
40
41 x = torch.randn(4, 6) # 形状为 [4, 6]
42 y = torch.randn(2, 12) # 形状为 [2, 12]
43
44 # 使用 view_as 将 x 重塑为与 y 相同的形状
45 x_reshaped = x.view_as(y)
46
47 print("原始 x 形状:", x.shape)
48 print("y 形状:", y.shape)
49 print("重塑后的 x 形状:", x_reshaped.shape)
50 # 原始 x 形状: torch.Size([4, 6])
51 # y 形状: torch.Size([2, 12])
52 # 重塑后的 x 形状: torch.Size([2, 12])

```

uniform_, uniform, uniform.out

参考信息

Pytorch 文档

- https://pytorch.org/docs/stable/generated/torch.Tensor.uniform_.html

Pytorch 实现

- 调用了 curand 相关 api, generator impl 会调用到 cuda graph 相关 api
- Kernel 实现:
<https://github.com/pytorch/pytorch/blob/v2.5.0/aten/src/ATen/native/cuda/DistributionTemplates.h#L66>
- api 会需要支持传一个 generator, kernel 实现里会调用 generator 去具体生成, 建议等支持第二阶段联合编译后再实现该 kernel

代码块

```

1 Tensor & uniform_(Tensor & self, double from, double to,
  ::std::optional<Generator> generator); // {"schema":
  "aten::uniform_(Tensor(a!) self, float from=0, float to=1, *, Generator?
  generator=None) -> Tensor(a!)", "dispatch": "True", "default": "False"}

```

代码块

```

1 x = torch.empty(2, 3)
2 x.uniform_(from=0, to=1)
3 x

```

```

4  tensor([[0.2154, 0.5179, 0.8284],
5         [0.4732, 0.9251, 0.3088]])
6  g = torch.Generator().manual_seed(42)
7
8  y = torch.empty(2, 3)
9  y.uniform_(from_=10, to=20, generator=g)
10 y
11 tensor([[16.6065, 18.1294, 17.8853],
12         [18.9614, 13.6132, 14.2703]])

```

ISSUE

sym_size.int

把符号形状 (symbolic size) 转换成实际的整数 (int64_t) 值。

代码块

```

1  x = torch.zeros(2, 3, 4)
2  sym_dims = x.sym_size()
3  sym_dims
4  [2, 3, 4]

```

参考信息

Pytorch 实现

<https://github.com/pytorch/pytorch/blob/9a0f7a3bb01b235ea04581ee540970a098071b72/aten/src/ATen/core/TensorBase.h#L165>

```

243     c10::SymIntArrayRef sym_sizes() const {
244         return impl_->sym_sizes();

```

<https://github.com/pytorch/pytorch/blob/9a0f7a3bb01b235ea04581ee540970a098071b72/aten/src/ATen/core/TensorBase.h#L1053>

```

1053     template <typename T, typename = enable_if_symint<T>>
1054     c10::SymIntArrayRef sizes(const TensorBase& t) { return t.sym_sizes(); }

```

<https://github.com/pytorch/pytorch/blob/9a0f7a3bb01b235ea04581ee540970a098071b72/aten/src/ATen/core/TensorBase.h#L243>

```

165     c10::SymInt sym_size(int64_t dim) const {
166         return impl_->sym_size(dim);
167     }

```

<https://github.com/pytorch/pytorch/blob/9a0f7a3bb01b235ea04581ee540970a098071b72/aten/src/ATen/core/TensorBase.h#L1058>

```
... 1058     template <typename T, typename = enable_if_symint<T>>
      1059         c10::SymInt size(const TensorBase& t, int64_t dim) { return t.sym_size(dim); }
```

【NOTE】不需要实现，功能由pytorch内部完成

polar, polar.out

参考信息

pytorch文档: <https://pytorch.org/docs/stable/generated/torch.polar.html>

参考实现:

<https://github.com/pytorch/pytorch/blob/183bca41de89825107dce1d1ecc1502d9993a684/aten/src/ATen/native/cuda/ComplexKernel.cu#L22>

代码块

```
1 >>> abs = torch.tensor([1.0, 2.0, 3.0])
2 >>> angle = torch.tensor([0.0, torch.pi/2, torch.pi])
3 >>> complex_tensor = torch.polar(abs, angle)
4 >>> print(complex_tensor)
5 tensor([ 1.0000e+00+0.0000e+00j, -8.7423e-08+2.0000e+00j,
6         -3.0000e+00-2.6227e-07j])
```

Function 接口

代码块

```
1 /**
2  * @brief Converts polar coordinates (magnitude, angle) to complex numbers
3  *
4  * Performs element-wise conversion from polar coordinates (r,  $\theta$ ) to complex
5  * numbers
6  * using the formula:  $z = r * (\cos(\theta) + i\sin(\theta))$ 
7  *
8  * @param dev Device handle for CUDA device management
9  * @param out Output complex tensor to store results
10 * @param abs Input tensor containing magnitude values (r)
11 * @param angle Input tensor containing angle values ( $\theta$ ) in radians
12 * @param out_sizes Output tensor sizes
13 * @param out_strides Output tensor strides
14 * @param abs_sizes abs tensor sizes
15 * @param abs_strides abs tensor strides
16 * @param angle_sizes Angle tensor sizes
```

```

16 * @param angle_strides Angle tensor strides
17 * @param ndims Number of dimensions in tensors
18 * @param out_dtype Output complex data type (complex32/64/128)
19 * @param input_dtype Input floating point data type (float16/32/64)
20 *
21 * @note Supported Data Types:
22 * Input (magnitude, angle):
23 * - float16 (half)
24 * - bfloat16
25 * - float32 (float)
26 * - float64 (double)
27 *
28 * Output (complex):
29 * - complex32 (half complex)
30 * - complex64 (single complex)
31 * - complex128 (double complex)
32 *
33 * @note Memory Layout:
34 * - Supports non-contiguous tensors
35 * - Handles arbitrary strides efficiently
36 * - Input tensors must have matching sizes
37 * - Output complex tensor has same shape as inputs
38 *
39 * @note Implementation Details:
40 * - Uses fast math functions for trigonometric operations
41 * - Vectorized operations where possible
42 * - Handles edge cases (inf, nan) according to IEEE standards
43 * - Thread-safe and optimized for GPU execution
44 */
45 void polar_kernel_cuda(
46     SiDeviceHandle *dev, // Device handle
47     void *out, // Output complex tensor
48     const void *abs, // Input magnitude tensor
49     const void *angle, // Input angle tensor
50     const int64_t* out_sizes, // Output sizes
51     const int64_t* out_strides, // Output strides
52     const int64_t* abs_sizes, // Magnitude sizes
53     const int64_t* abs_strides, // Magnitude strides
54     const int64_t* angle_sizes, // Angle sizes
55     const int64_t* angle_strides, // Angle strides
56     int ndims, // Number of dimensions
57     int out_dtype, // Output complex dtype
58     int input_dtype // Input float dtype
59 );

```

使用TensorIterator 传递参数，支持非contiguous tensor

rms_norm

参考信息

参考实现:

https://github.com/pytorch/pytorch/blob/a05cc9f4941bd3b93e44b21cb7dac47d37003b32/aten/src/ATen/native/layer_norm.cpp#L264

Pytorch资料:

<https://pytorch.org/docs/stable/generated/torch.nn.modules.normalization.RMSNorm.html>

代码块

```
1 func: rms_norm(Tensor input, int[] normalized_shape, Tensor? weight=None,
2               float? eps=None) -> Tensor
3 ##基本用法
4 import torch
5 import torch.nn.functional as F
6
7 input = torch.randn(2, 3, 4, 5)
8 normalized_shape = [4, 5] # 最后两个维度
9 weight = torch.ones(4, 5) # 可选的缩放参数
10
11 output = torch.nn.functional.rms_norm(
12     input,
13     normalized_shape,
14     weight=weight, # 可选
15     eps=1e-6      # 可选
16 )
```

Function 接口

代码块

```
1 /**
2  * @brief Performs Root Mean Square (RMS) Normalization on input tensor
3  *
4  * RMS Normalization normalizes the input tensor by dividing it by the root
5  * mean square,
6  * computed over the specified dimensions. It differs from Layer Normalization
7  * by not
8  * centering the inputs (no mean subtraction), which makes it more
9  * computationally efficient.
10  *
```

```

8  * Formula:  $y = [x / \sqrt{\text{mean}(x^2) + \text{eps}}] * \text{weight}$  (if provided)
9  *
10 * @param out Rms_norm result
11 * @param input Input tensor to be normalized
12 * @param weight Optional weight tensor for scaling the normalized
    output
13 *
    Must be of shape normalized_shape if provided
14 * @param eps Small constant added to denominator for numerical
    stability
15 *
    If not provided, uses default epsilon for the data
    type
16 * @param batch_size numel of input / numel of normalized tensor
17 * @param normalized_size padded numel of normalize tensor
18 *
    Normalization will be performed over the last dims
    of size normalized_shape
19 * @param original_normalized_size the original size of normalized_size
20 * @param ndims input/output ndims
21 * @param normalized_shape_ndims how many dims will be normalized
22 * @param weight_opt weither has weight input
23 * @param eps_opt weither has eps input
24 *
25 */
26 void rms_norm_kernel(
27     SiDeviceHandle *dev, // Device handle
28     void *output, // Output
29     const void *input, // Input tensor
30     const void *weight, // Optional weight tensor for scaling
    the normalized output
31     const float eps, // Small constant added to denominator
    for numerical stability If not provided, uses default epsilon for the data type
32     const int64_t batch_size, // numel of input / numel of normalized
    tensor 为非normalize维度的乘积
33     const int64_t normalized_size, // padded numel of normalize tensor, 按照
    512pad之后的original_normalized_size
34     const int64_t original_normalized_size, //original numel of normalize
    tensor 一般为input的最后几个维度的乘积
35     int weight_opt, // 0 for false, 1 for true
36     int eps_opt, // 0 for false, 1 for true
37 );
38
39 rms_norm的input size = output size
40
41 weight和eps都是可选参数, 其中weight的shape是normalized_shape, eps Defaults to 1e-6
42
43 传入input shape=[batch_size, normalized_size]
44 weight shape = [normalized_size]
45 out shape = [batch_size, normalized_size]

```

```
46
47 input cuda支持的输入类型
48 1. float
49 2. double
50 3. complex<float>
51 4. complex<double>
52 5. Half (float16)
53 6. BFloat16
```

ISSUE

虽然cuda的实现中可以接收非contiguous的input, 但是对于contiguous的支持是非必要的, 可以在pytorch侧先将输入转换为contiguous。

amax, amax.out

参考信息

pytorch文档: <https://pytorch.org/docs/stable/generated/torch.amax.html>

参考实现: 和max共用一个kernel

与max的不同:

- 1. 维度reduce方面:** amax/amin能够支持在多个维度上进行**reduce**操作, 而max/min在这一点上不具备这样的功能。
- 2. 是否返回索引:** amax/amin在执行操作时不会返回对应的索引, 而max/min则会返回索引。
- 3. 梯度传播方式:** 当存在多个相等的值时, amax/amin会将梯度均匀地分配到这些相等的值上, 而max(dim)/min(dim)在传播梯度时, 只会将梯度传播到源张量中的一个索引位置。

使用方法:

代码块

```
1 >>> x = torch.tensor([[1, 2, 3], [4, 5, 6]])
2 >>> max_val = torch.amax(x)
3 >>> max_val
4 tensor(6)
5
6 >>> x = torch.tensor([[1, 2, 3], [4, 5, 6]])
7 >>> max_dim0 = torch.amax(x, dim=0)
8 >>> max_dim0
9 tensor([4, 5, 6])
10
11 >>> x = torch.tensor([[1, 2, 3], [4, 5, 6]])
12 >>> max_dim1_keep = torch.amax(x, dim=1, keepdim=True)
```

```

13 >>> max_dim1_keep
14 tensor([[3],
15         [6]])
16
17 >>> x = torch.randn(2, 3, 4)
18 >>> x
19 tensor([[[[ 1.3958e+00,  1.0015e+00, -3.7019e-01, -1.5914e+00],
20           [-7.9525e-01,  3.4663e-02, -1.8195e+00, -1.3115e-03],
21           [ 5.5305e-01,  2.5448e+00,  7.1271e-01,  9.1277e-02]],
22
23          [[ 1.0165e+00,  2.1791e-01,  1.3918e-01, -1.5510e-01],
24           [-8.1696e-01,  3.8472e-01,  1.5919e+00,  4.4249e-01],
25           [-5.2152e-01,  6.0341e-01,  9.4258e-02,  2.8002e-01]]]])
26 >>> max_multi = torch.amax(x, dim=(0, 1))
27 >>> max_multi
28 tensor([1.3958, 2.5448, 1.5919, 0.4425])

```

Function 接口

参考: [目OP 的开发文档](#)

与max不同的Reduce dim的处理和output处理在cpu torch侧处理

bincount, bincount.out

参考信息

pytorch文档: <https://pytorch.org/docs/2.6/generated/torch.bincount.html#torch-bincount>

cuda参考实现:

<https://github.com/pytorch/pytorch/blob/1a6d50d40799ba05b17254b236a6bf8817dd7d06/aten/src/ATen/native/cuda/SummaryOps.cu#L248>

注意: 只能处理一维的非负整数张量

代码块

```

1 >>> x = torch.tensor([1, 2, 2, 3, 3, 3])
2 >>> counts = torch.bincount(x)
3 >>> counts
4 tensor([0, 1, 2, 3])
5
6 >>> x = torch.tensor([1, 1, 2])
7 >>> counts = torch.bincount(x, minlength=5)
8 >>> counts
9 tensor([0, 2, 1, 0, 0])
10

```

```

11 >>> x = torch.tensor([1, 1, 2, 3])
12 >>> weights = torch.tensor([0.1, 0.2, 0.3, 0.4])
13 >>> weighted_counts = torch.bincount(x, weights=weights)
14 >>> print(weighted_counts)
15 tensor([0.0000, 0.3000, 0.3000, 0.4000])

```

Function 接口

代码块

```

1  /**
2   * @brief Counts the frequency of each value in a non-negative integer tensor
3   *
4   * Bincount computes the occurrence count of each value in the input tensor.
5   * It can optionally apply weights to compute weighted sums instead of counts.
6   *
7   * Formula:
8   * - Without weights: out[i] = count of i in input
9   * - With weights: out[i] = sum(weights[j] where input[j] == i)
10  *
11  * @param output Output tensor containing the counts/sums
12  * @param self Input tensor (must be 1D non-negative integers)
13  * @param weights Optional weights tensor for weighted sums
14  * Must be same size as input if provided
15  * @param minlength Minimum length of output tensor
16  * @param input_size Number of elements in input tensor
17  * @param output_size Size of output tensor (max(input) + 1 or minlength)
18  * @param has_weights Flag indicating if weights are provided (0=false,
19  * 1=true)
20  */
21 void bincount_kernel(
22     SiDeviceHandle* dev, // Device handle
23     void* output, // Output tensor buffer
24     const void* self, // Input tensor data
25     const void* weights, // Optional weights data
26     int64_t minlength, // Minimum output length
27     int64_t input_size, // Number of elements in input
28     int64_t output_size, // Size of output tensor
29     int has_weights // Whether weights are provided
30 );
31 self支持的数据类型:
32 uint8_t (通过特例检查)
33 int8_t
34 int16_t
35 int32_t
36 int64_t

```

```
37
38 weights支持的数据类型:
39 float
40 double
```

首个版本开发说明

- (1) 数制支持上，初版需要支持int64_t类型
 - (2) 初版暂时不用支持weights参数
 - (3) 输出tensor的长度上，初版先依赖min len，后续考虑是否支持max + 1的情况
 - (4) 新情况变更，由于硬件限制，初版支持int32_t类型，由框架将数据进行64-32转换
- 内部接口调用形式：

代码块

```
1 template <typename T>
2 void bincount(void *d_in, void *d_out, T dim_size, T min_len);
```

ISSUE

- 输入必须是1维连续的Strided布局张量,不支持任何形式的非连续或特殊布局
- 只能处理一维的非负整数张量

nonzero,nonzero.out

参考信息

pytorch文档: <https://pytorch.org/docs/stable/generated/torch.nonzero.html>

参考实现:

<https://github.com/pytorch/pytorch/blob/a936d596f6f7d2bc2dc47b4b2320208b4908e7f2/aten/src/ATen/native/cuda/Nonzero.cu#L170>

代码块

```
1 >>> x = torch.tensor([[1, 0, 0], [0, 2, 0], [0, 0, 3]])
2 >>> indices = torch.nonzero(x)
3 >>> indices
4 tensor([[0, 0],
5         [1, 1],
6         [2, 2]])
7
8 >>> x = torch.tensor([[1, 0, 0], [0, 2, 0], [0, 0, 3]])
```

```

9  >>> rows, cols = torch.nonzero(x, as_tuple=True)
10 >>> rows
11 tensor([0, 1, 2])
12 >>> cols
13 tensor([0, 1, 2])
14
15 >>> x = torch.tensor([[[[1, 0], [0, 2]], [[0, 0], [3, 0]]]])
16 >>> indices = torch.nonzero(x)
17 >>> indices
18 tensor([[0, 0, 0],
19         [0, 1, 1],
20         [1, 1, 0]])
21
22 mask = torch.tensor([[True, False], [False, True]])
23 indices = torch.nonzero(mask)
24 print(indices)
25 # 输出:
26 # tensor([[0, 0],
27 #         [1, 1]])
28
29 x = torch.tensor([1, 2, 3, 4, 5])
30 indices = torch.nonzero(x > 3)
31 print(indices) # tensor([[3], [4]])
32
33 >>> x = torch.tensor([[[[0, 0], [0, 0]], [[0, 0], [0, 0]]]])
34 >>> indices = torch.nonzero(x)
35 >>> indices
36 tensor([], size=(0, 3), dtype=torch.int64)

```

在cuda实现中，该算子是动态的，需要分两次调用cuda算子，

1. 第一次计算self中包含多少非0元素num_nonzeros_h,
 - a. `cub::TransformInputIterator<bool, NonZeroOp<scalar_t>, const scalar_t*>`
`itr(self._const_data_ptr<scalar_t>(), NonZeroOp<scalar_t>());`将input全部转为0/1，然后通过sum来计算出来非0元素
2. 根据第一次调用返回的数据，初始化output的size为{num_nonzeros, ndim}
3. 第二次调用cuda，得到indicates，如果使用tuple，需要在pytorch侧进行后续处理，不在kernel中处理

Function 接口

这里设计的是第二次调用cuda使用的接口，第一次调用可参考[目OP 的开发文档](#)sum 接口

代码块

```

1  /**
2   * @brief Finds indices of non-zero elements in a tensor
3   *
4   * Nonzero computes the indices of all non-zero elements in the input tensor.
5   * The output is a 2D tensor where each row contains the indices of a non-zero
6   * element.
7   *
8   * @param out          Output tensor containing the indices
9   * @param self         Input tensor
10  * @param self_sizes   input shape
11  * @param self_strides input strides
12  * @param num_nonzeros Number of non-zero elements in input
13  * @param ndim          Number of dimensions of input tensor
14  */
15 void nonzero_kernel(
16     SiDeviceHandle* dev,           // Device handle
17     void* out,                    // Output tensor buffer
18     const void* self,             // Input tensor data
19     const int64_t* self_sizes,     // input shape
20     int64_t num_nonzeros,         // Number of non-zero elements
21     int64_t ndim,                // Number of dimensions
22 );
23 支持的数据类型
24  float32
25  float64
26  int8/uint8
27  int16/int32/int64
28  bool
29  bfloat16
30  complex64/complex128

```

首个版本开发说明

(1) 数制支持上，需要支持int32_t, bool

(2) 支持as_tuple参数

as_tuple为true时，返回tensor为d个，分别存储不同维度的坐标值（二维输入即存储坐标的x值，y值）

as_tuple为false时，返回tensor为1个，shape为[MAX_NUM, DIM_NUM]

(3) 由于硬件限制，对于int64_t类型，由框架将数据进行64-32转换

(4) 初版仅实现二维输入接口，输出也为二维，按照最大值配置空间，假设input是[m, n]，输出每个tensor的大小为[m * n]，而每个输出tensor的真实大小，将存储于out_size返回参数中

as_tuple为true时，返回tensor为2个

as_tuple为false时, 返回tensor为1个, shape为[m * n, 2]

内部接口调用形式:

代码块

```
1  /// @brief Returns 2 tensor containing the indices of all non-zero elements of
    input
2  /// @tparam T Datatype support [int32_t, bool]
3  /// @param[in] d_in Input tensor with shape [dim0_size, dim1_size]
4  /// @param[out] d_out0 Output tensor 0 with shape [*out_size]
5  /// @param[out] d_out1 Output tensor 1 with shape [*out_size]
6  /// @param[in] dim0_size Tensor shape param 0
7  /// @param[in] dim1_size Tensor shape param 1
8  /// @param[out] out_size Record nonzero element number, this param should be
    on DEVICE MEMORY
9  template <typename T>
10 void nonzero_2d(void *d_in, void *d_out0, void *d_out1, int dim0_size, int
    dim1_size, int *out_size);
11
12 /// @brief Returns 1 tensor containing the indices of all non-zero elements of
    input
13 /// @tparam T Datatype support [int32_t, bool]
14 /// @param[in] d_in Input tensor with shape [dim0_size, dim1_size]
15 /// @param[out] d_out Output tensor with shape [*out_size, 2]
16 /// @param[in] dim0_size Tensor shape param 0
17 /// @param[in] dim1_size Tensor shape param 1
18 /// @param[out] out_size Record nonzero element number, this param should be
    on DEVICE MEMORY
19 template <typename T>
20 void nonzero_2d_no_tuple(void *d_in, void *d_out, int dim0_size, int
    dim1_size, int *out_size);
```

ISSUE

输入到kernel之前, self会先做contiguous处理, kernel接受的均为contiguous数据

vllm_flash_attn_c::varlen_fwd, vllm_flash_attn_c::fwd_kvcache[vLLM]

参考信息

vLLM 调用的地方在: https://github.com/vllm-project/vllm/blob/e92694b6fe264a85371317295bca6643508034ef/vllm/v1/attention/backends/flash_attn.py#L194

在这里分为 FA2 和 FA3 (hopper架构):

https://github.com/vllm-project/flash-attention/blob/main/vllm_flash_attn/flash_attn_interface.py

flash-attention/vllm_flash_attn/flash_attn_interface.py at main · vllm-project/flash-attent...

```
1 def flash_attn_varlen_func(  
2     q,  
3     k,  
4     v,  
5     max_seqlen_q,  
6     cu_seqlens_q,  
7     max_seqlen_k,  
8     cu_seqlens_k=None, # only used for non-paged prefill  
9     seqused_k=None,  
10    dropout_p=0.0,  
11    softmax_scale=None,  
12    causal=False,  
13    window_size: Optional[List[int]] = None,  
14    softcap=0.0, # 0.0 means deactivated  
15    alibi_slopes=None,  
16    deterministic=False,  
17    return_attn_probs=False,  
18    block_table=None,  
19    *,  
20    return_softmax_lse=False,  
21    out=None,  
22    fa_version: int = DEFAULT_FA_VERSION,  
23 ):  
24     """dropout_p should be set to 0.0 during evaluation  
25     Supports multi-query and grouped-query attention (MQA/GQA) by passing in  
26     K, V with fewer heads  
27     than Q. Note that the number of heads in Q must be divisible by the number  
28     of heads in KV.  
29     For example, if Q has 6 heads and K, V have 2 heads, head 0, 1, 2 of Q  
30     will attention to head  
31     0 of K, V, and head 3, 4, 5 of Q will attention to head 1 of K, V.  
32     If causal=True, the causal mask is aligned to the bottom right corner of  
33     the attention matrix.  
34     For example, if seqlen_q = 2 and seqlen_k = 5, the causal mask (1 = keep,  
35     0 = masked out) is:  
36     1 1 1 1 0  
37     1 1 1 1 1  
38     If seqlen_q = 5 and seqlen_k = 2, the causal mask is:  
39     0 0  
40     0 0  
41     0 0
```

```
38     1 0
39     1 1
40     If the row of the mask is all zero, the output will be zero.
41
42     If window_size != (-1, -1), implements sliding window local attention.
Query at position i
43     will only attend to keys between
44     [i + seqlen_k - seqlen_q - window_size[0], i + seqlen_k - seqlen_q +
window_size[1]] inclusive.
45
46     Arguments:
47         q: (total_q, nheads, headdim), where total_q = total number of query
tokens in the batch.
48         k: (total_k, nheads_k, headdim), where total_k = total number of key
tokens in the batch.
49         v: (total_k, nheads_k, headdim), where total_k = total number of key
tokens in the batch.
50         cu_seqLens_q: (batch_size + 1,), dtype torch.int32. The cumulative
sequence lengths
51         of the sequences in the batch, used to index into q.
52         cu_seqLens_k: (batch_size + 1,), dtype torch.int32. The cumulative
sequence lengths
53         of the sequences in the batch, used to index into kv.
54         max_seqLen_q: int. Maximum query sequence length in the batch.
55         max_seqLen_k: int. Maximum key sequence length in the batch.
56         dropout_p: float. Dropout probability.
57         softmax_scale: float. The scaling of QK^T before applying softmax.
58         Default to 1 / sqrt(headdim).
59         causal: bool. Whether to apply causal attention mask (e.g., for auto-
regressive modeling).
60         window_size: (left, right). If not (-1, -1), implements sliding window
local attention.
61         softcap: float. Anything > 0 activates softcapping attention.
62         alibi_slopes: (nheads,) or (batch_size, nheads), fp32. A bias of
63         (-alibi_slope * |i + seqLen_k - seqLen_q - j|)
64         is added to the attention score of query i and key j.
65         deterministic: bool. Whether to use the deterministic implementation
of the backward pass,
66         which is slightly slower and uses more memory. The forward pass is
always deterministic.
67         return_attn_probs: bool. Whether to return the attention
probabilities. This option is for
68         testing only. The returned probabilities are not guaranteed to be
correct
69         (they might not have the right scaling).
70     Return:
71     out: (total, nheads, headdim).
```

```

72     softmax_lse [optional, if return_softmax_lse=True]: (nheads,
total_q_seqlen). The
73     logsumexp of each row of the matrix  $QK^T$  * scaling (e.g., log of
the softmax
74     normalization factor).
75     """
76     assert cu_seqlens_k is not None or seqused_k is not None, \
77         "cu_seqlens_k or seqused_k must be provided"
78     assert cu_seqlens_k is None or seqused_k is None, \
79         "cu_seqlens_k and seqused_k cannot be provided at the same time"
80     assert block_table is None or seqused_k is not None, \
81         "seqused_k must be provided if block_table is provided"
82
83     if softmax_scale is None:
84         softmax_scale = q.shape[-1] ** (-0.5)
85     # custom op does not support non-tuple input
86     real_window_size: Tuple[int, int]
87     if window_size is None:
88         real_window_size = (-1, -1)
89     else:
90         assert len(window_size) == 2
91         real_window_size = (window_size[0], window_size[1])
92     q, k, v = [maybe_contiguous(x) for x in (q, k, v)]
93
94     dummy_cu_seqlens_k = torch.empty_like(cu_seqlens_q)
95
96     if fa_version == 2:
97         out, softmax_lse = torch.ops._vllm_fa2_C.varlen_fwd(
98             q, k, v,
99             out,
100             cu_seqlens_q,
101             # cu_seqlens_k not used since we use seqused_k, but flash_api.cpp
102             # still wants it so we pass all zeros
103             dummy_cu_seqlens_k if cu_seqlens_k is None else cu_seqlens_k,
104             seqused_k,
105             None,
106             block_table,
107             alibi_slopes,
108             max_seqlen_q,
109             max_seqlen_k,
110             dropout_p,
111             softmax_scale,
112             False,
113             causal,
114             real_window_size[0],
115             real_window_size[1],
116             softcap,

```

```

117         return_softmax_lse and dropout_p > 0,
118         None,
119     )
120     elif fa_version == 3:
121         out, softmax_lse, _, _ = torch.ops._vllm_fa3_C.fwd(
122             q, k, v,
123             None, None, # k_new, v_new
124             out,
125             cu_seqlens_q,
126             cu_seqlens_k, # cu_seqlens_k
127             None, # cu_seqlens_k_new
128             None, sequested_k, # sequested_q, sequested_k
129             max_seqlen_q, max_seqlen_k,
130             block_table,
131             alibi_slopes,
132             None, # kv_batch_idx
133             None, None, # rotary_cos, rotary_sin
134             None, None, None, # q_descale, k_descale, v_descale
135             softmax_scale,
136             causal,
137             real_window_size[0], real_window_size[1],
138             0, # sink_token_length
139             softcap,
140             True, # rotary_interleaved
141             0, # num_splits
142             None, # pack_gqa
143             0, # sm_margin
144         )
145     else:
146         raise ValueError(f"Unsupported FA version: {fa_version}")
147     return (out, softmax_lse) if return_softmax_lse else out
148
149
150
151 def flash_attn_with_kvcache(
152     q,
153     k_cache,
154     v_cache,
155     k=None,
156     v=None,
157     rotary_cos=None,
158     rotary_sin=None,
159     cache_seqlens: Optional[Union[(int, torch.Tensor)]] = None,
160     cache_batch_idx: Optional[torch.Tensor] = None,
161     cache_leftpad: Optional[torch.Tensor] = None,
162     block_table: Optional[torch.Tensor] = None,
163     softmax_scale=None,

```

```

164     causal=False,
165     window_size=(-1, -1), # -1 means infinite context window
166     softcap=0.0, # 0.0 means deactivated
167     rotary_interleaved=True,
168     alibi_slopes=None,
169     num_splits=0,
170     return_softmax_lse=False,
171     *,
172     out=None,
173     # FA3 Only
174     scheduler_metadata=None,
175     q_descale=None,
176     k_descale=None,
177     v_descale=None,
178     # Version selector
179     fa_version: int = DEFAULT_FA_VERSION,
180 ):
181     """
182     If k and v are not None, k_cache and v_cache will be updated *inplace*
183     with the new values from
184     k and v. This is useful for incremental decoding: you can pass in the
185     cached keys/values from
186     the previous step, and update them with the new keys/values from the
187     current step, and do
188     attention with the updated cache, all in 1 kernel.
189
190     If you pass in k / v, you must make sure that the cache is large enough to
191     hold the new values.
192     For example, the KV cache could be pre-allocated with the max sequence
193     length, and you can use
194     cache_seqLens to keep track of the current sequence lengths of each
195     sequence in the batch.
196
197     Also apply rotary embedding if rotary_cos and rotary_sin are passed in.
198     The key @k will be
199     rotated by rotary_cos and rotary_sin at indices cache_seqLens,
200     cache_seqLens + 1, etc.
201     If causal or local (i.e., window_size != (-1, -1)), the query @q will be
202     rotated by rotary_cos
203     and rotary_sin at indices cache_seqLens, cache_seqLens + 1, etc.
204     If not causal and not local, the query @q will be rotated by rotary_cos
205     and rotary_sin at
206     indices cache_seqLens only (i.e. we consider all tokens in @q to be at
207     position cache_seqLens).
208
209     See tests/test_flash_attn.py::test_flash_attn_kvcache for examples of how
210     to use this function.

```

199 Supports multi-query and grouped-query attention (MQA/GQA) by passing in
200 KV with fewer heads
201 than Q. Note that the number of heads in Q must be divisible by the number
of heads in KV.
202 For example, if Q has 6 heads and K, V have 2 heads, head 0, 1, 2 of Q
will attention to head
203 0 of K, V, and head 3, 4, 5 of Q will attention to head 1 of K, V.
204
205 If `causal=True`, the `causal` mask is aligned to the bottom right corner of
the attention matrix.
206 For example, if `seqlen_q = 2` and `seqlen_k = 5`, the `causal` mask (1 = keep,
0 = masked out) is:
207 1 1 1 1 0
208 1 1 1 1 1
209 If `seqlen_q = 5` and `seqlen_k = 2`, the `causal` mask is:
210 0 0
211 0 0
212 0 0
213 1 0
214 1 1
215 If the row of the mask is all zero, the output will be zero.
216
217 If `window_size != (-1, -1)`, implements sliding window local attention.
Query at position `i`
218 will only attend to keys between
219 `[i + seqlen_k - seqlen_q - window_size[0], i + seqlen_k - seqlen_q +`
`window_size[1]]` inclusive.
220
221 Note: Does not support backward pass.
222
223 Arguments:
224 `q`: (batch_size, seqlen, nheads, headdim)
225 `k_cache`: (batch_size_cache, seqlen_cache, nheads_k, headdim) if
there's no `block_table`,
226 or (num_blocks, page_block_size, nheads_k, headdim) if there's a
`block_table` (i.e. paged KV cache)
227 `page_block_size` must be a multiple of 256.
228 `v_cache`: (batch_size_cache, seqlen_cache, nheads_k, headdim) if
there's no `block_table`,
229 or (num_blocks, page_block_size, nheads_k, headdim) if there's a
`block_table` (i.e. paged KV cache)
230 `k` [optional]: (batch_size, seqlen_new, nheads_k, headdim). If not
None, we concatenate
231 `k` with `k_cache`, starting at the indices specified by `cache_seqLens`.
232 `v` [optional]: (batch_size, seqlen_new, nheads_k, headdim). Similar to
`k`.

233 rotary_cos [optional]: (seqlen_ro, rotary_dim / 2). If not None, we
apply rotary embedding
234 to k and q. Only applicable if k and v are passed in. rotary_dim
must be divisible by 16.
235 rotary_sin [optional]: (seqlen_ro, rotary_dim / 2). Similar to
rotary_cos.
236 cache_seqLens: int, or (batch_size,), dtype torch.int32. The sequence
lengths of the
237 KV cache.
238 block_table [optional]: (batch_size, max_num_blocks_per_seq), dtype
torch.int32.
239 cache_batch_idx: (batch_size,), dtype torch.int32. The indices used to
index into the KV cache.
240 If None, we assume that the batch indices are [0, 1, 2, ...,
batch_size - 1].
241 If the indices are not distinct, and k and v are provided, the
values updated in the cache
242 might come from any of the duplicate indices.
243 softmax_scale: float. The scaling of QK^T before applying softmax.
244 Default to 1 / sqrt(headDim).
245 causal: bool. Whether to apply causal attention mask (e.g., for auto-
regressive modeling).
246 window_size: (left, right). If not (-1, -1), implements sliding window
local attention.
247 softcap: float. Anything > 0 activates softcapping attention.
248 rotary_interleaved: bool. Only applicable if rotary_cos and rotary_sin
are passed in.
249 If True, rotary embedding will combine dimensions 0 & 1, 2 & 3,
etc. If False,
250 rotary embedding will combine dimensions 0 & rotary_dim / 2, 1 &
rotary_dim / 2 + 1
251 (i.e. GPT-NeoX style).
252 alibi_slopes: (nheads,) or (batch_size, nheads), fp32. A bias of
253 (-alibi_slope * |i + seqLen_k - seqLen_q - j|)
254 is added to the attention score of query i and key j.
255 num_splits: int. If > 1, split the key/value into this many chunks
along the sequence.
256 If num_splits == 1, we don't split the key/value. If num_splits ==
0, we use a heuristic
257 to automatically determine the number of splits.
258 Don't change this unless you know what you are doing.
259 return_softmax_lse: bool. Whether to return the logsumexp of the
attention scores.
260
261 Return:
262 out: (batch_size, seqLen, nheads, headDim).

```

263     softmax_lse [optional, if return_softmax_lse=True]: (batch_size,
264     nheads, seqlen). The
265     logsumexp of each row of the matrix  $QK^T * \text{scaling}$  (e.g., log of
266     the softmax
267     normalization factor).
268     """
269     assert k_cache.stride(-1) == 1, "k_cache must have contiguous last
270     dimension"
271     assert v_cache.stride(-1) == 1, "v_cache must have contiguous last
272     dimension"
273     q, k, v = [maybe_contiguous(x) for x in (q, k, v)]
274     if softmax_scale is None:
275         softmax_scale = q.shape[-1] ** (-0.5)
276     if cache_seqLens is not None and isinstance(cache_seqLens, int):
277         cache_seqLens = torch.full(
278             (k_cache.shape[0],), cache_seqLens, dtype=torch.int32,
279             device=k_cache.device
280         )
281     cache_seqLens = maybe_contiguous(cache_seqLens)
282     cache_batch_idx = maybe_contiguous(cache_batch_idx)
283     block_table = maybe_contiguous(block_table)
284     if fa_version == 2:
285         if scheduler_metadata is not None and q_descale is not None \
286             and k_descale is not None and v_descale is not None:
287             raise NotImplementedError(
288                 "FA2 does not support scheduler_metadata, q_descale, "
289                 "k_descale, v_descale"
290             )
291     out, softmax_lse = torch.ops._vllm_fa2_C.fwd_kvcache(
292         q, k_cache, v_cache,
293         k, v,          # k_new, v_new
294         cache_seqLens,
295         rotary_cos,
296         rotary_sin,
297         cache_batch_idx,
298         cache_leftpad,
299         block_table,
300         alibi_slopes,
301         out,
302         softmax_scale,
303         causal,
304         window_size[0],
305         window_size[1],
306         softcap,
307         rotary_interleaved,
308         num_splits,

```

```

305     )
306     elif fa_version == 3:
307         assert alibi_slopes is None, "Alibi is not supported in FA3"
308         out, softmax_lse, _, _ = torch.ops.vllm_fa3_C.fwd(
309             q, k_cache, v_cache, # q, k, v
310             k, v,                # k_new, v_new
311             None,                # q_v
312             out,
313             None, None,          # cu_seqlens_q, cu_seqlens_k
314             None,                # cu_seqlens_k_new
315             None, cache_seqlens, # seqused_q, seqused_k
316             None, None,         # max_seqlen_q, max_seqlen_k
317             block_table,
318             cache_batch_idx,    # kv_batch_idx
319             None,               # leftpad_k
320             None, None,        # rotary_cos, rotary_sin
321             q_descale, k_descale, v_descale,
322             softmax_scale,
323             causal,
324             window_size[0], window_size[1],
325             softcap,
326             rotary_interleaved, # rotary_interleaved
327             scheduler_metadata,
328             num_splits,         # num_splits
329             None,               # pack_gqa
330             0,                  # sm_margin
331         )
332     else:
333         raise ValueError(f"Unsupported FA version: {fa_version}")
334     return (out, softmax_lse) if return_softmax_lse else out

```

FA2 定义在: https://github.com/vllm-project/flash-attention/blob/d637d8927a35922ce6f6c0dff6dd3f765ed71f3c/flash_attn/flash_attn_interface.py#L1375

hopper 架构的 FA3 定义位置不同: https://github.com/vllm-project/flash-attention/blob/d637d8927a35922ce6f6c0dff6dd3f765ed71f3c/hopper/flash_attn_interface.py#L541

参考实现: 虽然使用的 kernel 不同, 但是 FA2 和 FA3 的实现都是通过模板实例化

FA2: https://github.com/vllm-project/flash-attention/tree/main/csrc/flash_attn/src_flash_fwd_launch_template.h [flash_fwd_kernel.h](https://github.com/vllm-project/flash-attention/tree/main/csrc/flash_attn/src_flash_fwd_kernel.h)

FA3: https://github.com/vllm-project/flash-attention/blob/main/hopper/flash_fwd_launch_template.h https://github.com/vllm-project/flash-attention/blob/main/hopper/flash_api.cpp <https://github.com/vllm-project/flash-attention/tree/main/hopper/instantiations>

FA2接口可以参考

https://github.com/vllm-project/flash-attention/blob/main/csrc/flash_attn/flash_api.cpp

csrc/flash_attn/flash_api.cpp

```
1  PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
2      m.doc() = "FlashAttention";
3      m.def("fwd", &FLASH_NAMESPACE::mha_fwd, "Forward pass");
4      m.def("varlen_fwd", &FLASH_NAMESPACE::mha_varlen_fwd, "Forward pass
(variable length)");
5      m.def("bwd", &FLASH_NAMESPACE::mha_bwd, "Backward pass");
6      m.def("varlen_bwd", &FLASH_NAMESPACE::mha_varlen_bwd, "Backward pass
(variable length)");
7      m.def("fwd_kvcache", &FLASH_NAMESPACE::mha_fwd_kvcache, "Forward pass,
with KV-cache");
8  }
```

FA3接口可以参考

https://github.com/vllm-project/flash-attention/blob/main/hopper/flash_api.cpp

代码块

```
1  PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
2      m.doc() = "FlashAttention";
3      m.def("fwd", &mha_fwd, "Forward pass");
4      m.def("bwd", &mha_bwd, "Backward pass");
5      m.def("fwd_combine", &mha_combine, "Combine partial attention outputs");
6      m.def("get_scheduler_metadata", &mha_fwd_get_scheduler_metadata, "Get
scheduler metadata for varlen forward pass");
7  }
```

FA的forward方法均是通过一个结构体Flash_fwd_params来传递参数，调用对应的kernel，但是FA2和FA3的Flash_fwd_params结构体定义是不同的

FA2 https://github.com/vllm-project/flash-attention/blob/d637d8927a35922ce6f6c0dff6dd3f765ed71f3c/csrc/flash_attn/src/flash.h#L48

FA3 <https://github.com/vllm-project/flash-attention/blob/d637d8927a35922ce6f6c0dff6dd3f765ed71f3c/hopper/flash.h#L37>

如果接口支持结构体传递，可以直接使用上面链接中的结构体

Function 接口

FA2

Without kvcache

代码块

```
1  /**
2   * @brief Executes the Flash Attention forward pass with fully explicit
   parameters
3   *
4   * @param dev          Device handle for CUDA/CPU device management
5   *                    Used to identify the device where tensors are
   allocated and computation will occur
6   *
7   * @param q_ptr       Pointer to query tensor data (total_q, nheads,
   headdim) or (batch_size, seqlen_q, nheads, headdim)
8   *                    For variable length sequences, total_q = sum of all
   sequence lengths in batch
9   *
10  * @param k_ptr       Pointer to key tensor data (total_k, nheads_k,
   headdim) or (batch_size, seqlen_k, nheads_k, headdim)
11  *                    For KV cache, this points to the cached keys
12  *
13  * @param v_ptr       Pointer to value tensor data (total_k, nheads_k,
   headdim) or (batch_size, seqlen_k, nheads_k, headdim)
14  *                    For KV cache, this points to the cached values
15  *
16  * @param o_ptr       Pointer to output tensor data (total_q, nheads,
   headdim) or (batch_size, seqlen_q, nheads, headdim)
17  *                    Will store the result of attention computation
18  *
19  * @param softmax_lse_ptr Pointer to store log-sum-exp values from softmax
   computation (nheads, total_q) or (batch_size, nheads, seqlen_q)
20  *                    Used for numerical stability and backward pass. Set
   to nullptr if not needed.
21  *
22  * @param p_ptr       Pointer to store attention matrix (batch_size,
   nheads, seqlen_q, seqlen_k)
23  *                    Primarily used for debugging or visualization. Set
   to nullptr if not needed.
```

24 *
25 * *@param batch_size* *Number of sequences in the batch*
26 * *For fixed length: all sequences have the same length*
27 * *For variable length: sequences can have different*
lengths
28 *
29 * *@param seqlen_q* *Length of each query sequence (for fixed length) or*
maximum query sequence length (for variable length)
30 * *In decoder scenarios, typically 1 for*
autoregressive generation
31 *
32 * *@param seqlen_k* *Length of each key/value sequence (for fixed*
length) or maximum key/value sequence length (for variable length)
33 * *In KV cache scenarios, this is the current length*
of the cached context
34 *
35 * *@param num_heads* *Number of attention heads for query tensor*
36 * *Each head performs independent attention computation*
37 *
38 * *@param num_heads_k* *Number of attention heads for key/value tensors*
39 * *For standard attention: equals num_heads*
40 * *For multi-query attention (MQA): equals 1*
41 * *For grouped-query attention (GQA): equals num_heads*
/ group_size
42 *
43 * *@param head_size* *Dimension of each attention head*
44 * *Typically 64, 80, or 128 in modern architectures*
45 * *Must be ≤ 256 for Flash Attention implementation*
46 *
47 * *@param q_batch_stride* *Stride (in elements) between consecutive batches in*
query tensor
48 * *For contiguous tensor: $seqlen_q * num_heads * head_size$*
head_size
49 *
50 * *@param q_row_stride* *Stride (in elements) between consecutive rows*
(sequences) in query tensor
51 * *For contiguous tensor: $num_heads * head_size$*
52 *
53 * *@param q_head_stride* *Stride (in elements) between consecutive heads in*
query tensor
54 * *For contiguous tensor: head_size*
55 *
56 * *@param k_batch_stride* *Stride (in elements) between consecutive batches in*
key tensor
57 * *For contiguous tensor: $seqlen_k * num_heads_k * head_size$*
head_size
58 *

59 * @param k_row_stride Stride (in elements) between consecutive rows
(sequences) in key tensor

60 * For contiguous tensor: num_heads_k * head_size

61 *

62 * @param k_head_stride Stride (in elements) between consecutive heads in
key tensor

63 * For contiguous tensor: head_size

64 *

65 * @param v_batch_stride Stride (in elements) between consecutive batches in
value tensor

66 * For contiguous tensor: seqlen_k * num_heads_k *
head_size

67 *

68 * @param v_row_stride Stride (in elements) between consecutive rows
(sequences) in value tensor

69 * For contiguous tensor: num_heads_k * head_size

70 *

71 * @param v_head_stride Stride (in elements) between consecutive heads in
value tensor

72 * For contiguous tensor: head_size

73 *

74 * @param o_batch_stride Stride (in elements) between consecutive batches in
output tensor

75 * For contiguous tensor: seqlen_q * num_heads *
head_size

76 *

77 * @param o_row_stride Stride (in elements) between consecutive rows
(sequences) in output tensor

78 * For contiguous tensor: num_heads * head_size

79 *

80 * @param o_head_stride Stride (in elements) between consecutive heads in
output tensor

81 * For contiguous tensor: head_size

82 *

83 * @param softmax_scale Scale factor applied to attention scores before
softmax

84 * Typically $1/\sqrt{\text{head_size}}$ to normalize variance

85 * Affects the sharpness of attention distribution

86 *

87 * @param p_dropout Dropout probability (0.0 to 1.0)

88 * Fraction of attention connections to randomly zero
out during training

89 * Set to 0.0 for inference or to disable dropout

90 *

91 * @param softcap Softcapping value for attention scores (0.0 to
disable)

92 * *If > 0.0 , applies a soft cap to attention scores to prevent extreme values*

93 * *Helps with training stability for very long sequences*

94 *

95 * *@param is_causal Whether to apply causal attention mask*

96 * *If true, each position can only attend to previous positions (including itself)*

97 * *Used in autoregressive models to prevent looking at future tokens*

98 *

99 * *@param window_size_left Left window size for local attention (-1 for unlimited)*

100 * *If ≥ 0 , limits how far back each query can attend*

101 * *Useful for long sequence efficiency or implementing sliding window attention*

102 *

103 * *@param window_size_right Right window size for local attention (-1 for unlimited)*

104 * *If ≥ 0 , limits how far forward each query can attend*

105 * *Only relevant when is_causal=false*

106 *

107 * *@param cu_seqlens_q Cumulative sequence lengths for queries (batch_size+1)*

108 * *Format: [0, len1, len1+len2, len1+len2+len3, ...]*

109 * *Used for variable length sequences. Set to nullptr for fixed length.*

110 *

111 * *@param cu_seqlens_k Cumulative sequence lengths for keys/values (batch_size+1)*

112 * *Format: [0, len1, len1+len2, len1+len2+len3, ...]*

113 * *Used for variable length sequences. Set to nullptr for fixed length.*

114 *

115 * *@param seqused_k Actual sequence usage for keys (batch_size)*

116 * *Specifies how much of each sequence in the KV cache is actually used*

117 * *Useful when allocated length exceeds used length. Set to nullptr if not needed.*

118 *

119 * *@param leftpad_k Left padding for each sequence in key (batch_size)*

120 * *Specifies number of padding tokens at the beginning of each sequence*

121 * *Affects position calculations for rotary embeddings. Set to nullptr if no left padding.*

122 *

123 * @param *is_seqlens_k_cumulative* Whether sequence lengths in *cu_seqlens_k* are
cumulative

124 * *If true, cu_seqlens_k contains cumulative*
lengths

125 * *If false, cu_seqlens_k contains actual*
sequence lengths

126 *

127 * @param *knew_ptr* New key tensor data pointer (*batch_size*,
seqlen_knew, *nheads_k*, *headdim*)

128 * *Contains new keys to be appended to the KV cache*

129 * *Set to nullptr if not appending to cache*

130 *

131 * @param *vnew_ptr* New value tensor data pointer (*batch_size*,
seqlen_knew, *nheads_k*, *headdim*)

132 * *Contains new values to be appended to the KV cache*

133 * *Set to nullptr if not appending to cache*

134 *

135 * @param *seqlen_knew* Length of new keys/values to append to the cache

136 * *Typically 1 for autoregressive generation, or chunk*
size for batch processing

137 *

138 * @param *knew_batch_stride* Stride between batches in new key tensor

139 * *For contiguous tensor: seqlen_knew * num_heads_k **
head_size

140 *

141 * @param *knew_row_stride* Stride between rows in new key tensor

142 * *For contiguous tensor: num_heads_k * head_size*

143 *

144 * @param *knew_head_stride* Stride between heads in new key tensor

145 * *For contiguous tensor: head_size*

146 *

147 * @param *vnew_batch_stride* Stride between batches in new value tensor

148 * *For contiguous tensor: seqlen_knew * num_heads_k **
head_size

149 *

150 * @param *vnew_row_stride* Stride between rows in new value tensor

151 * *For contiguous tensor: num_heads_k * head_size*

152 *

153 * @param *vnew_head_stride* Stride between heads in new value tensor

154 * *For contiguous tensor: head_size*

155 *

156 * @param *rotary_cos_ptr* Cosine table for rotary embeddings (*seqlen_ro*,
rotary_dim/2)

157 * *Contains precomputed cosine values for rotary*
position embeddings

158 * *seqlen_ro should be >= maximum sequence length.*
Set to nullptr to disable.

159 *
160 * `@param rotary_sin_ptr` Sine table for rotary embeddings (`seqlen_ro`,
`rotary_dim/2`)
161 *
Contains precomputed sine values for rotary
position embeddings
162 *
`seqlen_ro` should be \geq maximum sequence length.
Set to `nullptr` to disable.
163 *
164 * `@param rotary_dim` Dimension for rotary embeddings (must be \leq
`head_size` and divisible by 16)
165 *
Only the first `rotary_dim` dimensions of each head
will have rotary embeddings
166 *
Set to 0 if not using rotary embeddings
167 *
168 * `@param is_rotary_interleaved` Whether rotary embeddings are interleaved
169 *
If true, dimensions are paired as (0,1),
(2,3), etc. (GPT-J style)
170 *
If false, dimensions are paired as
(0, `rotary_dim/2`), (1, `rotary_dim/2+1`), etc. (GPT-NeoX style)
171 *
172 * `@param cache_batch_idx` Indices for batch mapping (`batch_size`)
173 *
Maps from logical batch indices to physical
indices in the KV cache
174 *
Useful for reordering or sharing KV cache across
requests. Set to `nullptr` for direct mapping.
175 *
176 * `@param block_table` Block table for paged KV cache (`batch_size`,
`max_blocks_per_seq`)
177 *
Contains block indices for each sequence in paged
KV cache
178 *
Each row maps logical sequence positions to
physical memory blocks. Set to `nullptr` for contiguous KV cache.
179 *
180 * `@param block_table_batch_stride` Stride between batches in block table
181 *
For contiguous tensor: `max_blocks_per_seq`
182 *
183 * `@param page_block_size` Size of each block in paged KV cache
184 *
Typically 16 or 32 tokens per block
185 *
Only relevant when using paged KV cache
(`block_table != nullptr`)
186 *
187 * `@param alibi_slopes_ptr` ALiBi position bias slopes (`nheads`) or
(`batch_size`, `nheads`)
188 *
Contains slope values for ALiBi positional encoding
189 *
Different slopes per head create different
attention patterns. Set to `nullptr` to disable ALiBi.
190 *

```

191 * @param alibi_slopes_batch_stride Stride between batches in ALiBi slopes
192 *
193 * For contiguous tensor: nheads
194 * Set to 0 if alibi_slopes is shared across
195 * batch (shape is just (nheads))
196 *
197 * @param num_splits Number of splits for computation (0 for auto)
198 *
199 * Splits the key/value tensors into chunks along
200 * sequence dimension
201 * Helps with memory efficiency for long sequences. 0
202 * means auto-select based on GPU.
203 *
204 * @param oaccum_ptr Accumulated output pointer for split computation
205 * (total_q, nheads, headdim)
206 *
207 * Used internally for split-KV optimization
208 * Set to nullptr unless manually implementing split
209 * computation
210 *
211 * @param softmax_lseaccum_ptr Accumulated LSE pointer for split computation
212 * (nheads, total_q)
213 *
214 * Used internally for split-KV optimization
215 * Set to nullptr unless manually implementing
216 * split computation
217 *
218 * @param is_bf16 Whether tensors are in BF16 format (true) or FP16
219 * format (false)
220 *
221 * BF16 has less precision but more range than FP16
222 * Must accurately reflect the data type of input
223 * tensors
224 *
225 * @param unpadded_lse Whether LSE is in unpadded format
226 *
227 * If true, softmax_lse has shape [nheads, total_q]
228 * for variable length sequences
229 *
230 * If false, softmax_lse has shape [batch_size,
231 * nheads, seqlen_q]
232 *
233 * @param seqlenq_ngroups_swapped Whether q has been transposed for
234 * optimization
235 *
236 * If true, q has been transposed from (b, 1,
237 * (nheads_k, ngroups), d) to (b, ngroups, nheads_k, d)
238 *
239 * Used for specialized optimizations in
240 * certain model architectures
241 *
242 * @param stream CUDA stream to execute the kernel on
243 *
244 * Allows asynchronous execution and integration with
245 * other CUDA operations
246 *
247 * Set to nullptr to use the current CUDA stream
248 *
249 *

```

```

223  * @param force_split_kernel Whether to force using the split kernel
implementation
224  *                                If true, always uses the split-KV optimization
regardless of sequence length
225  *                                Required for KV cache, paged KV cache, or when
using cache_batch_idx
226  *
227  * @param return_softmax_lse Whether to return the log-sum-exp values
228  *                                If true, softmax_lse_ptr must be a valid pointer
229  *
230  * @return Status code (0 for success, non-zero for failure)
231  *     Specific error codes depend on the implementation
232  */
233  int run_mha_fwd(
234      void* dev, // Device handle
235      // Core tensor pointers
236      void* q_ptr,
237      void* k_ptr,
238      void* v_ptr,
239      void* o_ptr,
240      void* softmax_lse_ptr = nullptr,
241      void* p_ptr = nullptr,
242
243      // Dimensions
244      int batch_size = 1,
245      int seqlen_q = 0,
246      int seqlen_k = 0,
247      int num_heads = 0,
248      int num_heads_k = 0,
249      int head_size = 0,
250
251      // Stride information (in elements, not bytes)
252      int64_t q_batch_stride = 0,
253      int64_t q_row_stride = 0,
254      int64_t q_head_stride = 0,
255      int64_t k_batch_stride = 0,
256      int64_t k_row_stride = 0,
257      int64_t k_head_stride = 0,
258      int64_t v_batch_stride = 0,
259      int64_t v_row_stride = 0,
260      int64_t v_head_stride = 0,
261      int64_t o_batch_stride = 0,
262      int64_t o_row_stride = 0,
263      int64_t o_head_stride = 0,
264
265      // Scaling parameters
266      float softmax_scale = 0.0f,

```

```
267 float p_dropout = 0.0f,
268 float softcap = 0.0f,
269
270 // Masking parameters
271 bool is_causal = false,
272 int window_size_left = -1,
273 int window_size_right = -1,
274
275 // Variable length parameters
276 int* cu_seqlens_q = nullptr,
277 int* cu_seqlens_k = nullptr,
278 int* seqused_k = nullptr,
279 int* leftpad_k = nullptr,
280 bool is_seqlens_k_cumulative = true,
281
282 // KV cache parameters
283 void* knew_ptr = nullptr,
284 void* vnew_ptr = nullptr,
285 int seqlen_knew = 0,
286 int64_t knew_batch_stride = 0,
287 int64_t knew_row_stride = 0,
288 int64_t knew_head_stride = 0,
289 int64_t vnew_batch_stride = 0,
290 int64_t vnew_row_stride = 0,
291 int64_t vnew_head_stride = 0,
292
293 // Rotary embedding parameters
294 void* rotary_cos_ptr = nullptr,
295 void* rotary_sin_ptr = nullptr,
296 int rotary_dim = 0,
297 bool is_rotary_interleaved = true,
298
299 // Cache indexing
300 int* cache_batch_idx = nullptr,
301
302 // Paged KV cache parameters
303 int* block_table = nullptr,
304 int64_t block_table_batch_stride = 0,
305 int page_block_size = 0,
306
307 // ALiBi parameters
308 void* alibi_slopes_ptr = nullptr,
309 int64_t alibi_slopes_batch_stride = 0,
310
311 // Split computation parameters
312 int num_splits = 0,
313 void* oaccum_ptr = nullptr,
```

```

314     void* softmax_lseaccum_ptr = nullptr,
315
316     // Miscellaneous flags
317     bool is_bf16 = false,
318     bool unpadded_lse = false,
319     bool seqlenq_ngroups_swapped = false,
320
321     // Execution parameters
322     cudaStream_t stream = nullptr,
323     bool force_split_kernel = false
324     // Return options
325     bool return_softmax_lse = false
326 );

```

With kvcache

代码块

```

1  /**
2   * @brief Executes Flash Attention forward pass with KV cache
3   *
4   * @param dev Device handle for CUDA/CPU device management
5   *          Used to identify the device where tensors are
6   *          allocated and computation will occur
7   *
8   * @param q_ptr Pointer to query tensor data (batch_size, seqlen_q,
9   *          nheads, headdim)
10  *           Current query tokens to compute attention for
11  *           For decoder scenarios, typically has seqlen_q=1 for
12  *           autoregressive generation
13  *
14  * @param k_cache_ptr Pointer to key cache tensor data
15  *                   For standard KV cache: (batch_size_cache,
16  *                   seqlen_cache, nheads_k, headdim)
17  *                   For paged KV cache: (num_blocks, block_size,
18  *                   nheads_k, headdim)
19  *                   Contains previously processed keys
20  *
21  * @param v_cache_ptr Pointer to value cache tensor data
22  *                   For standard KV cache: (batch_size_cache,
23  *                   seqlen_cache, nheads_k, headdim)
24  *                   For paged KV cache: (num_blocks, block_size,
25  *                   nheads_k, headdim)
26  *                   Contains previously processed values
27  *
28  * @param o_ptr Pointer to output tensor data (batch_size,
29  *          seqlen_q, nheads, headdim)

```

| | | |
|----|-------------------------------|---|
| 22 | * | <i>Will store the result of attention computation</i> |
| 23 | * | |
| 24 | * @param softmax_lse_ptr | <i>Pointer to store log-sum-exp values (batch_size, nheads, seqlen_q)</i> |
| 25 | * | <i>Used for numerical stability. Set to nullptr if not needed.</i> |
| 26 | * | |
| 27 | * @param batch_size | <i>Number of sequences in the batch</i> |
| 28 | * | <i>For inference scenarios, often set to 1</i> |
| 29 | * | |
| 30 | * @param seqlen_q | <i>Length of each query sequence</i> |
| 31 | * | <i>For autoregressive generation, typically 1</i> |
| 32 | * | |
| 33 | * @param seqlen_cache | <i>Maximum sequence length in the KV cache</i> |
| 34 | * | <i>Represents the capacity of the cache, not necessarily fully used</i> |
| 35 | * | |
| 36 | * @param num_heads | <i>Number of attention heads for query tensor</i> |
| 37 | * | <i>Each head performs independent attention computation</i> |
| 38 | * | |
| 39 | * @param num_heads_k | <i>Number of attention heads for key/value tensors</i> |
| 40 | * | <i>For standard attention: equals num_heads</i> |
| 41 | * | <i>For multi-query attention (MQA): equals 1</i> |
| 42 | * | <i>For grouped-query attention (GQA): equals num_heads / group_size</i> |
| 43 | * | |
| 44 | * @param head_size | <i>Dimension of each attention head</i> |
| 45 | * | <i>Typically 64, 80, or 128 in modern architectures</i> |
| 46 | * | <i>Must be ≤ 256 for Flash Attention implementation</i> |
| 47 | * | |
| 48 | * @param q_batch_stride | <i>Stride (in elements) between consecutive batches in query tensor</i> |
| 49 | * | <i>For contiguous tensor: $seqlen_q * num_heads * head_size$</i> |
| 50 | * | |
| 51 | * @param q_row_stride | <i>Stride (in elements) between consecutive rows in query tensor</i> |
| 52 | * | <i>For contiguous tensor: $num_heads * head_size$</i> |
| 53 | * | |
| 54 | * @param q_head_stride | <i>Stride (in elements) between consecutive heads in query tensor</i> |
| 55 | * | <i>For contiguous tensor: head_size</i> |
| 56 | * | |
| 57 | * @param k_cache_batch_stride | <i>Stride between batches in key cache tensor</i> |
| 58 | * | <i>For contiguous tensor: $seqlen_cache * num_heads_k * head_size$</i> |
| 59 | * | <i>Ignored for paged KV cache</i> |

60 *
61 * **@param** *k_cache_row_stride* *Stride between rows in key cache tensor*
62 * *For contiguous tensor: num_heads_k * head_size*
63 * *For paged KV cache: block_size * num_heads_k * head_size*
64 *
65 * **@param** *k_cache_head_stride* *Stride between heads in key cache tensor*
66 * *For contiguous tensor: head_size*
67 *
68 * **@param** *v_cache_batch_stride* *Stride between batches in value cache tensor*
69 * *For contiguous tensor: seqlen_cache * num_heads_k * head_size*
70 * *Ignored for paged KV cache*
71 *
72 * **@param** *v_cache_row_stride* *Stride between rows in value cache tensor*
73 * *For contiguous tensor: num_heads_k * head_size*
74 * *For paged KV cache: block_size * num_heads_k * head_size*
75 *
76 * **@param** *v_cache_head_stride* *Stride between heads in value cache tensor*
77 * *For contiguous tensor: head_size*
78 *
79 * **@param** *o_batch_stride* *Stride between batches in output tensor*
80 * *For contiguous tensor: seqlen_q * num_heads * head_size*
81 *
82 * **@param** *o_row_stride* *Stride between rows in output tensor*
83 * *For contiguous tensor: num_heads * head_size*
84 *
85 * **@param** *o_head_stride* *Stride between heads in output tensor*
86 * *For contiguous tensor: head_size*
87 *
88 * **@param** *softmax_scale* *Scale factor applied to attention scores before softmax*
89 * *Typically 1/sqrt(head_size) to normalize variance*
90 * *Set to 0.0 to use default scale of 1/sqrt(head_size)*
91 *
92 * **@param** *softcap* *Softcapping value for attention scores (0.0 to disable)*
93 * *If > 0.0, applies a soft cap to attention scores to prevent extreme values*
94 *
95 * **@param** *is_causal* *Whether to apply causal attention mask*
96 * *If true, each position can only attend to previous positions (including itself)*
97 * *Used in autoregressive models to prevent looking at future tokens*

98 *
99 * **@param** *window_size_left* Left window size for local attention (-1 for unlimited)
100 * If ≥ 0 , limits how far back each query can attend
101 *
102 * **@param** *window_size_right* Right window size for local attention (-1 for unlimited)
103 * If ≥ 0 , limits how far forward each query can attend
104 * Only relevant when *is_causal*=false
105 *
106 * **@param** *k_new_ptr* New key tensor data pointer (batch_size, seq_len_knew, nheads_k, headdim)
107 * Contains new keys to be appended to the KV cache
108 * Set to nullptr if not appending to cache
109 *
110 * **@param** *v_new_ptr* New value tensor data pointer (batch_size, seq_len_knew, nheads_k, headdim)
111 * Contains new values to be appended to the KV cache
112 * Set to nullptr if not appending to cache
113 *
114 * **@param** *seq_len_knew* Length of new keys/values to append to the cache
115 * Typically 1 for autoregressive generation
116 *
117 * **@param** *k_new_batch_stride* Stride between batches in new key tensor
118 * For contiguous tensor: $seq_len_knew * num_heads_k$
119 * *head_size*
120 * **@param** *k_new_row_stride* Stride between rows in new key tensor
121 * For contiguous tensor: $num_heads_k * head_size$
122 *
123 * **@param** *k_new_head_stride* Stride between heads in new key tensor
124 * For contiguous tensor: *head_size*
125 *
126 * **@param** *v_new_batch_stride* Stride between batches in new value tensor
127 * For contiguous tensor: $seq_len_knew * num_heads_k$
128 * *head_size*
129 * **@param** *v_new_row_stride* Stride between rows in new value tensor
130 * For contiguous tensor: $num_heads_k * head_size$
131 *
132 * **@param** *v_new_head_stride* Stride between heads in new value tensor
133 * For contiguous tensor: *head_size*
134 *
135 * **@param** *cache_seqLens* Current sequence lengths in KV cache (batch_size)
136 * Specifies the current length of each sequence in the cache

137 * *Required when appending new keys/values to the cache*

138 *

139 * `@param rotary_cos_ptr` *Cosine table for rotary embeddings (seqLen_ro, rotary_dim/2)*

140 * *Contains precomputed cosine values for rotary position embeddings*

141 * *seqLen_ro should be \geq maximum sequence length. Set to nullptr to disable.*

142 *

143 * `@param rotary_sin_ptr` *Sine table for rotary embeddings (seqLen_ro, rotary_dim/2)*

144 * *Contains precomputed sine values for rotary position embeddings*

145 * *seqLen_ro should be \geq maximum sequence length. Set to nullptr to disable.*

146 *

147 * `@param rotary_dim` *Dimension for rotary embeddings (must be \leq head_size and divisible by 16)*

148 * *Only the first rotary_dim dimensions of each head will have rotary embeddings*

149 * *Set to 0 if not using rotary embeddings*

150 *

151 * `@param is_rotary_interleaved` *Whether rotary embeddings are interleaved*

152 * *If true, dimensions are paired as (0,1), (2,3), etc. (GPT-J style)*

153 * *If false, dimensions are paired as (0, rotary_dim/2), (1, rotary_dim/2+1), etc. (GPT-NeoX style)*

154 *

155 * `@param cache_batch_idx` *Indices for batch mapping (batch_size)*

156 * *Maps from logical batch indices to physical indices in the KV cache*

157 * *Useful for reordering or sharing KV cache across requests. Set to nullptr for direct mapping.*

158 *

159 * `@param cache_leftpad` *Left padding for each sequence in cache (batch_size)*

160 * *Specifies number of padding tokens at the beginning of each sequence*

161 * *Affects position calculations for rotary embeddings. Set to nullptr if no left padding.*

162 *

163 * `@param block_table` *Block table for paged KV cache (batch_size, max_blocks_per_seq)*

164 * *Contains block indices for each sequence in paged KV cache*

165 * *Each row maps logical sequence positions to physical memory blocks. Set to nullptr for contiguous KV cache.*

166 *

167 * **@param** *block_table_batch_stride* Stride between batches in block table
168 * *For contiguous tensor: max_blocks_per_seq*
169 *
170 * **@param** *page_block_size* Size of each block in paged KV cache
171 * *Typically 16 or 32 tokens per block*
172 * *Only relevant when using paged KV cache*
(block_table != nullptr)
173 *
174 * **@param** *alibi_slopes_ptr* ALiBi position bias slopes (nheads) or (batch_size, nheads)
175 * *Contains slope values for ALiBi positional encoding*
176 * *Different slopes per head create different attention patterns. Set to nullptr to disable ALiBi.*
177 *
178 * **@param** *alibi_slopes_batch_stride* Stride between batches in ALiBi slopes
179 * *For contiguous tensor: nheads*
180 * *Set to 0 if alibi_slopes is shared across batch (shape is just (nheads))*
181 *
182 * **@param** *num_splits* Number of splits for computation (0 for auto)
183 * *Splits the key/value tensors into chunks along sequence dimension*
184 * *Helps with memory efficiency for long sequences. 0 means auto-select based on GPU.*
185 *
186 * **@param** *oaccum_ptr* Accumulated output pointer for split computation
187 * *Used internally for split-KV optimization*
188 * *Set to nullptr unless manually implementing split computation*
189 *
190 * **@param** *softmax_lseaccum_ptr* Accumulated LSE pointer for split computation
191 * *Used internally for split-KV optimization*
192 * *Set to nullptr unless manually implementing split computation*
193 *
194 * **@param** *is_bf16* Whether tensors are in BF16 format (true) or FP16 format (false)
195 * *BF16 has less precision but more range than FP16*
196 * *Must accurately reflect the data type of input tensors*
197 *
198 * **@param** *stream* CUDA stream to execute the kernel on
199 * *Allows asynchronous execution and integration with other CUDA operations*
200 * *Set to nullptr to use the current CUDA stream*
201 *
202 * **@param** *return_softmax_lse* Whether to return the log-sum-exp values

```

203  *                                     If true, softmax_lse_ptr must be a valid pointer
204  *
205  * @return Status code (0 for success, non-zero for failure)
206  */
207  int run_mha_fwd_kvcache(
208      void* dev,                               // Device handle
209
210      // Core tensor pointers
211      void* q_ptr,
212      void* k_cache_ptr,
213      void* v_cache_ptr,
214      void* o_ptr,
215      void* softmax_lse_ptr = nullptr,
216
217      // Dimensions
218      int batch_size = 1,
219      int seqlen_q = 1,
220      int seqlen_cache = 0,
221      int num_heads = 0,
222      int num_heads_k = 0,
223      int head_size = 0,
224
225      // Stride information (in elements, not bytes)
226      int64_t q_batch_stride = 0,
227      int64_t q_row_stride = 0,
228      int64_t q_head_stride = 0,
229      int64_t k_cache_batch_stride = 0,
230      int64_t k_cache_row_stride = 0,
231      int64_t k_cache_head_stride = 0,
232      int64_t v_cache_batch_stride = 0,
233      int64_t v_cache_row_stride = 0,
234      int64_t v_cache_head_stride = 0,
235      int64_t o_batch_stride = 0,
236      int64_t o_row_stride = 0,
237      int64_t o_head_stride = 0,
238
239      // Scaling parameters
240      float softmax_scale = 0.0f,
241      float softcap = 0.0f,
242
243      // Masking parameters
244      bool is_causal = true,
245      int window_size_left = -1,
246      int window_size_right = -1,
247
248      // New KV parameters
249      void* k_new_ptr = nullptr,

```

```
250 void* v_new_ptr = nullptr,
251 int seqlen_knew = 0,
252 int64_t k_new_batch_stride = 0,
253 int64_t k_new_row_stride = 0,
254 int64_t k_new_head_stride = 0,
255 int64_t v_new_batch_stride = 0,
256 int64_t v_new_row_stride = 0,
257 int64_t v_new_head_stride = 0,
258
259 // KV cache parameters
260 int* cache_seqLens = nullptr,
261
262 // Rotary embedding parameters
263 void* rotary_cos_ptr = nullptr,
264 void* rotary_sin_ptr = nullptr,
265 int rotary_dim = 0,
266 bool is_rotary_interleaved = true,
267
268 // Cache indexing
269 int* cache_batch_idx = nullptr,
270 int* cache_leftpad = nullptr,
271
272 // Paged KV cache parameters
273 int* block_table = nullptr,
274 int64_t block_table_batch_stride = 0,
275 int page_block_size = 0,
276
277 // ALiBi parameters
278 void* alibi_slopes_ptr = nullptr,
279 int64_t alibi_slopes_batch_stride = 0,
280
281 // Split computation parameters
282 int num_splits = 0,
283 void* oaccum_ptr = nullptr,
284 void* softmax_lseaccum_ptr = nullptr,
285
286 // Data type flags
287 bool is_bf16 = false,
288
289 // Execution parameters
290 cudaStream_t stream = nullptr,
291
292 // Return options
293 bool return_softmax_lse = false
294 );
```

Without kvcache

代码块

```

1  /**
2   * @brief Executes the Flash Attention 3 forward pass with fully explicit
   parameters
3   *
4   * @param dev           Device handle for CUDA/CPU device management
5   *                       Used to identify the device where tensors are
   allocated and computation will occur
6   *
7   * @param q_ptr         Pointer to query tensor data (total_q, nheads,
   headdim) or (batch_size, seqlen_q, nheads, headdim)
8   *                       For variable length sequences, total_q = sum of all
   sequence lengths in batch
9   *
10  * @param k_ptr         Pointer to key tensor data (total_k, nheads_k,
   headdim) or (batch_size, seqlen_k, nheads_k, headdim)
11  *                       For KV cache, this points to the cached keys
12  *
13  * @param v_ptr         Pointer to value tensor data (total_k, nheads_k,
   headdim_v) or (batch_size, seqlen_k, nheads_k, headdim_v)
14  *                       For KV cache, this points to the cached values
15  *                       Note: headdim_v can be different from headdim in FA3
16  *
17  * @param o_ptr         Pointer to output tensor data (total_q, nheads,
   headdim_v) or (batch_size, seqlen_q, nheads, headdim_v)
18  *                       Will store the result of attention computation
19  *
20  * @param softmax_lse_ptr Pointer to store log-sum-exp values from softmax
   computation (nheads, total_q) or (batch_size, nheads, seqlen_q)
21  *                       Used for numerical stability and backward pass. Set
   to nullptr if not needed.
22  *
23  * @param batch_size    Number of sequences in the batch
24  *                       For fixed length: all sequences have the same length
25  *                       For variable length: sequences can have different
   lengths
26  *
27  * @param seqlen_q      Length of each query sequence (for fixed length) or
   maximum query sequence length (for variable length)
28  *                       In decoder scenarios, typically 1 for
   autoregressive generation
29  *

```

| | | |
|----|--------------------------------------|---|
| 30 | * <code>@param seqLen_k</code> | Length of each key/value sequence (for fixed length) or maximum key/value sequence length (for variable length) |
| 31 | * | In KV cache scenarios, this is the current length of the cached context |
| 32 | * | |
| 33 | * <code>@param num_heads</code> | Number of attention heads for query tensor |
| 34 | * | Each head performs independent attention computation |
| 35 | * | |
| 36 | * <code>@param num_heads_k</code> | Number of attention heads for key/value tensors |
| 37 | * | For standard attention: equals <code>num_heads</code> |
| 38 | * | For multi-query attention (MQA): equals 1 |
| 39 | * | For grouped-query attention (GQA): equals <code>num_heads / group_size</code> |
| 40 | * | |
| 41 | * <code>@param head_size</code> | Dimension of each attention head for query/key |
| 42 | * | Typically 64, 80, or 128 in modern architectures |
| 43 | * | Must be ≤ 256 for Flash Attention implementation |
| 44 | * | |
| 45 | * <code>@param head_size_v</code> | Dimension of each attention head for value/output |
| 46 | * | Can be different from <code>head_size</code> in FA3 |
| 47 | * | Set to 0 to use the same value as <code>head_size</code> |
| 48 | * | |
| 49 | * <code>@param q_batch_stride</code> | Stride (in elements) between consecutive batches in query tensor |
| 50 | * | For contiguous tensor: <code>seqLen_q * num_heads * head_size</code> |
| 51 | * | |
| 52 | * <code>@param q_row_stride</code> | Stride (in elements) between consecutive rows (sequences) in query tensor |
| 53 | * | For contiguous tensor: <code>num_heads * head_size</code> |
| 54 | * | |
| 55 | * <code>@param q_head_stride</code> | Stride (in elements) between consecutive heads in query tensor |
| 56 | * | For contiguous tensor: <code>head_size</code> |
| 57 | * | |
| 58 | * <code>@param k_batch_stride</code> | Stride (in elements) between consecutive batches in key tensor |
| 59 | * | For contiguous tensor: <code>seqLen_k * num_heads_k * head_size</code> |
| 60 | * | |
| 61 | * <code>@param k_row_stride</code> | Stride (in elements) between consecutive rows (sequences) in key tensor |
| 62 | * | For contiguous tensor: <code>num_heads_k * head_size</code> |
| 63 | * | |
| 64 | * <code>@param k_head_stride</code> | Stride (in elements) between consecutive heads in key tensor |
| 65 | * | For contiguous tensor: <code>head_size</code> |

66 *
67 * **@param** `v_batch_stride` *Stride (in elements) between consecutive batches in value tensor*
68 * *For contiguous tensor: $seq_len_k * num_heads_k * head_size_v$*
69 *
70 * **@param** `v_row_stride` *Stride (in elements) between consecutive rows (sequences) in value tensor*
71 * *For contiguous tensor: $num_heads_k * head_size_v$*
72 *
73 * **@param** `v_head_stride` *Stride (in elements) between consecutive heads in value tensor*
74 * *For contiguous tensor: $head_size_v$*
75 *
76 * **@param** `v_dim_stride` *Stride (in elements) between consecutive dimensions in value tensor*
77 * *For row-major: 1, for column-major: v_head_stride*
78 * *FA3 supports column-major layout for value tensor*
79 *
80 * **@param** `o_batch_stride` *Stride (in elements) between consecutive batches in output tensor*
81 * *For contiguous tensor: $seq_len_q * num_heads * head_size_v$*
82 *
83 * **@param** `o_row_stride` *Stride (in elements) between consecutive rows (sequences) in output tensor*
84 * *For contiguous tensor: $num_heads * head_size_v$*
85 *
86 * **@param** `o_head_stride` *Stride (in elements) between consecutive heads in output tensor*
87 * *For contiguous tensor: $head_size_v$*
88 *
89 * **@param** `softmax_scale` *Scale factor applied to attention scores before softmax*
90 * *Typically $1/\sqrt{head_size}$ to normalize variance*
91 * *Affects the sharpness of attention distribution*
92 *
93 * **@param** `softcap` *Softcapping value for attention scores (0.0 to disable)*
94 * *If > 0.0, applies a soft cap to attention scores to prevent extreme values*
95 * *Helps with training stability for very long sequences*
96 *
97 * **@param** `is_causal` *Whether to apply causal attention mask*
98 * *If true, each position can only attend to previous positions (including itself)*

99 * *Used in autoregressive models to prevent looking at future tokens*

100 *

101 * *@param window_size_left Left window size for local attention (-1 for unlimited)*

102 * *If ≥ 0 , limits how far back each query can attend*

103 * *Useful for long sequence efficiency or implementing sliding window attention*

104 *

105 * *@param window_size_right Right window size for local attention (-1 for unlimited)*

106 * *If ≥ 0 , limits how far forward each query can attend*

107 * *Only relevant when `is_causal=false`*

108 *

109 * *@param cu_seqlens_q Cumulative sequence lengths for queries (batch_size+1)*

110 * *Format: [0, len1, len1+len2, len1+len2+len3, ...]*

111 * *Used for variable length sequences. Set to nullptr for fixed length.*

112 *

113 * *@param cu_seqlens_k Cumulative sequence lengths for keys/values (batch_size+1)*

114 * *Format: [0, len1, len1+len2, len1+len2+len3, ...]*

115 * *Used for variable length sequences. Set to nullptr for fixed length.*

116 *

117 * *@param cu_seqlens_knew Cumulative sequence lengths for new keys/values (batch_size+1)*

118 * *Used when appending to KV cache with variable length sequences*

119 *

120 * *@param seqused_q Actual sequence usage for queries (batch_size)*

121 * *Specifies how much of each sequence is actually used*

122 * *New in FA3: allows specifying active sequence lengths*

123 *

124 * *@param seqused_k Actual sequence usage for keys (batch_size)*

125 * *Specifies how much of each sequence in the KV cache is actually used*

126 * *Useful when allocated length exceeds used length*

127 *

128 * *@param leftpad_k Left padding for each sequence in key (batch_size)*

129 * *Specifies number of padding tokens at the beginning of each sequence*

130 * *Affects position calculations for rotary embeddings*

131 *

132 * *@param* `qv_ptr` *Pointer to fused QV tensor for QV attention*
(total_q, nheads, headdim_v) or nullptr

133 * *New in FA3: allows fused QV computation for certain*
architectures

134 *

135 * *@param* `qv_batch_stride` *Stride between batches in QV tensor*

136 * *@param* `qv_row_stride` *Stride between rows in QV tensor*

137 * *@param* `qv_head_stride` *Stride between heads in QV tensor*

138 *

139 * *@param* `q_descale_ptr` *Pointer to query descaling factors (batch_size,*
nheads) or nullptr

140 * *Used for FP8 precision scaling in FA3*

141 *

142 * *@param* `k_descale_ptr` *Pointer to key descaling factors (batch_size,*
nheads_k) or nullptr

143 * *Used for FP8 precision scaling in FA3*

144 *

145 * *@param* `v_descale_ptr` *Pointer to value descaling factors (batch_size,*
nheads_k) or nullptr

146 * *Used for FP8 precision scaling in FA3*

147 *

148 * *@param* `q_descale_batch_stride` *Stride between batches in query descaling*
factors

149 * *@param* `q_descale_head_stride` *Stride between heads in query descaling*
factors

150 * *@param* `k_descale_batch_stride` *Stride between batches in key descaling*
factors

151 * *@param* `k_descale_head_stride` *Stride between heads in key descaling*
factors

152 * *@param* `v_descale_batch_stride` *Stride between batches in value descaling*
factors

153 * *@param* `v_descale_head_stride` *Stride between heads in value descaling*
factors

154 *

155 * *@param* `knew_ptr` *New key tensor data pointer (batch_size,*
seqlen_knew, nheads_k, headdim)

156 * *Contains new keys to be appended to the KV cache*

157 *

158 * *@param* `vnew_ptr` *New value tensor data pointer (batch_size,*
seqlen_knew, nheads_k, headdim_v)

159 * *Contains new values to be appended to the KV cache*

160 *

161 * *@param* `seqlen_knew` *Length of new keys/values to append to the cache*

162 * *Typically 1 for autoregressive generation, or chunk*
size for batch processing

163 *

164 * *@param* `knew_batch_stride` *Stride between batches in new key tensor*

165 * *@param* `knew_row_stride` *Stride between rows in new key tensor*

166 * `@param` `knew_head_stride` Stride between heads in new key tensor

167 * `@param` `vnew_batch_stride` Stride between batches in new value tensor

168 * `@param` `vnew_row_stride` Stride between rows in new value tensor

169 * `@param` `vnew_head_stride` Stride between heads in new value tensor

170 *

171 * `@param` `rotary_cos_ptr` Cosine table for rotary embeddings (`seqLen_ro`,
`rotary_dim/2`)

172 * Contains precomputed cosine values for rotary
position embeddings

173 *

174 * `@param` `rotary_sin_ptr` Sine table for rotary embeddings (`seqLen_ro`,
`rotary_dim/2`)

175 * Contains precomputed sine values for rotary
position embeddings

176 *

177 * `@param` `rotary_dim` Dimension for rotary embeddings (must be \leq
`head_size` and divisible by 16)

178 * Only the first `rotary_dim` dimensions of each head
will have rotary embeddings

179 *

180 * `@param` `is_rotary_interleaved` Whether rotary embeddings are interleaved

181 * If true, dimensions are paired as (0,1),
(2,3), etc. (GPT-J style)

182 * If false, dimensions are paired as
(0, `rotary_dim/2`), (1, `rotary_dim/2+1`), etc. (GPT-NeoX style)

183 *

184 * `@param` `kv_batch_idx` Indices for batch mapping (`batch_size`)

185 * Maps from logical batch indices to physical
indices in the KV cache

186 *

187 * `@param` `page_table` Block table for paged KV cache (`batch_size`,
`max_blocks_per_seq`)

188 * Contains block indices for each sequence in paged
KV cache

189 *

190 * `@param` `page_table_batch_stride` Stride between batches in block table

191 * `@param` `page_size` Size of each page in paged KV cache

192 * `@param` `num_pages` Total number of pages in paged KV cache

193 * `@param` `pagedkv_tma` Whether to use TMA (Tensor Memory Access) for
paged KV cache

194 * TMA is a hardware feature in NVIDIA Hopper
architecture

195 *

196 * `@param` `num_splits` Number of splits for computation (0 for auto)

197 * Splits the key/value tensors into chunks along
sequence dimension

198 *

```

199 * @param pack_gqa Whether to pack grouped-query attention for
optimization
200 * New in FA3: allows specialized GQA implementation
201 *
202 * @param scheduler_metadata Pointer to scheduler metadata for varlen
optimization
203 * New in FA3: allows precomputed scheduling
information
204 *
205 * @param sm_margin Number of SMs to reserve for communication (0 for
none)
206 * New in FA3: allows reserving SMs for communication
in multi-GPU scenarios
207 *
208 * @param is_bf16 Whether tensors are in BF16 format (true) or FP16
format (false)
209 * @param is_fp32 Whether tensors are in FP32 format
210 * @param is_e4m3 Whether tensors are in E4M3 (FP8) format
211 *
212 * @param stream CUDA stream to execute the kernel on
213 *
214 * @return Status code (0 for success, non-zero for failure)
215 */
216 int run_mha_fwd(
217     void* dev, // Device handle
218
219     // Core tensor pointers
220     void* q_ptr,
221     void* k_ptr,
222     void* v_ptr,
223     void* o_ptr,
224     void* softmax_lse_ptr = nullptr,
225
226     // Dimensions
227     int batch_size = 1,
228     int seqlen_q = 0,
229     int seqlen_k = 0,
230     int num_heads = 0,
231     int num_heads_k = 0,
232     int head_size = 0,
233     int head_size_v = 0, // New in FA3: can be different from
head_size
234
235     // Stride information (in elements, not bytes)
236     int64_t q_batch_stride = 0,
237     int64_t q_row_stride = 0,
238     int64_t q_head_stride = 0,

```

```
239     int64_t k_batch_stride = 0,
240     int64_t k_row_stride = 0,
241     int64_t k_head_stride = 0,
242     int64_t v_batch_stride = 0,
243     int64_t v_row_stride = 0,
244     int64_t v_head_stride = 0,
245     int64_t v_dim_stride = 1, // New in FA3: supports column-major
layout
246     int64_t o_batch_stride = 0,
247     int64_t o_row_stride = 0,
248     int64_t o_head_stride = 0,
249
250     // Scaling parameters
251     float softmax_scale = 0.0f,
252     float softcap = 0.0f,
253
254     // Masking parameters
255     bool is_causal = false,
256     int window_size_left = -1,
257     int window_size_right = -1,
258
259     // Variable length parameters
260     int* cu_seqlens_q = nullptr,
261     int* cu_seqlens_k = nullptr,
262     int* cu_seqlens_knew = nullptr, // New in FA3
263     int* seqused_q = nullptr, // New in FA3
264     int* seqused_k = nullptr,
265     int* leftpad_k = nullptr,
266
267     // QV fusion parameters (new in FA3)
268     void* qv_ptr = nullptr,
269     int64_t qv_batch_stride = 0,
270     int64_t qv_row_stride = 0,
271     int64_t qv_head_stride = 0,
272
273     // Descaling parameters for FP8 (new in FA3)
274     float* q_descale_ptr = nullptr,
275     float* k_descale_ptr = nullptr,
276     float* v_descale_ptr = nullptr,
277     int64_t q_descale_batch_stride = 0,
278     int64_t q_descale_head_stride = 0,
279     int64_t k_descale_batch_stride = 0,
280     int64_t k_descale_head_stride = 0,
281     int64_t v_descale_batch_stride = 0,
282     int64_t v_descale_head_stride = 0,
283
284     // KV cache parameters
```

```

285     void* knew_ptr = nullptr,
286     void* vnew_ptr = nullptr,
287     int seqlen_knew = 0,
288     int64_t knew_batch_stride = 0,
289     int64_t knew_row_stride = 0,
290     int64_t knew_head_stride = 0,
291     int64_t vnew_batch_stride = 0,
292     int64_t vnew_row_stride = 0,
293     int64_t vnew_head_stride = 0,
294
295     // Rotary embedding parameters
296     void* rotary_cos_ptr = nullptr,
297     void* rotary_sin_ptr = nullptr,
298     int rotary_dim = 0,
299     bool is_rotary_interleaved = true,
300
301     // Cache indexing
302     int* kv_batch_idx = nullptr,
303
304     // Paged KV cache parameters
305     int* page_table = nullptr,
306     int64_t page_table_batch_stride = 0,
307     int page_size = 0,
308     int num_pages = 0,
309     bool pagedkv_tma = false,           // New in FA3: TMA support
310
311     // Split computation parameters
312     int num_splits = 1,
313     bool pack_gqa = false,             // New in FA3
314     void* scheduler_metadata = nullptr, // New in FA3
315     int sm_margin = 0,                 // New in FA3
316
317     // Data type flags
318     bool is_bf16 = false,
319     bool is_fp32 = false,
320     bool is_e4m3 = false,             // New in FA3: FP8 support
321
322     // Execution parameters
323     cudaStream_t stream = nullptr
324 );

```

With kvcache

代码块

```

1  /**
2  * @brief Executes the Flash Attention 3 forward pass with KV cache

```

3 *
4 * @param dev Device handle for CUDA/CPU device management
5 *
6 * @param q_ptr Pointer to query tensor data (batch_size, seqLen_q, nheads, headdim)
7 * Current query tokens to compute attention for
8 *
9 * @param k_cache_ptr Pointer to key cache tensor data
10 * For standard KV cache: (batch_size_cache, seqLen_cache, nheads_k, headdim)
11 * For paged KV cache: (num_blocks, page_size, nheads_k, headdim)
12 *
13 * @param v_cache_ptr Pointer to value cache tensor data
14 * For standard KV cache: (batch_size_cache, seqLen_cache, nheads_k, headdim_v)
15 * For paged KV cache: (num_blocks, page_size, nheads_k, headdim_v)
16 *
17 * @param o_ptr Pointer to output tensor data (batch_size, seqLen_q, nheads, headdim_v)
18 * Will store the result of attention computation
19 *
20 * @param softmax_lse_ptr Pointer to store log-sum-exp values (batch_size, nheads, seqLen_q)
21 * Used for numerical stability. Set to nullptr if not needed.
22 *
23 * @param batch_size Number of sequences in the batch
24 * @param seqLen_q Length of each query sequence (typically 1 for autoregressive generation)
25 * @param seqLen_cache Maximum sequence length in the KV cache
26 * @param num_heads Number of attention heads for query tensor
27 * @param num_heads_k Number of attention heads for key/value tensors (for MQA/GQA)
28 * @param head_size Dimension of each attention head for query/key
29 * @param head_size_v Dimension of each attention head for value/output (can be different from head_size)
30 *
31 * @param q_batch_stride Stride between batches in query tensor
32 * @param q_row_stride Stride between rows in query tensor
33 * @param q_head_stride Stride between heads in query tensor
34 * @param k_cache_batch_stride Stride between batches in key cache tensor (ignored for paged KV cache)
35 * @param k_cache_row_stride Stride between rows in key cache tensor
36 * @param k_cache_head_stride Stride between heads in key cache tensor

37 * @param v_cache_batch_stride Stride between batches in value cache tensor
(ignored for paged KV cache)

38 * @param v_cache_row_stride Stride between rows in value cache tensor

39 * @param v_cache_head_stride Stride between heads in value cache tensor

40 * @param v_cache_dim_stride Stride between dimensions in value cache tensor
(for column-major layout)

41 * @param o_batch_stride Stride between batches in output tensor

42 * @param o_row_stride Stride between rows in output tensor

43 * @param o_head_stride Stride between heads in output tensor

44 *

45 * @param softmax_scale Scale factor for attention scores

46 * @param softcap Softcapping value (0.0 to disable)

47 *

48 * @param is_causal Whether to apply causal attention mask

49 * @param window_size_left Left window size for local attention (-1 for
unlimited)

50 * @param window_size_right Right window size for local attention (-1 for
unlimited)

51 *

52 * @param k_new_ptr New key tensor data pointer (batch_size,
seqlen_k_new, nheads_k, headdim)

53 * Contains new keys to be appended to the KV cache

54 *

55 * @param v_new_ptr New value tensor data pointer (batch_size,
seqlen_k_new, nheads_k, headdim_v)

56 * Contains new values to be appended to the KV cache

57 *

58 * @param seqlen_k_new Length of new keys/values to append to the cache

59 * Typically 1 for autoregressive generation

60 *

61 * @param k_new_batch_stride Stride between batches in new key tensor

62 * @param k_new_row_stride Stride between rows in new key tensor

63 * @param k_new_head_stride Stride between heads in new key tensor

64 * @param v_new_batch_stride Stride between batches in new value tensor

65 * @param v_new_row_stride Stride between rows in new value tensor

66 * @param v_new_head_stride Stride between heads in new value tensor

67 *

68 * @param qv_ptr Pointer to fused QV tensor for QV attention
(batch_size, seqlen_q, nheads, headdim_v)

69 * New in FA3: allows fused QV computation for certain
architectures

70 *

71 * @param qv_batch_stride Stride between batches in QV tensor

72 * @param qv_row_stride Stride between rows in QV tensor

73 * @param qv_head_stride Stride between heads in QV tensor

74 *

75 * @param q_descale_ptr Pointer to query descaling factors (batch_size, nheads) or nullptr

76 * @param k_descale_ptr Pointer to key descaling factors (batch_size, nheads_k) or nullptr

77 * @param v_descale_ptr Pointer to value descaling factors (batch_size, nheads_k) or nullptr

78 *

79 * @param cache_seqLens Current sequence lengths in KV cache (batch_size)

80 * Specifies the current length of each sequence in the cache

81 *

82 * @param rotary_cos_ptr Cosine table for rotary embeddings (seqLen_ro, rotary_dim/2)

83 * @param rotary_sin_ptr Sine table for rotary embeddings (seqLen_ro, rotary_dim/2)

84 * @param rotary_dim Dimension for rotary embeddings (must be <= head_size and divisible by 16)

85 * @param is_rotary_interleaved Whether rotary embeddings are interleaved

86 *

87 * @param cache_batch_idx Indices for batch mapping (batch_size)

88 * Maps from logical batch indices to physical indices in the KV cache

89 *

90 * @param cache_leftpad Left padding for each sequence in cache (batch_size)

91 * Specifies number of padding tokens at the beginning of each sequence

92 *

93 * @param page_table Block table for paged KV cache (batch_size, max_blocks_per_seq)

94 * Contains block indices for each sequence in paged KV cache

95 *

96 * @param page_table_batch_stride Stride between batches in block table

97 * @param page_size Size of each page in paged KV cache

98 * @param pagedkv_tma Whether to use TMA (Tensor Memory Access) for paged KV cache

99 *

100 * @param scheduler_metadata Pointer to scheduler metadata for varlen optimization

101 *

102 * @param num_splits Number of splits for computation (0 for auto)

103 * @param pack_gqa Whether to pack grouped-query attention for optimization

104 * @param sm_margin Number of SMs to reserve for communication (0 for none)

105 *

```

106  * @param is_bf16          Whether tensors are in BF16 format (true) or FP16
format (false)
107  * @param is_fp32          Whether tensors are in FP32 format
108  * @param is_e4m3          Whether tensors are in E4M3 (FP8) format
109  *
110  * @param stream           CUDA stream to execute the kernel on
111  *
112  * @param return_softmax_lse Whether to return the log-sum-exp values
113  *
114  * @return Status code (0 for success, non-zero for failure)
115  */
116  int run_mha_fwd_kvcache(
117      void* dev,                // Device handle
118
119      // Core tensor pointers
120      void* q_ptr,
121      void* k_cache_ptr,
122      void* v_cache_ptr,
123      void* o_ptr,
124      void* softmax_lse_ptr = nullptr,
125
126      // Dimensions
127      int batch_size = 1,
128      int seqlen_q = 1,
129      int seqlen_cache = 0,
130      int num_heads = 0,
131      int num_heads_k = 0,
132      int head_size = 0,
133      int head_size_v = 0,
134
135      // Stride information (in elements, not bytes)
136      int64_t q_batch_stride = 0,
137      int64_t q_row_stride = 0,
138      int64_t q_head_stride = 0,
139      int64_t k_cache_batch_stride = 0,
140      int64_t k_cache_row_stride = 0,
141      int64_t k_cache_head_stride = 0,
142      int64_t v_cache_batch_stride = 0,
143      int64_t v_cache_row_stride = 0,
144      int64_t v_cache_head_stride = 0,
145      int64_t v_cache_dim_stride = 1,
146      int64_t o_batch_stride = 0,
147      int64_t o_row_stride = 0,
148      int64_t o_head_stride = 0,
149
150      // Scaling parameters
151      float softmax_scale = 0.0f,

```

```
152     float softcap = 0.0f,
153
154     // Masking parameters
155     bool is_causal = false,
156     int window_size_left = -1,
157     int window_size_right = -1,
158
159     // New KV parameters
160     void* k_new_ptr = nullptr,
161     void* v_new_ptr = nullptr,
162     int seq_len_k_new = 0,
163     int64_t k_new_batch_stride = 0,
164     int64_t k_new_row_stride = 0,
165     int64_t k_new_head_stride = 0,
166     int64_t v_new_batch_stride = 0,
167     int64_t v_new_row_stride = 0,
168     int64_t v_new_head_stride = 0,
169
170     // QV fusion parameters
171     void* qv_ptr = nullptr,
172     int64_t qv_batch_stride = 0,
173     int64_t qv_row_stride = 0,
174     int64_t qv_head_stride = 0,
175
176     // Descaling parameters for FP8
177     float* q_descale_ptr = nullptr,
178     float* k_descale_ptr = nullptr,
179     float* v_descale_ptr = nullptr,
180
181     // KV cache parameters
182     int* cache_seq_lens = nullptr,
183
184     // Rotary embedding parameters
185     void* rotary_cos_ptr = nullptr,
186     void* rotary_sin_ptr = nullptr,
187     int rotary_dim = 0,
188     bool is_rotary_interleaved = true,
189
190     // Cache indexing
191     int* cache_batch_idx = nullptr,
192     int* cache_leftpad = nullptr,
193
194     // Paged KV cache parameters
195     int* page_table = nullptr,
196     int64_t page_table_batch_stride = 0,
197     int page_size = 0,
198     bool pagedkv_tma = false,
```

```

199
200 // Scheduler metadata
201 void* scheduler_metadata = nullptr,
202
203 // Split computation parameters
204 int num_splits = 0,
205 bool pack_gqa = false,
206 int sm_margin = 0,
207
208 // Data type flags
209 bool is_bf16 = false,
210 bool is_fp32 = false,
211 bool is_e4m3 = false,
212
213 // Execution parameters
214 cudaStream_t stream = nullptr,
215
216 // Return options
217 bool return_softmax_lse = false
218 );

```

ISSUE

传入 cuda kernel 之前会将 tensor 处理为 contiguous

代码块

```

1 @staticmethod
2 def forward(ctx, q, k, v, cu_seqlens_q, cu_seqlens_k, max_seq_len_q,
3 max_seq_len_k, ...):
4     # ...
5     q, k, v = [maybe_contiguous(x) for x in (q, k, v)]
6     # ...

```

cuda FA2 支持 FP16、BF16 和 FP32 浮点数据类型，

cuda FA3 (SM90+架构)，额外支持 E4M3 (FP8 格式，4 位指数，3 位尾数)

_C::fused_add_rms_norm[vLLM]

参考信息

vLLM 对应的接口实现：https://github.com/vllm-project/vllm/blob/b0ccfc565ac263b4097175b3ba33ad1d5dc5df8a/csrc/layernorm_kernels.cu#L52

Function 接口

代码块

```
1  /**
2   * @brief Performs fused add and RMS normalization
3   *
4   * @param dev Device handle for CUDA/CPU device management
5   * @param input First input tensor
6   * @param residual Second input tensor to be added to the first
7   * @param weight Scale factors applied after normalization
8   * @param out_sizes Output tensor sizes
9   * @param out_strides Output tensor strides
10  * @param input_sizes First input tensor sizes
11  * @param input_strides First input tensor strides
12  * @param residual_sizes Second input tensor sizes
13  * @param residual_strides Second input tensor strides
14  * @param weight_sizes Weight tensor sizes
15  * @param weight_strides Weight tensor strides
16  * @param ndims Number of dimensions in the tensors
17  * @param epsilon Small constant added for numerical stability
18  * @param num_tokens Number of tokens in the input tensor
19  * @param hidden_size Size of the last dimension in the input tensor
20  * @param out_dtype Data type of output tensor
21  * @param input_dtype Data type of input tensors
22  * @param residual_dtype Data type of residual tensor
23  * @param weight_dtype Data type of weight tensor
24  *
25  * @note This function:
26  *     1. Adds input and residual tensors
27  *     2. Computes RMS normalization on the sum
28  *     3. Applies weight scaling
29  *     4. Stores the result in the input tensor
30  *     5. Optionally updates the residual tensor with the sum
31  */
32 void fused_add_rms_norm_kernel(
33     void* dev, // Device handle
34     void* input, // First input tensor data pointer
35     void* residual, // Second input tensor data pointer
36     const void* weight, // Weight tensor data pointer
37     const int64_t* input_sizes, // First input tensor sizes
38     const int64_t* input_strides, // First input tensor strides
39     const int64_t* residual_sizes, // Second input tensor sizes
40     const int64_t* residual_strides, // Second input tensor strides
41     const int64_t* weight_sizes, // Weight tensor sizes
42     const int64_t* weight_strides, // Weight tensor strides
43     int ndims, // Number of dimensions in tensors
44     float epsilon, // Small constant for numerical stability
```

```

45     int num_tokens,           // Number of tokens in the input tensor
46     int hidden_size,         // Size of the last dimension in the input
    tensor
47     int input_dtype,         // Data type of input tensors
48     int residual_dtype,      // Data type of residual tensor
49     int weight_dtype,        // Data type of weight tensor
50     bool inplace_residual    // Whether to update residual in-place
51 );

```

ISSUE

- 支持的数据类型 float32 float16 bfloat16
- 从 vllm/model_executor/layers/layernorm.py 调用看，input, residual, weight 都为 contiguous 布局，但是具体实现也支持非 contiguous

_C::rms_norm[vLLM]

参考信息

实现参考: https://github.com/vllm-project/vllm/blob/cec8c7d7f8753d13737427ceb5cebe987f5f0549/csrc/layernorm_kernels.cu#L17

Function 接口

代码块

```

1  /**
2   * @brief Performs RMS (Root Mean Square) normalization on input tensor
3   *
4   * @param dev           Device handle for CUDA/CPU device management
5   * @param out           Output tensor to store the normalized result
6   * @param input         Input tensor to be normalized
7   * @param weight        Scale factors applied after normalization
8   * @param out_sizes     Output tensor sizes
9   * @param out_strides   Output tensor strides (can be non-contiguous)
10  * @param input_sizes   Input tensor sizes
11  * @param input_strides Input tensor strides (can be non-contiguous)
12  * @param weight_sizes  Weight tensor sizes
13  * @param weight_strides Weight tensor strides
14  * @param ndims         Number of dimensions in the input/output tensors
15  * @param epsilon       Small constant added for numerical stability
16  * @param num_tokens    Number of tokens in the input tensor
17  * @param hidden_size   Size of the last dimension in the input tensor
18  * @param out_dtype     Data type of output tensor
19  * @param input_dtype   Data type of input tensor

```

```

20 * @param weight_dtype Data type of weight tensor
21 *
22 * @note Memory Layout:
23 * - Fully supports non-contiguous tensors
24 * - Handles arbitrary strides efficiently
25 * - Optimized for both contiguous and non-contiguous cases
26 *
27 * @note Behavior:
28 * - Normalizes using RMS (sqrt of mean of squares)
29 * - Applies weight scaling after normalization
30 * - Supports different data types for input, output, and weight
31 * - Optimized implementations for different hardware and data types
32 */
33 void rms_norm_kernel(
34     void* dev, // Device handle
35     void* out, // Output tensor data pointer
36     const void* input, // Input tensor data pointer
37     const void* weight, // Weight tensor data pointer
38     const int64_t* out_sizes, // Output sizes
39     const int64_t* out_strides, // Output strides
40     const int64_t* input_sizes, // Input tensor sizes
41     const int64_t* input_strides, // Input tensor strides
42     const int64_t* weight_sizes, // Weight tensor sizes
43     const int64_t* weight_strides, // Weight tensor strides
44     int ndims, // Number of dimensions in tensors
45     const float epsilon, // Small constant for numerical stability
46     int num_tokens, // Number of tokens in the input tensor
47     int hidden_size, // Size of the last dimension in the input
48     tensor
49     int out_dtype, // Data type of output tensor
50     int input_dtype, // Data type of input tensor
51     int weight_dtype // Data type of weight tensor
52 );

```

ISSUE

支持的数据类型 float32 float16 bfloat16

从 vllm/model_executor/layers/layernorm.py 调用看, input, residual, weight 都为 contiguous 布局, 但是具体实现也支持非 contiguous

_C::rotary_embedding[vLLM]

参考信息

参考 cuda 实现: [https://github.com/vllm-](https://github.com/vllm-project/vllm/blob/18ed3132d2bfe1df9a74729457b69243955221e8/csrc/pos_encoding_kernels.cu#L124)

[project/vllm/blob/18ed3132d2bfe1df9a74729457b69243955221e8/csrc/pos_encoding_kernels.cu#L124](https://github.com/vllm-project/vllm/blob/18ed3132d2bfe1df9a74729457b69243955221e8/csrc/pos_encoding_kernels.cu#L124)

Function 接口

代码块

```
1  /**
2   * @brief Applies rotary position embeddings to query and key tensors
3   *
4   * @param dev Device handle for CUDA/CPU device management
5   * @param positions Position indices tensor [batch_size, seq_len]
6   * or [num_tokens]
7   * @param query Query tensor to apply rotary embeddings to
8   * @param key Key tensor to apply rotary embeddings to
9   * @param cos_sin_cache Precomputed cosine and sine values
10  * [max_position, rot_dim]
11  * @param positions_sizes Sizes of positions tensor
12  * @param positions_strides Strides of positions tensor
13  * @param query_sizes Sizes of query tensor
14  * @param query_strides Strides of query tensor
15  * @param key_sizes Sizes of key tensor
16  * @param key_strides Strides of key tensor
17  * @param cache_sizes Sizes of cos_sin_cache tensor
18  * @param cache_strides Strides of cos_sin_cache tensor
19  * @param positions_ndims Number of dimensions in positions tensor
20  * @param query_ndims Number of dimensions in query tensor
21  * @param key_ndims Number of dimensions in key tensor
22  * @param cache_ndims Number of dimensions in cos_sin_cache tensor
23  * @param head_size Size of each attention head
24  * @param is_neox Whether to use GPT-NeoX style rotary embeddings
25  * @param query_dtype Data type of query tensor
26  * @param key_dtype Data type of key tensor
27  * @param cache_dtype Data type of cos_sin_cache tensor
28  *
29  * @note Memory Layout:
30  * - Supports both [batch_size, seq_len, num_heads, head_size] and
31  * [num_tokens, num_heads, head_size] layouts
32  * - Handles non-contiguous tensors through explicit stride parameters
33  * - Optimized for different data types and hardware
34  *
35  * @note Behavior:
36  * - Applies rotary position embeddings in-place to query and key tensors
37  * - Supports both GPT-J style (default) and GPT-NeoX style rotary
38  * embeddings
```

```

35  *      - Efficiently handles grouped-query attention (GQA) where num_heads >
      num_kv_heads
36  */
37  void rotary_embedding_kernel(
38      void* dev,                // Device handle
39      const int64_t* positions,  // Position indices tensor
40      void* query,              // Query tensor
41      void* key,                // Key tensor
42      const void* cos_sin_cache, // Precomputed cosine and sine values
43      const int64_t* positions_sizes, // Sizes of positions tensor
44      const int64_t* positions_strides, // Strides of positions tensor
45      const int64_t* query_sizes, // Sizes of query tensor
46      const int64_t* query_strides, // Strides of query tensor
47      const int64_t* key_sizes, // Sizes of key tensor
48      const int64_t* key_strides, // Strides of key tensor
49      const int64_t* cache_sizes, // Sizes of cos_sin_cache tensor
50      const int64_t* cache_strides, // Strides of cos_sin_cache tensor
51      int positions_ndims, // Number of dimensions in positions
      tensor
52      int query_ndims, // Number of dimensions in query tensor
53      int key_ndims, // Number of dimensions in key tensor
54      int cache_ndims, // Number of dimensions in
      cos_sin_cache tensor
55      int head_size, // Size of each attention head
56      bool is_neox, // Whether to use GPT-NeoX style rotary
      embeddings
57      int query_dtype, // Data type of query tensor
58      int key_dtype, // Data type of key tensor
59      int cache_dtype // Data type of cos_sin_cache tensor
60  );

```

ISSUE

- 支持的数据类型: float32, float16, bfloat16
- cuda 的实现有对非 contiguous 的支持

_C::silu_and_mul[vLLM]

参考信息

参考 cuda 实现: https://github.com/vllm-project/vllm/blob/18ed3132d2bfe1df9a74729457b69243955221e8/csrc/activation_kernels.cu#L82

Function 接口

```

1  1.2.7.4
2  * @brief Applies SiLU activation to first half of input and multiplies with
   second half
3  *
4  * @param dev           Device handle for CUDA/CPU device management
5  * @param out          Output tensor to store the result
6  * @param input        Input tensor with twice the size of output in last
   dimension
7  * @param out_sizes    Output tensor sizes
8  * @param out_strides  Output tensor strides (can be non-contiguous)
9  * @param input_sizes  Input tensor sizes
10 * @param input_strides Input tensor strides (can be non-contiguous)
11 * @param ndims         Number of dimensions in the tensors
12 * @param out_dtype     Data type of output tensor
13 * @param input_dtype   Data type of input tensor
14 *
15 * @note Memory Layout:
16 *     - Input tensor has shape [..., 2*d] where last dimension contains
   pairs of values
17 *     - Output tensor has shape [..., d]
18 *     - Supports non-contiguous tensors through explicit stride parameters
19 *
20 * @note Behavior:
21 *     - Computes out[i] = SiLU(input[i]) * input[i+d] for i in [0,d)
22 *     - SiLU(x) = x * sigmoid(x)
23 *     - Optimized implementation for different hardware and data types
24 */
25 void silu_and_mul_kernel(
26     void* dev,           // Device handle
27     void* out,           // Output tensor data pointer
28     const void* input,   // Input tensor data pointer
29     const int64_t* out_sizes, // Output sizes
30     const int64_t* out_strides, // Output strides
31     const int64_t* input_sizes, // Input tensor sizes
32     const int64_t* input_strides, // Input tensor strides
33     int ndims,           // Number of dimensions in tensors
34     int out_dtype,       // Data type of output tensor
35     int input_dtype      // Data type of input tensor
36 );

```

ISSUE

支持非 contiguous 输入

Cuda kernel 支持的数据类型: float32, float16, bfloat16

首版开发说明

1. 数据类型支持bfloat16, float32和float16, 通过template设置
2. 目前llama和deepseek, 在dim1_size上已经满足1024B对齐的要求, 暂时不考虑不对齐pad的情况
3. dim1_size是output的dimension 1 size。input的dimension 1 size是2*dim1_size

代码块

```
1  /// @brief Applies SiLU activation to first half of input and multiplies with
    second half
2  /// @tparam T Datatype support [sifmt::bfloat16, sifmt::float16,
    sifmt::float32]
3  /// @param d_in Input tensor [dim0_size, 2*dim1_size]
4  /// @param d_out Output tensor [dim0_size, dim1_size]
5  /// @param dim0_size Size of dimension 0 for the input/output tensor
6  /// @param dim1_size Size of dimension 1 for the output tensor (total size
    must be 1024B aligned)
7  template <typename T>
8  void silu_and_mul(void *d_in, void *d_out, int dim0_size, int dim1_size);
```

_C_cache_ops::reshape_and_cache_flash[vLLM]

参考信息

reshape_and_cache_flash 是一个用于重塑和缓存注意力机制中的键值对的函数, 专为 Flash Attention 优化。通过分析 cache_kernels.cu 中的实现, 我可以确定其主要功能:

功能: 将模型生成的键 (key) 和值 (value) 张量重新排列并存储到 KV 缓存中

数据流: 从 [num_tokens, num_heads, head_size] 格式转换为 [num_blocks, block_size, num_heads, head_size] 格式

特殊处理:

支持数据类型转换 (如 FP16/BF16 到 FP8)

处理槽位映射 (slot mapping)

支持缩放因子 (用于量化)

参考实现: https://github.com/vllm-project/vllm/blob/18ed3132d2bfe1df9a74729457b69243955221e8/csrc/cache_kernels.cu#L411

Function 接口

代码块

```
1
2  /**
3  * @brief Reshapes and caches key-value tensors for Flash Attention
```

```

4 *
5 * @param dev Device handle for CUDA/CPU device management
6 * @param key Key tensor [num_tokens, num_heads, head_size]
7 * @param value Value tensor [num_tokens, num_heads, head_size]
8 * @param key_cache Key cache tensor [num_blocks, block_size,
num_heads, head_size]
9 * @param value_cache Value cache tensor [num_blocks, block_size,
num_heads, head_size]
10 * @param slot_mapping Mapping from token indices to cache slots
[num_tokens]
11 * @param key_sizes Sizes of key tensor
12 * @param key_strides Strides of key tensor
13 * @param value_sizes Sizes of value tensor
14 * @param value_strides Strides of value tensor
15 * @param key_cache_sizes Sizes of key cache tensor
16 * @param key_cache_strides Strides of key cache tensor
17 * @param value_cache_sizes Sizes of value cache tensor
18 * @param value_cache_strides Strides of value cache tensor
19 * @param slot_mapping_sizes Sizes of slot mapping tensor
20 * @param slot_mapping_strides Strides of slot mapping tensor
21 * @param key_ndims Number of dimensions in key tensor
22 * @param value_ndims Number of dimensions in value tensor
23 * @param key_cache_ndims Number of dimensions in key cache tensor
24 * @param value_cache_ndims Number of dimensions in value cache tensor
25 * @param slot_mapping_ndims Number of dimensions in slot mapping tensor
26 * @param num_tokens Number of tokens to process
27 * @param num_heads Number of attention heads
28 * @param head_size Size of each attention head
29 * @param block_size Size of each cache block
30 * @param key_dtype Data type of key tensor
31 * @param value_dtype Data type of value tensor
32 * @param key_cache_dtype Data type of key cache tensor
33 * @param value_cache_dtype Data type of value cache tensor
34 * @param kv_cache_type Type of KV cache quantization ("fp8", "fp8_e4m3",
"auto", etc.)
35 * @param k_scale Scale factor for key quantization [1] or [num_heads]
36 * @param v_scale Scale factor for value quantization [1] or
[num_heads]
37 * @param k_scale_size Size of k_scale tensor (must be 1 or num_heads)
38 * @param v_scale_size Size of v_scale tensor (must be 1 or num_heads)
39 *
40 * @note Memory Layout:
41 * - Input tensors (key, value) have shape [num_tokens, num_heads,
head_size]
42 * - Cache tensors have shape [num_blocks, block_size, num_heads,
head_size]
43 * - Supports non-contiguous tensors through explicit stride parameters

```

```

44 *
45 * @note Scale Factors:
46 * - k_scale can be either a single value [1] or per-head values
  [num_heads]
47 * - v_scale can be either a single value [1] or per-head values
  [num_heads]
48 * - When size is 1, the same scale is applied to all heads
49 * - When size is num_heads, each head uses its own scale factor
50 *
51 * @note Behavior:
52 * - Reshapes key and value tensors to match the cache format
53 * - Maps tokens to their corresponding cache slots using slot_mapping
54 * - Supports optional quantization to FP8 format with scaling
55 * - Ignores tokens with negative slot indices (padding tokens)
56 */
57 void reshape_and_cache_flash_kernel(
58     void* dev, // Device handle
59     const void* key, // Key tensor data pointer
60     const void* value, // Value tensor data pointer
61     void* key_cache, // Key cache tensor data pointer
62     void* value_cache, // Value cache tensor data pointer
63     const int64_t* slot_mapping, // Slot mapping data pointer
64     const int64_t* key_sizes, // Key tensor sizes
65     const int64_t* key_strides, // Key tensor strides
66     const int64_t* value_sizes, // Value tensor sizes
67     const int64_t* value_strides, // Value tensor strides
68     const int64_t* key_cache_sizes, // Key cache tensor sizes
69     const int64_t* key_cache_strides, // Key cache tensor strides
70     const int64_t* value_cache_sizes, // Value cache tensor sizes
71     const int64_t* value_cache_strides, // Value cache tensor strides
72     const int64_t* slot_mapping_sizes, // Slot mapping tensor sizes
73     const int64_t* slot_mapping_strides, // Slot mapping tensor strides
74     int key_ndims, // Number of dimensions in key tensor
75     int value_ndims, // Number of dimensions in value
  tensor
76     int key_cache_ndims, // Number of dimensions in key cache
  tensor
77     int value_cache_ndims, // Number of dimensions in value
  cache tensor
78     int slot_mapping_ndims, // Number of dimensions in slot
  mapping tensor
79     int num_tokens, // Number of tokens to process
80     int num_heads, // Number of attention heads
81     int head_size, // Size of each attention head
82     int block_size, // Size of each cache block
83     int key_dtype, // Data type of key tensor
84     int value_dtype, // Data type of value tensor

```

```

85     int key_cache_dtype,           // Data type of key cache tensor
86     int value_cache_dtype,        // Data type of value cache tensor
87     const char* kv_cache_type,    // Type of KV cache quantization
88     const float* k_scale,         // Scale factor for key quantization
89     const float* v_scale,         // Scale factor for value quantization
90     int k_scale_size,             // Size of k_scale tensor (1 or
    num_heads)
91     int v_scale_size             // Size of v_scale tensor (1 or
    num_heads)
92 );

```

ISSUE

- 支持非 contiguous 输入
- Cuda kernel 支持的数据类型：float32, float16, bfloat16, 量化类型
__NV_E4M3/__NV_E5M2

_C_cache_ops::concat_and_cache_mla

参考信息

concat_and_cache_mla 的核心功能是：

1. 连接 (Concatenate) 两个张量：

kv_c: KV 内容张量 [num_tokens, kv_lora_rank]

k_pe: 键位置编码张量 [num_tokens, pe_dim]

2. 将连接后的数据缓存：

将连接后的数据存储到 KV 缓存中 [num_blocks, block_size, (kv_lora_rank + pe_dim)]

使用 slot_mapping 将每个 token 映射到缓存中的正确位置

3. 可选的量化处理：

支持将数据量化为 FP8 格式以节省内存

使用 scale 参数作为量化的缩放因子

参考实现：https://github.com/vllm-project/vllm/blob/18ed3132d2bfe1df9a74729457b69243955221e8/csrc/cache_kernels.cu#L462

Function 接口

代码块

```

1  /**
2  * @brief Concatenates and caches tensors for Multi-head Latent Attention (MLA)

```

```

3 *
4 * @param dev Device handle for CUDA/CPU device management
5 * @param kv_c KV content tensor [num_tokens, kv_lora_rank]
6 * @param k_pe Key position encoding tensor [num_tokens,
pe_dim]
7 * @param kv_cache KV cache tensor [num_blocks, block_size,
(kv_lora_rank + pe_dim)]
8 * @param slot_mapping Mapping from token indices to cache slots
[num_tokens]
9 * @param kv_c_sizes Sizes of KV content tensor
10 * @param kv_c_strides Strides of KV content tensor
11 * @param k_pe_sizes Sizes of key position encoding tensor
12 * @param k_pe_strides Strides of key position encoding tensor
13 * @param kv_cache_sizes Sizes of KV cache tensor
14 * @param kv_cache_strides Strides of KV cache tensor
15 * @param slot_mapping_sizes Sizes of slot mapping tensor
16 * @param slot_mapping_strides Strides of slot mapping tensor
17 * @param kv_c_ndims Number of dimensions in KV content tensor
18 * @param k_pe_ndims Number of dimensions in key position encoding
tensor
19 * @param kv_cache_ndims Number of dimensions in KV cache tensor
20 * @param slot_mapping_ndims Number of dimensions in slot mapping tensor
21 * @param num_tokens Number of tokens to process
22 * @param kv_lora_rank Rank of KV LoRA
23 * @param pe_dim Dimension of position encoding
24 * @param block_size Size of each cache block
25 * @param kv_c_dtype Data type of KV content tensor
26 * @param k_pe_dtype Data type of key position encoding tensor
27 * @param kv_cache_dtype Data type of KV cache tensor
28 * @param kv_cache_type Type of KV cache quantization
29 * @param scale Scale factor for quantization
30 *
31 * @note Memory Layout:
32 * - KV content tensor has shape [num_tokens, kv_lora_rank]
33 * - Key position encoding tensor has shape [num_tokens, pe_dim]
34 * - KV cache tensor has shape [num_blocks, block_size, (kv_lora_rank +
pe_dim)]
35 *
36 * @note Behavior:
37 * - Concatenates KV content and key position encoding
38 * - Maps tokens to their corresponding cache slots using slot_mapping
39 * - Supports optional quantization with scaling
40 */
41 void concat_and_cache_mla_kernel(
42     void* dev, // Device handle
43     const void* kv_c, // KV content tensor data pointer

```

```

44     const void* k_pe, // Key position encoding tensor data
pointer
45     void* kv_cache, // KV cache tensor data pointer
46     const int64_t* slot_mapping, // Slot mapping data pointer
47     const int64_t* kv_c_sizes, // KV content tensor sizes
48     const int64_t* kv_c_strides, // KV content tensor strides
49     const int64_t* k_pe_sizes, // Key position encoding tensor sizes
50     const int64_t* k_pe_strides, // Key position encoding tensor strides
51     const int64_t* kv_cache_sizes, // KV cache tensor sizes
52     const int64_t* kv_cache_strides, // KV cache tensor strides
53     const int64_t* slot_mapping_sizes, // Slot mapping tensor sizes
54     const int64_t* slot_mapping_strides, // Slot mapping tensor strides
55     int kv_c_ndims, // Number of dimensions in KV content
tensor
56     int k_pe_ndims, // Number of dimensions in key
position encoding tensor
57     int kv_cache_ndims, // Number of dimensions in KV cache
tensor
58     int slot_mapping_ndims, // Number of dimensions in slot
mapping tensor
59     int num_tokens, // Number of tokens to process
60     int kv_lora_rank, // Rank of KV LoRA
61     int pe_dim, // Dimension of position encoding
62     int block_size, // Size of each cache block
63     int kv_c_dtype, // Data type of KV content tensor
64     int k_pe_dtype, // Data type of key position encoding
tensor
65     int kv_cache_dtype, // Data type of KV cache tensor
66     const char* kv_cache_type, // Type of KV cache quantization
67     const float* scale // Scale factor for quantization size
= [1]
68 );

```

ISSUE

支持非 contiguous 输入

Cuda kernel 支持的数据类型: float32, float16, bfloat16, 量化类型

__NV_E4M3/__NV_E5M2

_moe_C::moe_align_block_size[vLLM]

参考信息

moe_align_block_size 的核心功能是:

1. 重排序 tokens:

- a. 根据分配给它们的专家 ID 对输入 tokens 进行重排序
- b. 确保分配给同一专家的 tokens 在内存中连续存储

2. 对齐到块大小:

- a. 将每个专家处理的 tokens 数量填充到 `block_size` 的倍数, 这种对齐优化了 GPU 上的并行计算效率

3. 生成专家映射:

- a. 为每个块分配对应的专家 ID
- b. 记录重排序后的 token 总数 (包括填充)

参考实现: [https://github.com/vllm-](https://github.com/vllm-project/vllm/blob/18ed3132d2bfe1df9a74729457b69243955221e8/csrc/moe/moe_align_sum_kernels.cu#L296)

[project/vllm/blob/18ed3132d2bfe1df9a74729457b69243955221e8/csrc/moe/moe_align_sum_kernels.cu#L296](https://github.com/vllm-project/vllm/blob/18ed3132d2bfe1df9a74729457b69243955221e8/csrc/moe/moe_align_sum_kernels.cu#L296)

Function 接口

代码块

```
1  /**
2   * @brief Aligns tokens by expert and arranges them in blocks for efficient
   MoE computation
3   *
4   * @param dev Device handle for CUDA/CPU device management
5   * @param topk_ids Tensor containing expert IDs for each token
   [num_tokens]
6   * @param sorted_token_ids Output tensor to store sorted token indices
   [num_tokens_post_pad]
7   * @param expert_ids Output tensor to store expert ID for each block
   [num_blocks]
8   * @param total_tokens_post_pad Output scalar to store total number of tokens
   after padding
9   * @param topk_ids_sizes Sizes of topk_ids tensor
10  * @param topk_ids_strides Strides of topk_ids tensor
11  * @param sorted_token_ids_sizes Sizes of sorted_token_ids tensor
12  * @param sorted_token_ids_strides Strides of sorted_token_ids tensor
13  * @param expert_ids_sizes Sizes of expert_ids tensor
14  * @param expert_ids_strides Strides of expert_ids tensor
15  * @param topk_ids_ndims Number of dimensions in topk_ids tensor
16  * @param sorted_token_ids_ndims Number of dimensions in sorted_token_ids
   tensor
17  * @param expert_ids_ndims Number of dimensions in expert_ids tensor
18  * @param num_tokens Number of tokens to process
19  * @param num_experts Number of experts in the MoE model
20  * @param block_size Size of each processing block (for alignment)
```

```

21 * @param topk_ids_dtype Data type of topk_ids tensor
22 * @param use_global_memory Whether to use global memory instead of shared
memory
23 * @param use_16bit_counters Whether to use 16-bit counters to save shared
memory
24 *
25 * @note Memory Layout:
26 * - topk_ids: Contains expert ID assigned to each token
27 * - sorted_token_ids: Contains token indices sorted by expert ID, with
padding
28 * - expert_ids: Contains expert ID for each block of tokens
29 *
30 * @note Behavior:
31 * - Groups tokens by expert ID and arranges them in blocks of size
block_size
32 * - Pads each expert's tokens to a multiple of block_size
33 * - Returns the total number of tokens after padding
34 * - Automatically selects between shared memory and global memory
implementation
35 *
36 */
37 void moe_align_block_size_kernel(
38     void* dev, // Device handle
39     const void* topk_ids, // Expert IDs for each token
40     void* sorted_token_ids, // Output sorted token indices
41     void* expert_ids, // Output expert ID for each block
42     void* total_tokens_post_pad, // Output total tokens after padding
43     const int64_t* topk_ids_sizes, // Sizes of topk_ids tensor
44     const int64_t* topk_ids_strides, // Strides of topk_ids tensor
45     const int64_t* sorted_token_ids_sizes, // Sizes of sorted_token_ids
tensor
46     const int64_t* sorted_token_ids_strides, // Strides of sorted_token_ids
tensor
47     const int64_t* expert_ids_sizes, // Sizes of expert_ids tensor
48     const int64_t* expert_ids_strides, // Strides of expert_ids tensor
49     int topk_ids_ndims, // Number of dimensions in topk_ids
tensor
50     int sorted_token_ids_ndims, // Number of dimensions in
sorted_token_ids tensor
51     int expert_ids_ndims, // Number of dimensions in expert_ids
tensor
52     int num_tokens, // Number of tokens to process
53     int num_experts, // Number of experts in the MoE model
54     int block_size, // Size of each processing block
55     int topk_ids_dtype, // Data type of topk_ids tensor
56     bool use_global_memory, // Whether to use global memory
57     bool use_16bit_counters // Whether to use 16-bit counters

```

58);

ISSUE

对于 topk_ids，默认是 int32 类型，如果 topk_ids < 65535，可选择 int16 类型

代码块

```
1     if (shared_mem_i16 < device_max_shared_mem && topk_ids.numel() <= 65535) {
2         use_i16 = true;
3     }
```

其余_ids 使用 int32 类型

支持非 contiguous

_moe_C::moe_sum[vLLM]

参考信息

moe_sum 的核心功能是：

1. 专家输出求和：

a. 对每个 token，将其分配到的多个专家（通常是 top-k 个专家）的输出进行求和，生成最终的 token 表示

2. 高效实现：

a. 针对常见的 top-k 值（2、3、4）提供了优化的 CUDA 实现

b. 对于其他 top-k 值，回退到 PyTorch 的通用求和操作

参考实现：https://github.com/vllm-project/vllm/blob/18ed3132d2bfe1df9a74729457b69243955221e8/csrc/moe/moe_align_sum_kernels.cu#L422

Function 接口

代码块

```
1  /**
2   * @brief Sums outputs from multiple experts for each token in MoE models
3   *
4   * @param dev Device handle for CUDA/CPU device management
5   * @param out Output tensor to store summed results
6   * [num_tokens, hidden_size]
```

```

6   * @param input          Input tensor with expert outputs [num_tokens,
topk, hidden_size]
7   * @param out_sizes     Sizes of output tensor
8   * @param out_strides   Strides of output tensor
9   * @param input_sizes   Sizes of input tensor
10  * @param input_strides  Strides of input tensor
11  * @param out_ndim       Number of dimensions in output tensor
12  * @param input_ndim    Number of dimensions in input tensor
13  * @param num_tokens     Number of tokens to process
14  * @param topk           Number of experts per token (typically 1-4)
15  * @param hidden_size    Size of the hidden dimension
16  * @param out_dtype      Data type of output tensor
17  * @param input_dtype    Data type of input tensor
18  *
19  * @note Memory Layout:
20  *   - input: Contains outputs from topk experts for each token
21  *   - out: Contains the sum of expert outputs for each token
22  *
23  * @note Behavior:
24  *   - For each token, sums the outputs from its topk experts
25  *   - Has optimized implementations for topk=2, 3, and 4
26  *   - Falls back to PyTorch's sum for other topk values
27  */
28 void moe_sum_kernel(
29     void* dev,                // Device handle
30     void* out,                // Output tensor data pointer
31     const void* input,        // Input tensor data pointer
32     const int64_t* out_sizes, // Output tensor sizes
33     const int64_t* out_strides, // Output tensor strides
34     const int64_t* input_sizes, // Input tensor sizes
35     const int64_t* input_strides, // Input tensor strides
36     int out_ndim,            // Number of dimensions in output
tensor
37     int input_ndim,          // Number of dimensions in input tensor
38     int num_tokens,          // Number of tokens to process
39     int topk,                // Number of experts per token
40     int hidden_size,         // Size of the hidden dimension
41     int out_dtype,           // Data type of output tensor
42     int input_dtype,         // Data type of input tensor
43 );

```

ISSUE

- Cuda 优化的 top = 2, 3, 4 的 kernel 中需要保证 input, output 在最后两个维度上是 contiguous 的

- 如果需要处理非连续张量，可以先调用 `contiguous()` 方法或者使用 `at::sum(input, 1)` 代替，但可能牺牲一些性能
- 支持的浮点类型 `float32`、`float16`、`bfloat16`

`_vllm_fa2_C::varlen_fwd[vLLM]`

查看 [国OP 的开发文档](#) 中的 FA2 部分

`vllm::inplace_fused_experts[vLLM]`

参考信息

参考实现：[https://github.com/vllm-](https://github.com/vllm-project/vllm/blob/e92694b6fe264a85371317295bca6643508034ef/vllm/model_executor/layers/used_moe/fused_moe.py#L1119)

[project/vllm/blob/e92694b6fe264a85371317295bca6643508034ef/vllm/model_executor/layers/used_moe/fused_moe.py#L1119](https://github.com/vllm-project/vllm/blob/e92694b6fe264a85371317295bca6643508034ef/vllm/model_executor/layers/used_moe/fused_moe.py#L1119) [https://github.com/vllm-](https://github.com/vllm-project/vllm/blob/e92694b6fe264a85371317295bca6643508034ef/vllm/model_executor/layers/used_moe/fused_moe.py#L623C5-L623C28)
[project/vllm/blob/e92694b6fe264a85371317295bca6643508034ef/vllm/model_executor/layers/used_moe/fused_moe.py#L623C5-L623C28](https://github.com/vllm-project/vllm/blob/e92694b6fe264a85371317295bca6643508034ef/vllm/model_executor/layers/used_moe/fused_moe.py#L623C5-L623C28)

这是一个组合 op，`fused_experts_impl` 函数调用了以下关键操作：

1. `moe_align_block_size`：重排序和对齐 tokens [国OP 的开发文档](#)
2. `invoke_fused_moe_kernel`（两次）：执行高效的矩阵乘法
3. `silu_and_mul`：应用 SiLU 激活函数 [国OP 的开发文档](#)
4. `moe_sum`：合并多个专家的输出 [国OP 的开发文档](#)

`invoke_fused_moe_kernel` 在 cuda 中分为两种 kernel：**`fused_moe_kernel_gptq_awq`** 和 **`fused_moe_kernel`**

`fused_moe_kernel_gptq_awq`

1. 专为 GPTQ 和 AWQ 量化方案设计
2. 支持 INT4（4 位整数）和 INT8（8 位整数）量化的权重
3. 处理更复杂的量化参数，包括零点（zero point）

`fused_moe_kernel`

1. 支持更通用的量化方案
2. 处理 FP8（8 位浮点）和 INT8 量化
3. 更灵活的块大小和分组配置

| | fused_moe_kernel_gptq_awq | fused_moe_kernel |
|--------|---------------------------|-------------------------|
| 量化方案 | INT4, INT8 | FP8, INT8 |
| 零点处理 | 需要额外内存访问 | 不需要 |
| 特殊内存布局 | INT4 打包存储 | FP8 双向量化 |
| 分组策略 | 一维分组 (group_size) | 二维分组 (group_n, group_k) |

运行逻辑:

- 线程索引和输出位置计算
- 专家和 token 映射确定
- 权重加载和反量化
- 输入加载和处理
- 矩阵乘法计算
- 路由权重应用
- 结果存储

Function 接口

fused_moe_kernel_gptq_awq

代码块

```

1  /**
2   * @brief Performs fused MoE computation with GPTQ/AWQ quantized weights
3   *
4   * @param dev          Device handle for CUDA/CPU device management
5   * @param input        Input tensor [num_tokens, hidden_size]
6   * @param weights      Quantized expert weights [num_experts,
7   * @param output       Output tensor [num_tokens_post_pad, topk,
8   * @param weight_scales Weight scaling factors for dequantization
9   * @param weight_zero_points Weight zero points for dequantization (optional)
10  * @param topk_weights Routing weights for each token [num_tokens,
11  * @param topk_ids     Expert IDs for each token [num_tokens, topk]
12  * @param sorted_token_ids Sorted token indices [num_tokens_post_pad]
13  * @param expert_ids   Expert ID for each block [num_blocks]
14  * @param num_tokens_post_pad Total number of tokens after padding
15  * @param input_sizes  Sizes of input tensor
16  * @param input_strides Strides of input tensor
17  * @param weights_sizes Sizes of weights tensor
18  * @param weights_strides Strides of weights tensor

```

19 * @param output_sizes Sizes of output tensor
 20 * @param output_strides Strides of output tensor
 21 * @param weight_scales_sizes Sizes of weight_scales tensor
 22 * @param weight_scales_strides Strides of weight_scales tensor
 23 * @param weight_zero_points_sizes Sizes of weight_zero_points tensor (if not null)
 24 * @param weight_zero_points_strides Strides of weight_zero_points tensor (if not null)
 25 * @param topk_weights_sizes Sizes of topk_weights tensor
 26 * @param topk_weights_strides Strides of topk_weights tensor
 27 * @param topk_ids_sizes Sizes of topk_ids tensor
 28 * @param topk_ids_strides Strides of topk_ids tensor
 29 * @param sorted_token_ids_sizes Sizes of sorted_token_ids tensor
 30 * @param sorted_token_ids_strides Strides of sorted_token_ids tensor
 31 * @param expert_ids_sizes Sizes of expert_ids tensor
 32 * @param expert_ids_strides Strides of expert_ids tensor
 33 * @param input_ndim Number of dimensions in input tensor
 34 * @param weights_ndim Number of dimensions in weights tensor
 35 * @param output_ndim Number of dimensions in output tensor
 36 * @param weight_scales_ndim Number of dimensions in weight_scales tensor
 37 * @param weight_zero_points_ndim Number of dimensions in weight_zero_points tensor
 38 * @param topk_weights_ndim Number of dimensions in topk_weights tensor
 39 * @param topk_ids_ndim Number of dimensions in topk_ids tensor
 40 * @param sorted_token_ids_ndim Number of dimensions in sorted_token_ids tensor
 41 * @param expert_ids_ndim Number of dimensions in expert_ids tensor
 42 * @param num_tokens Number of tokens to process
 43 * @param num_experts Number of experts in the model
 44 * @param hidden_size Size of input hidden dimension
 45 * @param output_size Size of output dimension
 46 * @param topk Number of experts per token
 47 * @param block_size_m Block size for M dimension
 48 * @param block_size_n Block size for N dimension
 49 * @param block_size_k Block size for K dimension
 50 * @param group_size_m Group size for M dimension
 51 * @param group_size Quantization group size
 52 * @param mul_routed_weight Whether to multiply by routing weights
 53 * @param input_dtype Data type of input tensor
 54 * @param weights_dtype Data type of weights tensor
 55 * @param output_dtype Data type of output tensor
 56 * @param weight_scales_dtype Data type of weight_scales tensor
 57 * @param weight_zero_points_dtype Data type of weight_zero_points tensor
 58 * @param use_int4_w4a16 Whether to use INT4 weights with FP16 activations
 59 * @param use_int8_w8a16 Whether to use INT8 weights with FP16 activations
 60 *

```

61  * @note Memory Layout:
62  *      - input: Contains token hidden states
63  *      - weights: Contains quantized expert weights
64  *      - output: Contains expert outputs for each token
65  *
66  * @note Quantization:
67  *      - For INT4: Two 4-bit values are packed into one byte
68  *      - For INT8: Each weight is stored as a single byte
69  *      - Dequantization: (quant_value - zero_point) * scale
70  *
71  * @note Behavior:
72  *      - Performs matrix multiplication between token hidden states and
73  *      - Handles quantized weights with dequantization
74  *      - Optionally multiplies results by routing weights
75  *      - Outputs are arranged by token and expert
76  */
77  void fused_moe_kernel_gptq_awq(
78      void* dev, // Device handle
79      const void* input, // Input tensor data pointer
80      const void* weights, // Weights tensor data pointer
81      void* output, // Output tensor data pointer
82      const void* weight_scales, // Weight scales data pointer
83      const void* weight_zero_points, // Weight zero points data pointer
84      (can be null)
85      const void* topk_weights, // Topk weights data pointer
86      const void* topk_ids, // Topk IDs data pointer
87      const void* sorted_token_ids, // Sorted token IDs data pointer
88      const void* expert_ids, // Expert IDs data pointer
89      const void* num_tokens_post_pad, // Number of tokens post padding
90      const int64_t* input_sizes, // Input tensor sizes
91      const int64_t* input_strides, // Input tensor strides
92      const int64_t* weights_sizes, // Weights tensor sizes
93      const int64_t* weights_strides, // Weights tensor strides
94      const int64_t* output_sizes, // Output tensor sizes
95      const int64_t* output_strides, // Output tensor strides
96      const int64_t* weight_scales_sizes, // Weight scales tensor sizes
97      const int64_t* weight_scales_strides, // Weight scales tensor strides
98      const int64_t* weight_zero_points_sizes, // Weight zero points tensor
99      sizes (if not null)
100     const int64_t* weight_zero_points_strides, // Weight zero points tensor
101     strides (if not null)
102     const int64_t* topk_weights_sizes, // Topk weights tensor sizes
103     const int64_t* topk_weights_strides, // Topk weights tensor strides
104     const int64_t* topk_ids_sizes, // Topk IDs tensor sizes
105     const int64_t* topk_ids_strides, // Topk IDs tensor strides
106     const int64_t* sorted_token_ids_sizes, // Sorted token IDs tensor sizes

```

```

104     const int64_t* sorted_token_ids_strides, // Sorted token IDs tensor strides
105     const int64_t* expert_ids_sizes,       // Expert IDs tensor sizes
106     const int64_t* expert_ids_strides,    // Expert IDs tensor strides
107     int input_ndim,                       // Number of dimensions in input tensor
108     int weights_ndim,                     // Number of dimensions in weights
    tensor
109     int output_ndim,                       // Number of dimensions in output
    tensor
110     int weight_scales_ndim,               // Number of dimensions in weight
    scales tensor
111     int weight_zero_points_ndim,         // Number of dimensions in weight zero
    points tensor
112     int topk_weights_ndim,                // Number of dimensions in topk
    weights tensor
113     int topk_ids_ndim,                    // Number of dimensions in topk IDs
    tensor
114     int sorted_token_ids_ndim,           // Number of dimensions in sorted
    token IDs tensor
115     int expert_ids_ndim,                  // Number of dimensions in expert IDs
    tensor
116     int num_tokens,                       // Number of tokens to process
117     int num_experts,                       // Number of experts in the model
118     int hidden_size,                      // Size of input hidden dimension
119     int output_size,                       // Size of output dimension
120     int topk,                              // Number of experts per token
121     int block_size_m,                      // Block size for M dimension
122     int block_size_n,                      // Block size for N dimension
123     int block_size_k,                      // Block size for K dimension
124     int group_size_m,                      // Group size for M dimension
125     int group_size,                        // Quantization group size
126     bool mul_routed_weight,               // Whether to multiply by routing
    weights
127     int input_dtype,                       // Data type of input tensor
128     int weights_dtype,                     // Data type of weights tensor
129     int output_dtype,                      // Data type of output tensor
130     int weight_scales_dtype,               // Data type of weight scales tensor
131     int weight_zero_points_dtype,         // Data type of weight zero points
    tensor
132     bool use_int4_w4a16,                   // Whether to use INT4 weights with
    FP16 activations
133     bool use_int8_w8a16                   // Whether to use INT8 weights with
    FP16 activations
134 );

```

fused_moe_kernel


```

45 * @param output_size      Size of output dimension
46 * @param topk             Number of experts per token
47 * @param block_size_m    Block size for M dimension
48 * @param block_size_n    Block size for N dimension
49 * @param block_size_k    Block size for K dimension
50 * @param group_size_m    Group size for M dimension
51 * @param group_n         Group size for N dimension in block quantization
52 * @param group_k         Group size for K dimension in block quantization
53 * @param mul_routed_weight Whether to multiply by routing weights
54 * @param input_dtype     Data type of input tensor
55 * @param weights_dtype   Data type of weights tensor
56 * @param output_dtype    Data type of output tensor
57 * @param input_scales_dtype Data type of input_scales tensor
58 * @param weight_scales_dtype Data type of weight_scales tensor
59 * @param use_fp8_w8a8    Whether to use FP8 weights with FP8 activations
60 * @param use_int8_w8a16 Whether to use INT8 weights with FP16
activations
61 *
62 * @note Memory Layout:
63 *     - input: Contains token hidden states
64 *     - weights: Contains expert weights
65 *     - output: Contains expert outputs for each token
66 *
67 * @note Quantization:
68 *     - For FP8: Both weights and activations use 8-bit floating point
69 *     - For INT8: Weights use 8-bit integers, activations use 16-bit
floating point
70 *     - Block quantization: Different scale factors for different blocks
71 *
72 * @note Behavior:
73 *     - Performs matrix multiplication between token hidden states and
expert weights
74 *     - Handles quantized weights and activations
75 *     - Optionally multiplies results by routing weights
76 *     - Outputs are arranged by token and expert
77 */
78 void fused_moe_kernel(
79     void* dev,                // Device handle
80     const void* input,        // Input tensor data pointer
81     const void* weights,      // Weights tensor data pointer
82     void* output,            // Output tensor data pointer
83     const void* input_scales, // Input scales data pointer
84     const void* weight_scales, // Weight scales data pointer
85     const void* topk_weights, // Topk weights data pointer
86     const void* topk_ids,     // Topk IDs data pointer
87     const void* sorted_token_ids, // Sorted token IDs data pointer
88     const void* expert_ids,   // Expert IDs data pointer

```

```

89     const void* num_tokens_post_pad, // Number of tokens post padding
90     const int64_t* input_sizes, // Input tensor sizes
91     const int64_t* input_strides, // Input tensor strides
92     const int64_t* weights_sizes, // Weights tensor sizes
93     const int64_t* weights_strides, // Weights tensor strides
94     const int64_t* output_sizes, // Output tensor sizes
95     const int64_t* output_strides, // Output tensor strides
96     const int64_t* input_scales_sizes, // Input scales tensor sizes
97     const int64_t* input_scales_strides, // Input scales tensor strides
98     const int64_t* weight_scales_sizes, // Weight scales tensor sizes
99     const int64_t* weight_scales_strides, // Weight scales tensor strides
100    const int64_t* topk_weights_sizes, // Topk weights tensor sizes
101    const int64_t* topk_weights_strides, // Topk weights tensor strides
102    const int64_t* topk_ids_sizes, // Topk IDs tensor sizes
103    const int64_t* topk_ids_strides, // Topk IDs tensor strides
104    const int64_t* sorted_token_ids_sizes, // Sorted token IDs tensor sizes
105    const int64_t* sorted_token_ids_strides, // Sorted token IDs tensor strides
106    const int64_t* expert_ids_sizes, // Expert IDs tensor sizes
107    const int64_t* expert_ids_strides, // Expert IDs tensor strides
108    int input_ndim, // Number of dimensions in input tensor
109    int weights_ndim, // Number of dimensions in weights
    tensor
110    int output_ndim, // Number of dimensions in output
    tensor
111    int input_scales_ndim, // Number of dimensions in input
    scales tensor
112    int weight_scales_ndim, // Number of dimensions in weight
    scales tensor
113    int topk_weights_ndim, // Number of dimensions in topk
    weights tensor
114    int topk_ids_ndim, // Number of dimensions in topk IDs
    tensor
115    int sorted_token_ids_ndim, // Number of dimensions in sorted
    token IDs tensor
116    int expert_ids_ndim, // Number of dimensions in expert IDs
    tensor
117    int num_tokens, // Number of tokens to process
118    int num_experts, // Number of experts in the model
119    int hidden_size, // Size of input hidden dimension
120    int output_size, // Size of output dimension
121    int topk, // Number of experts per token
122    int block_size_m, // Block size for M dimension
123    int block_size_n, // Block size for N dimension
124    int block_size_k, // Block size for K dimension
125    int group_size_m, // Group size for M dimension
126    int group_n, // Group size for N dimension in block
    quantization

```

```

127     int group_k, // Group size for K dimension in block
quantization
128     bool mul_routed_weight, // Whether to multiply by routing
weights
129     int input_dtype, // Data type of input tensor
130     int weights_dtype, // Data type of weights tensor
131     int output_dtype, // Data type of output tensor
132     int input_scales_dtype, // Data type of input scales tensor
133     int weight_scales_dtype, // Data type of weight scales tensor
134     bool use_fp8_w8a8, // Whether to use FP8 weights with FP8
activations
135     bool use_int8_w8a16 // Whether to use INT8 weights with
FP16 activations
136 );

```

ISSUE

两个内核都支持张量非contiguous

支持的数据类型:

| | fused_moe_kernel_gptq_awq | fused_moe_kernel |
|--------|---------------------------|-------------------------------------|
| 输入数据类型 | FP32, FP16, BF16 | FP32, FP16, BF16, UINT8 (FP8) |
| 权重数据类型 | INT8, 打包 INT4 | FP32, FP16, BF16, INT8, UINT8 (FP8) |
| 计算精度 | FP32 或 FP16 | FP32 或 FP16 |

vllm::unified_attention[vLLM]

参考信息

源代码: <https://github.com/vllm-project/vllm/blob/e92694b6fe264a85371317295bca6643508034ef/vllm/attention/layer.py#L297>

这是一个用来调用attention layer forward()方法的统一接口

代码块

```

1  def unified_attention(
2      query: torch.Tensor,
3      key: torch.Tensor,
4      value: torch.Tensor,
5      layer_name: str,
6  ) -> torch.Tensor:
7      forward_context: ForwardContext = get_forward_context()
8      attn_metadata = forward_context.attn_metadata
9      self = forward_context.attn_layers[layer_name]

```

```

10 kv_cache = self.kv_cache[forward_context.virtual_engine]
11 return self.impl.forward(self, query, key, value, kv_cache, attn_metadata)

```

不涉及到kernel的开发，起到一个分发到具体attention实现的作用，建议此处使用一个更加高级的API，在分发至不同attention kernel的时候再针对具体attention需要的参数开发对应C++接口内容，比如说这类接口 [OP 的开发文档](#)

Function 接口

代码块

```

1  /**
2   * @brief Attention metadata structure containing all necessary information
   for attention computation
3   */
4  struct AttentionMetadata {
5      // Sequence information
6      torch::Tensor block_tables;           // [num_seqs, max_blocks_per_seq]
7      torch::Tensor context_lens;          // [num_seqs]
8      torch::Tensor max_context_len;       // scalar
9      torch::Tensor seq_lens;              // [num_seqs]
10     torch::Tensor seq_start_loc;          // [num_seqs]
11
12     // Batch information
13     int num_tokens;                       // Total number of tokens
14     int num_seqs;                         // Number of sequences
15     int max_seq_len;                     // Maximum sequence length
16
17     // Attention mechanism
18     torch::Tensor alibi_slopes;           // [num_heads] or empty
19     torch::Tensor attn_bias;              // Optional attention bias
20     torch::Tensor slot_mapping;           // [num_tokens]
21     int sliding_window_size;             // Size of sliding window (-1 if not
used)
22
23     // Quantization
24     bool enable_kv_scale_calc;           // Flag to enable KV scale calculation
25
26     // Prefill/decode information
27     int num_prefills;                     // Number of prefill requests
28     int num_prefill_tokens;              // Number of prefill tokens
29     int num_decode_tokens;               // Number of decode tokens
30 };
31
32 /**
33  * @brief KV Cache structure containing key and value cache tensors

```

```

34  */
35  struct KVCache {
36      torch::Tensor key_cache;           // [num_blocks, block_size,
num_kv_heads, head_size]
37      torch::Tensor value_cache;       // [num_blocks, block_size,
num_kv_heads, head_size]
38      int block_size;                  // Size of each block
39      int max_blocks;                  // Maximum number of blocks
40  };
41
42  /**
43   * @brief Performs attention computation with support for various attention
mechanisms
44   *
45   * @param output          Output tensor to store attention results
[batch_size, seq_len, hidden_size]
46   * @param query          Query tensor [batch_size, seq_len, hidden_size]
47   * @param key             Key tensor [batch_size, seq_len, hidden_size]
48   * @param value          Value tensor [batch_size, seq_len, hidden_size]
49   * @param kv_cache       KV cache structure containing key and value
cache tensors
50   * @param metadata       Attention metadata structure
51   * @param num_heads      Number of attention heads
52   * @param num_kv_heads  Number of key/value heads (for MQA/GQA)
53   * @param head_size     Size of each attention head
54   * @param scale         Scaling factor for attention scores
55   * @param k_scale       Scale factor for key (for FP8 quantization)
56   * @param v_scale       Scale factor for value (for FP8 quantization)
57   * @param attn_type     Type of attention (0=decoder, 1=encoder,
2=encoder_only, 3=encoder_decoder)
58   *
59   * @return torch::Tensor Output tensor with attention results
60   *
61   * @note This is a higher-level interface that uses PyTorch tensors directly
62   * @note It internally calls the lower-level attention_forward function
63   */
64  torch::Tensor attention_forward(
65      torch::Tensor query,           // Query tensor
66      torch::Tensor key,            // Key tensor
67      torch::Tensor value,         // Value tensor
68      const KVCache& kv_cache,     // KV cache structure
69      const AttentionMetadata& metadata, // Attention metadata
70      int num_heads,                // Number of attention heads
71      int num_kv_heads,             // Number of key/value heads
72      int head_size,                // Size of each attention head
73      float scale,                  // Scaling factor for attention scores
74      float k_scale,                // Scale factor for key (for FP8)

```

```

75     float v_scale, // Scale factor for value (for FP8)
76     int attn_type, // Type of attention
77     torch::optional<torch::Tensor> output = torch::nullopt // Optional pre-
    allocated output tensor
78 );

```

ISSUE

不需要kernel

vllm::unified_attention_with_output[vLLM]

参考信息

<https://github.com/vllm-project/vllm/blob/6e1fc61f0fb90c37f0d4a1a8f76235a6e4e1103c/vllm/attention/layer.py#L328>

这是一个用来调用attention layer forward()方法的统一接口，与上面的相比只有一点不同，参数中包含output，可以参考上一个中的高级API

代码块

```

1  def unified_attention_with_output(
2      query: torch.Tensor,
3      key: torch.Tensor,
4      value: torch.Tensor,
5      output: torch.Tensor,
6      layer_name: str,
7  ) -> None:
8      forward_context: ForwardContext = get_forward_context()
9      attn_metadata = forward_context.attn_metadata
10     self = forward_context.attn_layers[layer_name]
11     kv_cache = self.kv_cache[forward_context.virtual_engine]
12     self.impl.forward(self,
13                       query,
14                       key,
15                       value,
16                       kv_cache,
17                       attn_metadata,
18                       output=output)

```

__scaled_mm

参考信息

NV 文档: <https://docs.nvidia.com/cuda/cublas/#id12>

使用场景参考: <https://dev-discuss.pytorch.org/t/float8-in-pytorch-1-x/1815>

Pytorch 源码: 算子 C++地址:

<https://github.com/pytorch/pytorch/blob/0b0d28acdcabfeb3e60cac536749d86138222f8/aten/src/ATen/native/cuda/Blas.cpp#L1365>, cuBlasLt Kernel 调用地址:

<https://github.com/pytorch/pytorch/blob/0b0d28acdcabfeb3e60cac536749d86138222f8/aten/src/ATen/cuda/CUDABlas.cpp#L1513>

我们主要需要支持的格式: kFloat8_e5m2, kFloat8_e4m3fn

需要支持两种 scale: tensor wise&row wise

CUDA Kernel 对 Tensor 要求来源: <https://docs.nvidia.com/cuda/cublas/#id83>

对于 tensor wise scaling 的情况, 参考代码

前期的支持需求:

output_fp32 = A_fp8@B_fp8*a_scale_fp32*b_scale_fp32

out dtype: float32

scale a dtype: float32

scale b dtype: float32

scale res: 初期不需要支持

input a dtype: float8_e4m3fn

input b dtype: float8_e4m3fn

bias: 初期不需要支持

代码块

```
1 import os
2 import torch
3
4 a = torch.eye(16, dtype=torch.float32)
5 a[0, 1] = 2
6 a[1, 0] = 3
7 a[1, 1] = 4
8 b = torch.eye(16, dtype=torch.float32)
9 b[0, 1] = 2
10 b[1, 0] = 3
11 b[1, 1] = 4
12 res_scale = 1.
13 a = a.cuda() / res_scale
14 b = b.cuda() / res_scale
```

```

15
16 c = a.mm(b)[:2,:2]
17
18 """
19 tensor([[ 7., 10.],
20         [15., 22.]], device='cuda:0')
21 """
22 print(c)
23
24 a_scale = a.abs().max()
25 b_scale = b.abs().max()
26
27 a_fp8 = (a / a_scale).to(dtype=torch.float8_e4m3fn)
28 b_fp8 = (b / b_scale).to(dtype=torch.float8_e4m3fn)
29
30 # to column-major
31 b_fp8 = b_fp8.t().contiguous().t()
32
33 c2 = torch._scaled_mm(a_fp8, b_fp8, a_scale, b_scale)[:2,:2]
34 """
35 tensor([[ 7., 10.],
36         [15., 22.]], device='cuda:0', dtype=torch.float8_e4m3fn)
37 """
38 print(c2)
39
40 c3 = torch._scaled_mm(a_fp8, b_fp8, a_scale, b_scale, out_dtype=torch.float32)
41 [:2,:2]
42 """
43 tensor([[ 7., 10.],
44         [15., 22.]], device='cuda:0')
45 """
46 print(c3)
47
48 scale_res = a_scale * b_scale
49
50 c4 = torch._scaled_mm(a_fp8, b_fp8, torch.tensor([1.0]).cuda(),
51 torch.tensor([1.0]).cuda(), scale_result=scale_res)[:2,:2]
52 """
53 tensor([[ 7., 10.],
54         [15., 22.]], device='cuda:0', dtype=torch.float8_e4m3fn)
55 """
56 print(c4)
57
58 bias = torch.tensor([1.0] * 16, dtype=torch.bfloat16).cuda()
59 c5 = torch._scaled_mm(a_fp8, b_fp8, a_scale, b_scale,
60 out_dtype=torch.float8_e4m3fn, bias=bias)[:2,:2]

```

```

59  """
60  tensor([[ 8., 11.],
61         [16., 24.]], device='cuda:0', dtype=torch.float8_e4m3fn)
62  """
63  print(c5)
64
65

```

Tensorwise scaling is enabled when `CUBLASLT_MATMUL_DESC_X_SCALE_MODE` attributes (here `X` stands for `A`, `B`, `C`, `D`, or `AUX`; see `cublasLtMatmulDescAttributes_t`) for all FP8-precision tensors is set to `CUBLASLT_MATMUL_MATRIX_SCALE_SCALAR_32F` (this is the default value for FP8 tensors). In such case, the matmul operation in cuBLAS is defined in the following way (assuming, for exposition, that all tensors are using an FP8 precision):

$$D = scale_D \cdot (\alpha \cdot scale_A \cdot scale_B \cdot \text{op}(A)\text{op}(B) + \beta \cdot scale_C \cdot C).$$

Here A , B , and C are input tensors, and $scale_A$, $scale_B$, $scale_C$, $scale_D$, α , and β are input scalars. This differs from the other matrix multiplication routines because of this addition of scaling factors for each matrix. The $scale_A$, $scale_B$, and $scale_C$ are used for de-quantization, and $scale_D$ is used for quantization. Note that all the scaling factors are applied multiplicatively. This means that sometimes it is necessary to use a scaling factor or its reciprocal depending on the context in which it is applied. For more information on FP8, see `cublasLtMatmul()` and `cublasLtMatmulDescAttributes_t`.

For such matrix multiplications, epilogues and the absolute maximums of intermediate values are computed as follows:

$$\begin{aligned}
 Aux_{temp} &= \alpha \cdot scale_A \cdot scale_B \cdot \text{op}(A)\text{op}(B) + \beta \cdot scale_C \cdot C, \\
 D_{temp} &= \text{Epilogue}(Aux_{temp}), \\
 amax_D &= \text{absmax}(D_{temp}), \\
 amax_{Aux} &= \text{absmax}(Aux_{temp}), \\
 D &= scale_D * D_{temp}, \\
 Aux &= scale_{Aux} * Aux_{temp}.
 \end{aligned}$$

Here Aux is an auxiliary output of matmul consisting of the values that are passed to an epilogue function like GELU, $scale_{Aux}$ is an optional scaling factor that can be applied to Aux , and $amax_{Aux}$ is the maximum absolute value in Aux before scaling. For more information, see attributes `CUBLASLT_MATMUL_DESC_AMAX_D_POINTER` and `CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_AMAX_POINTER` in `cublasLtMatmulDescAttributes_t`.

See the table below when using FP8 kernels:

Table 3. When A, B, C, and D Use Layouts for FP8

| AType | BType | CType | DType | Bias Type |
|----------------|----------------|----------------|----------------|---------------|
| CUDA_R_8F_E4M3 | CUDA_R_8F_E4M3 | CUDA_R_16BF | CUDA_R_16BF | CUDA_R_16BF 8 |
| | | | CUDA_R_8F_E4M3 | CUDA_R_16BF 8 |
| | | CUDA_R_16F | CUDA_R_16F | CUDA_R_16F 8 |
| | | | CUDA_R_8F_E4M3 | CUDA_R_16F 8 |
| | | CUDA_R_32F | CUDA_R_32F | CUDA_R_16BF 8 |
| | | CUDA_R_8F_E5M2 | CUDA_R_16BF | CUDA_R_16BF |
| | CUDA_R_8F_E4M3 | | | CUDA_R_16BF 8 |
| | CUDA_R_8F_E5M2 | | | CUDA_R_16BF 8 |
| | CUDA_R_16F | | CUDA_R_16F | CUDA_R_16F 8 |
| | | | CUDA_R_8F_E4M3 | CUDA_R_16F 8 |
| | | | CUDA_R_8F_E5M2 | CUDA_R_16F 8 |
| | CUDA_R_32F | CUDA_R_32F | CUDA_R_16BF 8 | |
| CUDA_R_8F_E5M2 | CUDA_R_8F_E4M3 | CUDA_R_16BF | CUDA_R_16BF | CUDA_R_16BF 8 |
| | | | CUDA_R_8F_E4M3 | CUDA_R_16BF 8 |
| | | | CUDA_R_8F_E5M2 | CUDA_R_16BF 8 |
| | | CUDA_R_16F | CUDA_R_16F | CUDA_R_16F 8 |
| | | | CUDA_R_8F_E4M3 | CUDA_R_16F 8 |
| | | | CUDA_R_8F_E5M2 | CUDA_R_16F 8 |
| | CUDA_R_32F | CUDA_R_32F | CUDA_R_16BF 8 | |

| | |
|---|--|
| <code>CUBLASLT_EPILOGUE_BIAS = 4</code> | Apply (broadcast) bias from the bias vector. Bias vector length must match matrix D rows, and it must be packed (such as stride between vector elements is 1). Bias vector is broadcast to all columns and added before applying the final postprocessing. |
|---|--|

mm MXFormat

参考信息

实现方案: https://github.com/pytorch/ao/tree/main/torchao/prototype/mx_formats

NV 文档: <https://docs.nvidia.com/cuda/cublas/#id12>

1. MXTensor 继承 Tensor，通过 `__torch_dispatch__` 来处理：

对于 mm

a. `Tensor @ MXTensor -> convert(Tensor->MX) @ check_and_transform_shape([MX,MX]) -> return convert(MX -> Tensor)`

b. `MX @ MX -> check_and_transform_shape([MX,MX]) -> return MX`

注意：

a. `convert (Tensor->MX)` 时，软件需要选择 MXFormat 对应的正确 shape，并且 padding，详情参考：[SiPU 1.0 Tiled Format 架构设计方案 v1.0](#)

b. `check_and_transform_shape` 主要是因为有可能 MXTensor 是由一个 MXTensor B 矩/A 矩阵转置而来，需要 check，但是由于这一步 transform 开销较大，最好有 warning

c. B 矩阵只能采用 M32 的 tile shape

2. 更改 StorageImpl，参考 pytorch_ascend 项目，添加 desc 项，记录实际的内存 alignment 情况，实际内存大小等信息

ISSUE:

对于 Op 实现，是直接注册到 `aten::mm` 上还是创建一个新的算子

直接注册到 `aten::mm` 上不确定这个 op 对应的所有 dispatch key 是否能正确处理 MXTensor

自定义实现的话 AutoGrad 的实现还需要特别处理

mx_to_high_precision/high_precision_to_mx

`mx_ptr` 的 size 可以根据 `tileformat` 文档计算出来

• data size:

◦ `mxint8`: $M * N$

◦ `mxfp6`: $M * N * 3 / 4$

• Header size:

◦ `mxint8`: $M * N / 16$

◦ `mxfp6`: $M * N / 32$

• e.g. `mxfp6`, `supertile(M=32, N=128*4=512)`, data size=12288B, header size=512B, total=12800B

• e.g. `mxint8`, `supertile(M=32, N=32*4=128)`, data size=4096B, header size=256B, total=4352B

目前支持的版本为 `tile to tile`，框架会在上层做 `linear to tile/tile to linear`

代码块

```

1 // MX_T: sifmt::mxint8, sifmt::mxfp6e3m2, sifmt::mxfp6e2m3
2 // HP_T: sifmt::bfloat16, sifmt::float16, sifmt::float32
3 template<typename MX_T, typename HP_T>
4 void mx_to_high_precision(void* mx_ptr, HP_T* hp_ptr, int batch_size, int
dim0, int dim1)
5
6 template<typename MX_T, typename HP_T>
7 void high_precision_to_mx(void* mx_ptr, HP_T* hp_ptr, int batch_size, int
dim0, int dim1)

```

`_empty_affine_quantized`

示例:

代码块

```

1 # PYTHON例子
2 import torch
3
4 # 创建一个空的QTensor
5 a = torch._empty_affine_quantized(
6     size=(3, 4),
7     dtype=torch.quint8,
8     scale=0.02,
9     zero_point=64,
10    device='cpu'
11 )
12
13 print(a)
14
15 """
16 >>> tensor([[ 3.5200,  2.1400,  3.7200, -0.4200],
17            [-0.4000,  1.0600, -1.2800, -1.2800],
18            [-0.8000,  3.3200,  3.3200, -0.5400]], size=(3, 4), dtype=torch.quint8,
19            quantization_scheme=torch.per_tensor_affine, scale=0.02, zero_point=64)
20 """

```

[NOTE] 不需要算子实现

参考实现

Source file: `pytorch/aten/src/ATen/native/quantized/TensorFactories.cpp: 15`

```

// ~~~~~ empty ~~~~~
// We explicitly pass in scale and zero_point because we don't have the infra
// ready to support quantizer in python frontend, once that is ready, we'll
// change to use quantizer
Tensor empty_affine_quantized(
    IntArrayRef size,
    std::optional<ScalarType> dtype,
    std::optional<Layout> layout,
    std::optional<Device> device,
    std::optional<bool> pin_memory,
    double scale,
    int64_t zero_point,
    std::optional<c10::MemoryFormat> optional_memory_format) {
    // See [Note: hacky wrapper removal for TensorOptions]
    TensorOptions options_ = TensorOptions().dtype(dtype).layout(layout).device(device).pinned_memory(pinned_memory: pin_memory);

    TORCH_CHECK(
        !(options_.has_memory_format() && optional_memory_format.has_value()),
        args: "Cannot set memory_format both in TensorOptions and explicit argument; please delete "
        "the redundant setter.");
    auto options: TensorOptions = options_.merge_memory_format(optional_memory_format);
    TORCH_CHECK(
        options.has_dtype(),
        args: "Must provide data type for Tensor creation functions.");
    return new_qtensor(
        You, 3周前 • Uncommitted changes
        sizes: size,
        options,
        quantizer: make_per_tensor_affine_quantizer(
            scale, zero_point, scalar_type: typeMetaToScalarType(dtype: options.dtype()));
    }
}

```

empty_affine_quantized 是 `_empty_affine_quantized` 的内部实现接口，这里调用了 `new_qtensor`，创建一个空 QTensor，QTensor 里保存着创建的 Tensor 以及一个 Quantizer。

Source file: `pytorch/aten/src/ATen/quantized/Quantizer.cpp`: 108

```
inline Tensor new_qtensor(      You, 3周前 • Uncommitted changes
    IntArrayRef sizes,
    const TensorOptions& options,
    QuantizerPtr quantizer) {
    auto memory_format: MemoryFormat = options.memory_format_opt().value_or(u: MemoryFormat::Contiguous);
    auto device: Device = options.device();
    at::Allocator* allocator = nullptr;
    // TODO: why isn't this just using GetAllocator
    if (device.is_cuda()) {
        allocator = at::detail::getCUDAHooks().getCUDADeviceAllocator();
    } else if (device.is_cpu()) {
        allocator = at::getCPUAllocator();
    } else if (device.is_meta()) {
        allocator = GetAllocator(t: kMeta);
    } else if (device.is_privateuseone()) {
        allocator = GetAllocator(t: kPrivateUse1);
    } else {
        TORCH_INTERNAL_ASSERT(0, "unrecognized device for new_qtensor: ", device);
    }

#ifdef USE_PYTORCH_QNNPACK
    if (at::globalContext().qEngine() == at::QEngine::QNNPACK) {
        TORCH_CHECK(!device.is_cuda(), "It looks like you are trying to quantize a CUDA tensor ",
            "while QNNPACK backend is enabled. Although not expected to happen in ",
            "practice, you might have done it for testing purposes. ",
            "Please, either change the quantization engine or move the tensor to a CPU.");
        allocator = c10::GetDefaultMobileCPUAllocator();
    }
#endif

    at::DispatchKey tensorDispatchKey = options.computeDispatchKey();
    native::check_size_nonnegative(size: sizes);
    auto dtype: caffe2::TypeMeta = options.dtype();
    TORCH_CHECK(
        isQIntType(t: typeMetaToScalarType(dtype)),
        "ScalarType ",
        typeMetaToScalarType(dtype),
```

```

#ifdef USE_PYTORCH_QNNPACK
    if (at::globalContext().qEngine() == at::QEngine::QNNPACK) {
        TORCH_CHECK(!device.is_cuda(), "It looks like you are trying to quantize a CUDA tensor ",
                    "while QNNPACK backend is enabled. Although not expected to happen in ",
                    "practice, you might have done it for testing purposes. ",
                    "Please, either change the quantization engine or move the tensor to a CPU.");
        allocator = c10::GetDefaultMobileCPUAllocator();
    }
#endif

at::DispatchKey tensorDispatchKey = options.computeDispatchKey();
native::check_size_nonnegative(size: sizes);
auto dtype: caffe2::TypeMeta = options.dtype();
TORCH_CHECK(
    isQIntType(t: typeMetaToScalarType(dtype)),
    "ScalarType ",
    typeMetaToScalarType(dtype),
    " is not supported in new_qtensor.");
auto scalar_type: ScalarType = typeMetaToScalarType(dtype);
int64_t size_bytes = get_sub_byte_tensor_size(sizes, dtype_itemsize: dtype.itemsize(), t: scalar_type);

auto storage = make_storage_impl(
    use_byte_size: StorageImpl::use_byte_size_t(),
    size_bytes,
    data_ptr: allocator->allocate(n: size_bytes),
    allocator,
    /*resizable=*/true,
    device_opt: device);
auto tensor: Tensor = detail::make_tensor<QTensorImpl>(
    storage, at::DispatchKeySet(k: tensorDispatchKey), dtype, quantizer);
get_qtensorimpl(self: tensor)->set_sizes_contiguous(new_size: sizes);
get_qtensorimpl(self: tensor)->empty_tensor_restride(memory_format);
return tensor;
}

```

is_neg

判断一个Tensor的DispatchKeySet中的Negative位是否被设置。

示例

代码块

```

1  import torch
2
3  a = torch.tensor((3, 4))
4  # 此时Negative bit没有被设置
5  print(torch.is_neg(a))
6
7  """
8  >>> False
9  """
10
11 # 设置一下Negative bit
12 b = torch._C._set_neg(b, True)
13 print(torch.is_neg(a))
14

```

```
15  """
16  >>> True
17  """
```

源码

Source file: pytorch/aten/src/ATen/core/TensorBase.h: 392

```
inline bool is_neg() const {
    return impl_->is_neg();
}
```

Source file: pytorch/c10/core/TensorImpl: 1442

```
/**
 * Whether or not the tensor should be negated
 */
inline bool is_neg() const {
    constexpr auto negative_ks = DispatchKeySet(DispatchKey::Negative);
    return key_set_.has_all(negative_ks);
}
```

【NOTE】不需要实现，功能由pytorch内部完成

is_conj

判断一个Tensor的DispatchKeySet中的Conjugate位是否被设置。

示例

代码块

```
1  import torch
2
3  a = torch.tensor((3, 4))
4  # 此时Conjugate bit没有被设置
5  print(torch.is_conj(a))
6
7  """
8  >>> False
9  """
10
```

```
11 # 设置一下Conjugate bit
12 b = torch._C._set_conj(b, True)
13 print(torch.is_conj(a))
14
15 """
16 >>> True
17 """
```

源码

Source file: pytorch/aten/src/ATen/core/TensorBase.h: 380

```
inline bool is_conj() const {
    return impl_>is_conj();
}
```

Source file: pytorch/c10/core/TensorImpl: 1393

```
/**
 * Whether or not the imaginary part of the tensor should be negated
 */
inline bool is_conj() const {
    constexpr auto conjugate_ks = DispatchKeySet(DispatchKey::Conjugate);
    return key_set_.has_all(conjugate_ks);
}
```

【NOTE】不需要实现，功能由pytorch内部完成

contiguous

参考信息

pytorch文档: <https://docs.pytorch.org/docs/stable/generated/torch.Tensor.contiguous.html>

示例

代码块

```
1 import torch
2
3 # 直接创建的Tensor默认都是连续的
4 a = torch.randn((3, 4))
```

```

5 print(a.is_contiguous())
6
7 """
8 >>> True
9 """
10
11 # 转置操作只会交换size以及重新计算strides, 此时内存不连续
12 # 调用contiguous方法重新生成新的Tensor, 将原来的内存拷贝为连续内存
13 b = a.T
14 print(b.is_contiguous())
15 c = b.contiguous()
16 print(c.is_contiguous())
17
18 """
19 >>> False
20 >>> True
21 """

```

源码

Source file: pytorch/aten/src/ATen/native/TensorProperties.cpp: 115

```

Tensor contiguous(const Tensor& self, MemoryFormat memory_format) {
    if (self.is_contiguous(memory_format)) {
        return self;
    }
    TORCH_CHECK(
        memory_format != MemoryFormat::Preserve,
        "preserve memory format is unsupported by the contiguous operator");
    return self.clone(memory_format);
}

```

Source file: pytorch/aten/src/ATen/native/TensorProperties.cpp: 2131

通过contiguous方法调用clone方法传入的memory_format不会是Preserve, 因此直接看else逻辑, 这里首先创建空Tensor, 并调用copy_方法拷贝src Tensor。

```

Tensor clone(
    const Tensor& src,
    std::optional<c10::MemoryFormat> optional_memory_format) {
    auto memory_format = optional_memory_format.value_or(MemoryFormat::Preserve);
    Tensor self;
    if (memory_format == MemoryFormat::Preserve) {
        if (src.is_non_overlapping_and_dense()) {
            // Copy all strides, this is marginally faster than calling empty_like
            self = at::empty_strided_symint(
                src.sym_sizes(), src.sym_strides(), src.options());
        } else {
            self = at::empty_like(src);
        }
    } else {
        self = at::empty_like(src, src.options(), memory_format);
    }

    if (src._is_zerotensor()) {
        self.zero_();
    } else {
        self.copy_(src);
    }
    return self;
}

```

【NOTE】依赖copy_方法实现。

INT64模拟方案

其他厂商处理方法

1. 寒武纪，直接截断，转为int32或者float类型

9.2. Cambricon PyTorch训练

Q: 使用Cambricon PyTorch训练时，出现类似warning:

```
torch/nn/modules/batchnorm.py:113: UserWarning: cast int64(double) to int(float) implicitly due to known MLU restrictions.
```

A: 算子暂不支持64位。默认情况下，long或int64类型数据会截断为int32，double类型数据会截断为float类型。

用户可设置环境变量来禁止默认的截断：`export CATCH_ALL_OH_64BIT_TO_32BIT_CASTING=0`。如果网络中出现64位数据时，计算会报错，此时，可使用Python函数将64位数据强制转换为32位。

2. 燧原 torch_gcu2.5

7.6. 已知问题

1. 底层软硬件无法支持64位类型，如I64，F64等，torch_gcu会做隐式转换，转为使用32位数据类型。可能会导致精度问题或数值溢出问题。在使用时会输出warning提示信息。
2. 由于1，当tensor为64位类型时，使用tensor.data_ptr接口获取的地址可能有误

3. 摩尔线程GPU自身支持64位数据

- 昇腾在torch_npu源码和架构介绍中可以知道 [【昇腾学院】解密昇腾AI处理器--DaVinci架构（计算单元）-云社区-华为云](#) AI core不支持64位数据类型，会把一些涉及64位计算op放在AICPU(标量计算单元)上op_plugin/ops/base_ops/aclops/ArgsortKernelNpu.cpp · Ascend/op-plugin - Gitee.com，但是也有一些没有明显标识的算子支持int64 op_plugin/ops/base_ops/aclops/NonzeroKernelNpu.cpp · Ascend/op-plugin - 码云 - 开源中国 看不到算子内部实现，无法确认实现方案。
- NV <https://forums.developer.nvidia.com/t/question-about-64-bit-integer-performance/64147> 在这个回答中可以看到，在Ampere架构之前都不原生支持64位宽的数据计算，而是通过软件模拟，拆分为高32和低32位，利用已有的INT32单元分步计算。Ampere之后原生支持64位计算。

SIPU可参考方案

方案1：使用int32模拟int64（类似NV Ampere架构之前，需要在kernel中处理）

用两个int32变量拼接表示一个int64数值，通过算法实现四则运算和比较，问了算子那边可以做，但是比较麻烦，可以利用SIMD特性，但是指令数量会变多，性能会有损耗

方案2：限制数据范围（提前在框架中进行截断）

将 `int64` 张量替换为 `int32`（需确认数值范围不超过 $2^{31}-1$ ）

方案3：使用RV Core计算

类似华为的方案，RV core目前支持64位的scalar计算，使用cpu的逻辑方式编写kernel，但是无法利用SIMD的特性，性能较差。

目前场景分析

int32 最大取值-2147483648 到 2147483647

- bincount：统计每一个layer中被选中的专家编号，专家最大为256，可以控制在int32范围内
- add_func

代码块

```
1 /share_data/users/haoliu2/miniconda3/envs/SIPU/lib/python3.10/site-
  packages/transformers/generation/utils.py
2 if model_kwargs.get("use_cache", True):
3     model_kwargs["cache_position"] = model_kwargs["cache_position"]
  [-1:] + num_new_tokens
4 else:
5     past_positions = model_kwargs.pop("cache_position")
6     new_positions = torch.arange(
```

```

7         past_positions[-1] + 1, past_positions[-1] + num_new_tokens +
          1, dtype=past_positions.dtype
8     ).to(past_positions.device)
9     model_kwargs["cache_position"] = torch.cat((past_positions,
          new_positions))
10    return model_kwargs

```

`model_kwargs["cache_position"]` 和 `num_new_tokens` 中保存的数据最大值取决于模型的最大上下文长度 `config.max_position_embeddings` 与模型支持的上下文长度相关，目前支持上下文长度最大的模型是 Gemini，达到百万级别，完全在 `int32` 的范围以内，不存在溢出风险

3. mul_func:

代码块

```

1 /share_data/users/haoliu2/miniconda3/envs/SIPU/lib/python3.10/site-
  packages/transformers/generation/utils.py
2 # finished sentences should have their next token be a padding token
3 if has_eos_stopping_criteria:
4     next_tokens = next_tokens * unfinished_sequences + pad_token_id * (1 -
          unfinished_sequences)

```

- 对未完成的序列，保留模型预测的下一个令牌（`next_tokens`）。
- 对已完成的序列，用填充令牌（`pad_token_id`）替换预测结果。

理论最大值：`max(vocab_size - 1, pad_token_id)`

- GPT-2: 词汇表大小 50,257 → 最大令牌ID = 50,256
- LLaMA-2: 词汇表大小 32,000 → 最大令牌ID = 31,999
- BERT: 词汇表大小 30,522 → 最大令牌ID = 30,521
- DS V3: 词汇表大小 129,280 → 最大令牌ID = 129,279

`pad_token_id` 通常较小（如 0），不会显著影响最大值

可以分析得出 `max` 取值在 `int32` 的范围以内，不存在溢出风险

目前推荐在框架侧进行处理，转换到 `int32` 来处理

