

HW3

何重阳 学号：2025211210

2026 年 1 月 6 日

目录

1	算子实现	3
1.1	Add	3
1.1.1	Add 算子算法设计	3
1.1.2	核心代码实现	3
1.2	Linear	4
1.2.1	Linear 算子算法设计	4
1.2.2	核心代码实现	4
1.2.3	linear.py 运行结果	5
1.3	Argmax	6
1.3.1	Argmax 算子算法设计	6
1.3.2	核心代码实现	6
1.3.3	argmax.py 运行结果	7
1.4	Embedding	7
1.4.1	Embedding 算子算法设计	7
1.4.2	核心代码实现	7
1.4.3	embedding.py 运行结果	8
1.5	Rope	8
1.5.1	Rope 算子算法设计	8
1.5.2	核心代码实现	9
1.5.3	rope.py 运行结果	11
1.6	self_attention	11
1.6.1	self_attention 算子算法设计	11
1.6.2	核心代码实现	11
1.6.3	self_attention.py 运行结果	13

目录	2
1.7 swiglu	14
1.7.1 swiglu 算子算法设计	14
1.7.2 核心代码实现	14
1.7.3 swiglu.py 运行结果	15
1.8 rms_norm	15
1.8.1 rms_norm 算子算法设计	15
1.8.2 核心代码实现	16
1.8.3 rms_norm.py 运行结果	17
2 推理框架搭建	17
2.1 test_infer.py 运行结果展示	17
3 性能分析与算子优化	18
3.1 基础版分析	18
3.2 改进 1	18
3.3 改进 2	20
3.4 改进 3	21
A Triton Kernels for Quantized Linear Operations	22

1 算子实现

1.1 Add

1.1.1 Add 算子算法设计

Add 算子是基于 Triton 框架开发的 GPU 并行化逐元素加法算子，核心功能为完成两个任意形状张量的逐元素加法运算 ($output = input + other$)。

1.1.2 核心代码实现

Add 算子的核心代码如下：

```
1      import triton
2      import triton.language as tl
3
4      @triton.jit
5      def add_kernel(
6          input_ptr, other_ptr, output_ptr,
7          n_elements,
8          BLOCK_SIZE: tl.constexpr, # 编译期常量：控制并行块大小
9      ):
10         # 1. 解析程序ID，计算当前块的起始位置
11         pid = tl.program_id(axis=0)
12         block_start = pid * BLOCK_SIZE
13         # 2. 生成当前块内的元素偏移量
14         offsets = block_start + tl.arange(0, BLOCK_SIZE)
15         # 3. 生成掩码：避免访问超出张量长度的元素
16         mask = offsets < n_elements
17
18         # 4. 加载输入数据（全局内存 → GPU共享内存）
19         input = tl.load(input_ptr + offsets, mask=mask)
20         other = tl.load(other_ptr + offsets, mask=mask)
21         # 5. 核心计算：逐元素加法
22         output = input + other
23         # 6. 存储结果到输出张量
24         tl.store(output_ptr + offsets, output, mask=mask)
25
26
```

Listing 1: 基于 Triton 的 Add 算子核心实现代码

1.2 Linear

1.2.1 Linear 算子算法设计

本文基于 Triton 框架实现了适配大/小矩阵场景的 Linear 算子，核心完成全连接层计算 $Y = X \cdot W^T + b$ ($X \in \mathbb{R}^{M \times K}$, $W \in \mathbb{R}^{N \times K}$, $b \in \mathbb{R}^N$)。

算子核心执行流程如下：

1. 解析输入维度 M, N, K ;
2. 计算网格维度 $\text{grid} = \lceil M/\text{BLOCK}_M \rceil \times \lceil N/\text{BLOCK}_N \rceil$ ，每个网格对应一个输出子块；
3. 核函数内解析程序 ID (pid)，划分行/列块索引，计算块偏移量；
4. K 维度分块循环：加载输入/权重子块，按场景执行矩阵乘法并累加；
5. (可选) 加载偏置并广播到累加结果，掩码存储最终结果到输出矩阵。

1.2.2 核心代码实现

关键逻辑如下：

```

1      @triton.jit
2      def linear_kernel_large(
3          input_ptr, weight_ptr, bias_ptr, output_ptr,
4          M, N, K, # 维度参数
5          # 步长/块大小等constexpr参数 (省略)
6          BLOCK_SIZE_M: tl.constexpr, BLOCK_SIZE_N: tl.constexpr,
7          BLOCK_SIZE_K: tl.constexpr,
8      ):
9          # 1. 块索引与偏移计算 (核心流程)
10         pid = tl.program_id(axis=0)
11         num_pid_n = tl.cdiv(N, BLOCK_SIZE_N)
12         pid_m = pid // num_pid_n # 行块索引
13         pid_n = pid % num_pid_n # 列块索引
14         offs_m = pid_m * BLOCK_SIZE_M + tl.arange(0, BLOCK_SIZE_M)
15         offs_n = pid_n * BLOCK_SIZE_N + tl.arange(0, BLOCK_SIZE_N)
16         offs_k = tl.arange(0, BLOCK_SIZE_K)
17
18         # 2. 累加器初始化
19         accumulator = tl.zeros([BLOCK_SIZE_M, BLOCK_SIZE_N], dtype=tl.
float32)

```

```

20     # 3. K维度分块循环（核心：tl.dot矩阵乘法）
21     for k in range(0, K, BLOCK_SIZE_K):
22         k_offs = k + offs_k
23         # 越界掩码（保证内存安全）
24         a_mask = (offs_m[:, None] < M) & (k_offs[None, :] < K)
25         b_mask = (offs_n[None, :] < N) & (k_offs[:, None] < K)
26
27         # 加载子块并转换精度
28         a = tl.load(input_ptr + offs_m[:, None]*stride_im + k_offs[
None, :]*stride_in, mask=a_mask, other=0.0)
29         b = tl.load(weight_ptr + offs_n[None, :]*stride_wn + k_offs
[:, None]*stride_wk, mask=b_mask, other=0.0)
30         a_f32, b_f32 = a.to(tl.float32), b.to(tl.float32)
31
32         # ===== 核心：tl.dot实现子块矩阵乘法 =====
33         accumulator += tl.dot(a_f32, b_f32, allow_tf32=False)
34
35         # 偏置添加+结果存储（省略重复代码）
36         if bias_ptr is not None:
37             bias = tl.load(bias_ptr + offs_n, mask=offs_n < N, other=0.0)
38             accumulator += bias[None, :]
39         tl.store(output_ptr + offs_m[:, None]*stride_om + offs_n[None,
:]*stride_on, accumulator, mask=(offs_m[:, None]<M) & (offs_n[None,
:]<N))
40

```

Listing 2: Linear 算子核心逻辑

1.2.3 linear.py 运行结果

```

root@autodl-container-fcb946b0f6-6ca1dbb1:~/autodl-tmp/大模型计算_HW3/llaisys-thu# bash test_ops.sh
Testing Ops.linear on nvidia
out (2, 3), x (2, 4), w (3, 4), bias True, dtype <f32>
out (2, 3), x (2, 4), w (3, 4), bias True, dtype <f16>
out (2, 3), x (2, 4), w (3, 4), bias True, dtype <bf16>
out (512, 4096), x (512, 4096), w (4096, 4096), bias True, dtype <f32>
out (512, 4096), x (512, 4096), w (4096, 4096), bias True, dtype <f16>
out (512, 4096), x (512, 4096), w (4096, 4096), bias True, dtype <bf16>
Test passed!

```

图 1: Linear 算子运行结果

1.3 Argmax

1.3.1 Argmax 算子算法设计

该算子根据输入规模自动选择不同的执行路径：当张量长度较小时，使用单个 block 在一次 kernel 调用中完成最大值及其索引的计算；当张量规模较大时，则在 kernel 内部采用分块迭代的方式顺序扫描输入数据，从而避免多 kernel 启动和中间归约缓冲区带来的额外开销。

在实现层面，算子充分利用 Triton 提供的向量化加载与归约原语，通过 `tl.max` 计算局部最大值，并结合布尔掩码与 `tl.min` 操作获取最大值的首个出现位置。为保证数值稳定性，所有参与比较的中间结果统一提升至 `float32` 精度；同时，在 masked load 场景下使用负无穷作为填充值，以避免越界元素对归约结果产生干扰。该设计在保持实现简洁性的同时，实现了良好的性能与可扩展性，并能够与 LLAISys 张量表示无缝集成。

1.3.2 核心代码实现

```

1      @triton.jit
2      def argmax_kernel_single(
3          max_idx_ptr,      # 输出：最大值索引
4          max_val_ptr,      # 输出：最大值
5          vals_ptr,         # 输入：数据指针
6          n_elements,       # 输入：元素个数
7          BLOCK_SIZE: tl.constexpr,
8      ):
9          # 生成 block 内偏移
10         offsets = tl.arange(0, BLOCK_SIZE)
11         mask = offsets < n_elements
12
13         # 向量化加载数据，越界位置填充为 -inf
14         vals = tl.load(vals_ptr + offsets, mask=mask, other=float('-
15         inf'))
16
17         # block 内归约求最大值（返回 float32）
18         local_max = tl.max(vals, axis=0)
19
20         # 查找最大值的首个出现位置
21         is_max = (vals == local_max) & mask
22         local_idx = tl.where(is_max, offsets, n_elements)
23         local_idx = tl.min(local_idx, axis=0)

```

```

24         # 写回结果
25         tl.store(max_val_ptr, local_max)
26         tl.store(max_idx_ptr, local_idx)
27

```

Listing 3: Triton Argmax Kernel 核心实现

1.3.3 argmax.py 运行结果

```

root@autodl-container-fcb946b0f6-6ca1dbb1:~/autodl-tmp/大模型计算_HW3/llaisys-thu# bash test_ops.sh
Testing Ops.argmax on nvidia
shape (4,) dtype <f32>
shape (4,) dtype <f16>
shape (4,) dtype <bf16>
shape (4096,) dtype <f32>
shape (4096,) dtype <f16>
shape (4096,) dtype <bf16>
Test passed!

```

图 2: argmax 算子运行结果

1.4 Embedding

1.4.1 Embedding 算子算法设计

该 Triton embedding kernel 采用按序列分块的并行策略，将索引到向量的查表操作显式展开为内存加载与存储过程。每个 program 负责一个连续的 token block，并对词表索引进行合法性检查后完成对应嵌入向量的拷贝，从而实现对变长输入的安全支持。

1.4.2 核心代码实现

```

1      @triton.jit
2      def embedding_kernel(
3          weight_ptr, index_ptr, output_ptr,
4          vocab_size, hidden_size, seq_len,
5          BLOCK_SIZE: tl.constexpr,
6      ):
7          # program id, 对应序列维度上的一个 block
8          pid = tl.program_id(axis=0)
9          block_start = pid * BLOCK_SIZE
10
11         # block 内的 token 偏移
12         offsets = block_start + tl.arange(0, BLOCK_SIZE)
13         mask = offsets < seq_len
14

```

```

15     # 加载索引（越界位置填充为 0）
16     indices = tl.load(index_ptr + offsets, mask=mask, other=0)
17
18     # 逐 token 处理 embedding 查表
19     for i in range(BLOCK_SIZE):
20         if i + block_start < seq_len:
21             # 读取当前 token 的词表索引
22             idx = tl.load(index_ptr + block_start + i)
23
24             # 索引合法性检查
25             if idx >= 0 and idx < vocab_size:
26                 weight_offset = idx * hidden_size
27                 output_offset = (block_start + i) * hidden_size
28
29                 # 将对应 embedding 向量逐元素拷贝到输出
30                 for j in range(hidden_size):
31                     val = tl.load(weight_ptr + weight_offset + j)
32                     tl.store(output_ptr + output_offset + j, val)
33

```

Listing 4: Triton Embedding Kernel 核心实现

1.4.3 embedding.py 运行结果

```

root@autodl-container-fcb946b0f6-6ca1dbb1:~/autodl-tmp/大模型计算_HW3/llaisys-thu# bash test_ops.sh
Testing Ops.embedding on nvidia
idx_shape (1,) embd_shape (2, 3) dtype <f32>
idx_shape (1,) embd_shape (2, 3) dtype <f16>
idx_shape (1,) embd_shape (2, 3) dtype <bf16>
idx_shape (50,) embd_shape (512, 4096) dtype <f32>
idx_shape (50,) embd_shape (512, 4096) dtype <f16>
idx_shape (50,) embd_shape (512, 4096) dtype <bf16>
Test passed!

```

图 3: embedding 算子运行结果

1.5 Rope

1.5.1 Rope 算子算法设计

RoPE 算子以 (head, position) 作为并行粒度，每个 Triton program 负责处理一个注意力头在某一序列位置上的旋转计算。特征维度 d 被等分为两部分，对应旋转变换中的实部与虚部。对于每一对特征分量 (a, b) ，算子根据位置索引 p 和频率参数 θ 计算旋转角度，并

执行如下变换：

$$a' = a \cos(\omega_p) - b \sin(\omega_p),$$

$$b' = b \cos(\omega_p) + a \sin(\omega_p),$$

其中频率 ω_p 由位置索引与维度相关的幂函数共同确定。

1.5.2 核心代码实现

```

1      @triton.jit
2      def rope_kernel(
3          input_ptr, pos_ids_ptr, output_ptr,
4          seq_len, n_head, d, theta,
5          stride_im, stride_in, stride_ih,
6          stride_om, stride_on, stride_oh,
7          BLOCK_SIZE: tl.constexpr,
8      ):
9          # 每个 program 对应一个 (head, position)
10         pid = tl.program_id(axis=0)
11         head_id = pid // seq_len
12         pos = pid % seq_len
13
14         # 越界检查
15         if head_id >= n_head or pos >= seq_len:
16             return
17
18         # 读取位置索引
19         pos_id = tl.load(pos_ids_ptr + pos)
20         d_half = d // 2
21
22         # 按 block 遍历特征维度
23         for i in range(0, d_half, BLOCK_SIZE):
24             idx = i + tl.arange(0, BLOCK_SIZE)
25             mask = idx < d_half
26
27             # 实部与虚部索引
28             a_idx = idx
29             b_idx = idx + d_half
30
31             # 计算输入偏移

```

```

32     a_offset = pos * stride_im + head_id * stride_in + a_idx *
stride_ih
33     b_offset = pos * stride_im + head_id * stride_in + b_idx *
stride_ih
34
35     # 加载输入
36     a = tl.load(input_ptr + a_offset, mask=mask, other=0.0)
37     b = tl.load(input_ptr + b_offset, mask=mask, other=0.0)
38
39     # 转换为 float32 进行旋转计算
40     a = a.to(tl.float32)
41     b = b.to(tl.float32)
42
43     # 使用 float64 计算高精度频率
44     pos_id_f64 = pos_id.to(tl.float64)
45     idx_f64 = idx.to(tl.float64)
46     exponent = 2.0 * idx_f64 / d
47
48     #  $\theta^{\text{exponent}} = 2^{(\text{exponent} * \log_2(\theta))}$ 
49     theta_f64 = theta + 0.0
50     log2_theta = tl.log2(theta_f64).to(tl.float64)
51     log2_power = exponent * log2_theta
52     theta_power = tl.exp2(log2_power)
53
54     freqs_f64 = pos_id_f64 / theta_power
55     angle = freqs_f64.to(tl.float32)
56
57     cos_val = tl.cos(angle)
58     sin_val = tl.sin(angle)
59
60     # 应用 RoPE 旋转变换
61     a_out = a * cos_val - b * sin_val
62     b_out = b * cos_val + a * sin_val
63
64     # 输出偏移
65     a_out_offset = pos * stride_om + head_id * stride_on + a_idx *
stride_oh
66     b_out_offset = pos * stride_om + head_id * stride_on + b_idx *
stride_oh

```

```

67
68     # 写回结果
69     tl.store(output_ptr + a_out_offset, a_out, mask=mask)
70     tl.store(output_ptr + b_out_offset, b_out, mask=mask)
71

```

Listing 5: Triton RoPE Kernel 核心实现

1.5.3 rope.py 运行结果

```

root@autodl-container-fcb946b0f6-6ca1dbb1:~/autodl-tmp/大模型计算_HW3/llaisys-thu# bash test_ops.sh
Testing Ops.rope on nvidia
  shape (2, 1, 4) range (0, 2) dtype <f32>
  shape (2, 1, 4) range (0, 2) dtype <f16>
  shape (2, 1, 4) range (0, 2) dtype <bf16>
  shape (512, 4, 4096) range (512, 1024) dtype <f32>
  shape (512, 4, 4096) range (512, 1024) dtype <f16>
  shape (512, 4, 4096) range (512, 1024) dtype <bf16>
Test passed!

```

图 4: rope 算子运行结果

1.6 self_attention

1.6.1 self_attention 算子算法设计

Key 和 Value 序列沿序列维度按 BLOCK_SIZE_N 分块, 每个 block 依次参与与当前 Query block 的注意力计算。对于注意力分数的计算, 当特征维度较大时优先使用 `tl.dot` 原语, 并启用 IEEE 精度以提升数值稳定性; 在特征维度较小时, 则回退至显式逐元素乘加实现, 以降低 kernel 启动与调度开销。

为支持自回归生成场景, 算子在计算过程中同时引入边界掩码与因果 (causal) 掩码, 确保每个 Query 位置只能关注其当前位置之前的 Key。掩码后的注意力分数被设置为负无穷, 从而在 Softmax 计算中自然衰减为零权重。

在 Softmax 实现上, 该算子采用 online softmax 技术, 通过维护逐 block 更新的最大值 m_i 和归一化因子 d_i , 避免显式存储完整注意力矩阵, 有效降低显存占用并提升数值稳定性。所有中间统计量均在 float32 精度下进行更新, 以平衡计算精度与性能。

1.6.2 核心代码实现

```

1     for kv_block_start in range(0, kv_len, BLOCK_SIZE_N):
2         # 加载 Key / Value block, 并转换为 float32
3         k = tl.load(k_block_ptr, boundary_check=(0, 1)).to(tl.float32)
4         v = tl.load(v_block_ptr, boundary_check=(0, 1)).to(tl.float32)
5
6         # 计算注意力分数 QK^T

```

```

7      # 当维度较大时, 使用 tl.dot 以获得更好的数值精度
8      if USE_DOT:
9          attention_scores = tl.dot(q, k, input_precision="ieee") * scale
10     else:
11         q_expanded = q[:, :, None]
12         k_expanded = k[None, :, :]
13         product = q_expanded * k_expanded
14         product = tl.where(q_mask[None, :, None], product, 0.0)
15         attention_scores = tl.sum(product, axis=1) * scale
16
17     # 构造边界掩码
18     kv_idx = kv_block_start + tl.arange(0, BLOCK_SIZE_N)
19     kv_mask = kv_idx < kv_len
20
21     # 构造因果掩码 (causal mask)
22     causal_limits = (
23         BLOCK_SIZE_M * pid_m
24         + tl.arange(0, BLOCK_SIZE_M)
25         + (kv_len - q_len)
26         + 1
27     )
28     casual_mask = kv_idx[None, :] < causal_limits[:, None]
29
30     mask_m = (BLOCK_SIZE_M * pid_m + tl.arange(0, BLOCK_SIZE_M)) <
31     q_len
32
33     # 合并所有掩码
34     full_mask = mask_m[:, None] & kv_mask[None, :] & casual_mask
35
36     # 被掩码位置设置为 -inf, 使 Softmax 权重为 0
37     masked_attention_scores = tl.where(
38         full_mask, attention_scores, float("-inf")
39     )
40
41     # Online Softmax 更新 (float32)
42     m_i_new = tl.maximum(m_i, tl.max(masked_attention_scores, axis=1)
43 )
44
45     alpha = tl.exp2(log2_e * (m_i - m_i_new))
46     exp_attention_scores = tl.exp2(

```

```

44     log2_e * (masked_attention_scores - m_i_new[:, None])
45 )
46 d_i_k = tl.sum(exp_attention_scores, axis=1)
47
48 # 计算加权 Value 聚合
49 if USE_DOT and dv >= 16 and BLOCK_SIZE_M >= 16:
50     o = tl.dot(exp_attention_scores, v, input_precision="ieee")
51 else:
52     exp_scores_expanded = exp_attention_scores[:, :, None]
53     v_expanded = v[None, :, :]
54     o = tl.sum(exp_scores_expanded * v_expanded, axis=1)
55
56 # 更新累积输出
57 acc = acc * alpha[:, None] + o
58
59 # 更新 Softmax 统计量
60 m_i = m_i_new
61 d_i = d_i * alpha + d_i_k
62
63 # 指针前移, 处理下一个 KV block
64 k_block_ptr = tl.advance(k_block_ptr, (0, BLOCK_SIZE_N))
65 v_block_ptr = tl.advance(v_block_ptr, (BLOCK_SIZE_N, 0))
66

```

Listing 6: Triton Self-Attention Kernel 核心计算逻辑

1.6.3 self_attention.py 运行结果

```

root@autodl-container-fcb946b0f6-6ca1dbb1:~/autodl-tmp/大模型计算_HW3/llaisys-thu# bash test_ops.sh
Testing Ops.self_attention on nvidia
qlen=2 kvlen=2 nh=1 nkvh=1 hd=4 dtype <f32>
qlen=2 kvlen=2 nh=1 nkvh=1 hd=4 dtype <f16>
qlen=2 kvlen=2 nh=1 nkvh=1 hd=4 dtype <bf16>
qlen=5 kvlen=11 nh=4 nkvh=2 hd=8 dtype <f32>
qlen=5 kvlen=11 nh=4 nkvh=2 hd=8 dtype <f16>
qlen=5 kvlen=11 nh=4 nkvh=2 hd=8 dtype <bf16>
qlen=1 kvlen=1 nh=4 nkvh=2 hd=8 dtype <f32>
qlen=1 kvlen=1 nh=4 nkvh=2 hd=8 dtype <f16>
qlen=1 kvlen=1 nh=4 nkvh=2 hd=8 dtype <bf16>
Test passed!

```

图 5: self_attention 算子运行结果

1.7 swiglu

1.7.1 swiglu 算子算法设计

SwiGLU 算子以序列维度为并行轴，每个 Triton program 负责处理一个 token 对应的一整行中间特征。特征维度 (`intermediate_size`) 在 kernel 内部按 `BLOCK_SIZE` 分块处理，从而在保证并行度的同时支持任意长度的中间特征维度。算子显式接收输入与输出张量的 `stride` 信息，能够正确处理非连续内存布局。

在数值计算方面，算子在 kernel 内部将输入提升至 `float32` 精度完成 Sigmoid 与乘法运算，以对齐 PyTorch 中 `gate.float()` 的计算语义并提升数值稳定性，最终再转换回原始数据类型进行存储。

1.7.2 核心代码实现

```
1      @triton.jit
2      def swiglu_kernel(
3          gate_ptr, up_ptr, output_ptr,
4          seq_len, intermediate_size,
5          stride_gm, stride_gn,
6          stride_um, stride_un,
7          stride_om, stride_on,
8          BLOCK_SIZE: tl.constexpr,
9      ):
10         # 每个 program 负责一个序列位置
11         pid = tl.program_id(axis=0)
12         row = pid
13
14         # 当前 block 处理的列索引
15         cols = tl.arange(0, BLOCK_SIZE)
16         mask = cols < intermediate_size
17
18         # 越界检查
19         if row >= seq_len:
20             return
21
22         # 计算输入与输出的内存偏移
23         gate_offset = row * stride_gm + cols * stride_gn
24         up_offset = row * stride_um + cols * stride_un
25         output_offset = row * stride_om + cols * stride_on
26
```

```

27     # 加载 gate 与 up 分支数据
28     gate = tl.load(gate_ptr + gate_offset, mask=mask, other=0.0)
29     up = tl.load(up_ptr + up_offset, mask=mask, other=0.0)
30
31     # 将输入转换为 float32 以提升数值精度
32     gate_f32 = gate.to(tl.float32)
33     up_f32 = up.to(tl.float32)
34
35     # 计算 sigmoid(gate)
36     sigmoid_gate = 1.0 / (1.0 + tl.exp(-gate_f32))
37
38     # SwiGLU: up * gate * sigmoid(gate)
39     result_f32 = up_f32 * gate_f32 * sigmoid_gate
40
41     # 转回原始数据类型并写回输出
42     output = result_f32.to(gate.dtype)
43     tl.store(output_ptr + output_offset, output, mask=mask)
44

```

Listing 7: Triton SwiGLU Kernel 核心实现

1.7.3 swiglu.py 运行结果

```

root@autodl-container-fcb946b0f6-6ca1dbb1:~/autodl-tmp/大模型计算_HW3/llaisys-thu# bash test_ops.sh
Testing Ops.swiglu on nvidia
shape (2, 3) dtype <f32>
shape (2, 3) dtype <f16>
shape (2, 3) dtype <bf16>
shape (512, 4096) dtype <f32>
shape (512, 4096) dtype <f16>
shape (512, 4096) dtype <bf16>
Test passed!

```

图 6: swiglu 算子运行结果

1.8 rms_norm

1.8.1 rms_norm 算子算法设计

RMSNorm 算子以序列维度为并行轴，每个 Triton program 负责处理一个 token 对应的一整行特征向量。由于特征维度通常较大，算子在 kernel 内部采用两阶段（two-pass）计算策略：第一阶段以 block 为单位遍历特征维度，累积该行输入的平方和；第二阶段在得到归一化因子后，再次遍历特征维度完成归一化与缩放操作。该设计避免了额外的中间张量分配，同时能够支持任意长度的特征维度。

1.8.2 核心代码实现

```

1      @triton.jit
2      def rms_norm_kernel(
3          input_ptr, weight_ptr, output_ptr,
4          seq_len, hidden_size, eps,
5          stride_im, stride_in,
6          stride_om, stride_on,
7          BLOCK_SIZE: tl.constexpr,
8      ):
9          # 每个 program 负责一个 token (一行特征)
10         row = tl.program_id(axis=0)
11
12         # ----- 第一阶段: 累积平方和 -----
13         sum_squared = 0.0
14         for block_start in range(0, hidden_size, BLOCK_SIZE):
15             cols = tl.arange(0, BLOCK_SIZE) + block_start
16             mask = cols < hidden_size
17
18             input_offsets = row * stride_im + cols * stride_in
19             x = tl.load(input_ptr + input_offsets, mask=mask, other=0.0)
20
21             x_squared = x * x
22             sum_squared += tl.sum(x_squared, axis=0)
23
24         # 计算均方根的倒数
25         mean_squared = sum_squared / hidden_size
26         mean_squared_eps = mean_squared + eps
27         rsqrt_val = tl.rsqrt(mean_squared_eps)
28
29         # ----- 第二阶段: 归一化并缩放 -----
30         for block_start in range(0, hidden_size, BLOCK_SIZE):
31             cols = tl.arange(0, BLOCK_SIZE) + block_start
32             mask = cols < hidden_size
33
34             input_offsets = row * stride_im + cols * stride_in
35             output_offsets = row * stride_om + cols * stride_on
36
37             x = tl.load(input_ptr + input_offsets, mask=mask, other=0.0)

```



```

38     weight = tl.load(weight_ptr + cols, mask=mask, other=0.0)
39
40     # RMSNorm: x / sqrt(mean(x^2) + eps) * weight
41     x_norm = x * rsqrt_val
42     output = x_norm * weight
43
44     tl.store(output_ptr + output_offsets, output, mask=mask)
45

```

Listing 8: Triton RMSNorm Kernel 核心实现

1.8.3 rms_norm.py 运行结果

```

root@autodl-container-fcb946b0f6-6ca1dbb1:~/autodl-tmp/大模型计算_HW3/llaisys-thu# bash test_ops.sh
Testing Ops.rms_norm on nvidia
  shape (1, 4) dtype <f32>
  shape (1, 4) dtype <f16>
  shape (1, 4) dtype <bf16>
  shape (512, 4096) dtype <f32>
  shape (512, 4096) dtype <f16>
  shape (512, 4096) dtype <bf16>
Test passed!

```

图 7: rms_norm 算子运行结果

2 推理框架搭建

2.1 test_infer.py 运行结果展示

```

Time elapsed: 1.82s

=== Your Result ===

Tokens:
[151646, 151646, 151644, 15191, 525, 498, 30, 151645, 151648, 198, 91786, 0, 358, 2776, 18183, 39350, 10911, 16, 11,
458, 20443, 11229, 17847, 3465, 553, 18183, 39350, 13, 358, 2776, 518, 697, 2473, 323, 1035, 387, 33972, 311, 7789,
498, 448, 894, 43883, 476, 9079, 498, 1231, 614, 624, 151649, 271, 91786, 0, 358, 2776, 18183, 39350, 10911, 16, 11
, 458, 20443, 11229, 17847, 3465, 553, 18183, 39350, 13, 358, 2776, 518, 697, 2473, 323, 1035, 387, 33972, 311, 7789
, 498, 448, 894, 43883, 476, 9079, 498, 1231, 614, 13, 151643]

Contents:
<| User| >Who are you?<| Assistant| ><think>
Greetings! I'm DeepSeek-R1, an artificial intelligence assistant created by DeepSeek. I'm at your service and would
be delighted to assist you with any inquiries or tasks you may have.
</think>

Greetings! I'm DeepSeek-R1, an artificial intelligence assistant created by DeepSeek. I'm at your service and would
be delighted to assist you with any inquiries or tasks you may have.

Time elapsed: 8.17s
Test passed!

```

图 8: test_infer.py 运行结果

3 性能分析与算子优化

3.1 基础版分析

在未优化的版中，test_infer.py 的运行时间为 **8.17 s**。实验表明，linear_kernel_large、cudaFree、cudaMemcpy、以及 cudaLaunchKernel 等操作的耗时较长、时间占比较大。因此，可以针对这几个操作进行优化。

Time	Total Time	Instances	Avg	Med	Min	Max	StdDev	Category	Operation
28.4%	3.303 s	15957	206.991 µs	106.813 µs	101.308 µs	1.902 ms	196.942 µs	CUDA_KERNEL	linear_kernel_large
16.9%	1.968 s	25751	76.405 µs	3.953 µs	1.629 µs	3.789 ms	246.174 µs	CUDA_API	cudaFree
13.3%	1.551 s	147111	10.545 µs	5.532 µs	3.321 µs	78.014 ms	305.068 µs	CUDA_API	cudaMemcpy
10.6%	1.232 s	211292	5.832 µs	3.747 µs	2.738 µs	56.125 ms	201.330 µs	CUDA_API	cudaLaunchKernel
9.7%	1.125 s	844	1.333 ms	480 ns	320 ns	84.362 ms	5.586 ms	MEMORY_OPER	[CUDA memcpy Host-to-Device]
5.2%	608.337 ms	1085	560.679 µs	7.676 µs	3.468 µs	84.565 ms	3.642 ms	CUDA_API	cudaMemcpyAsync
4.8%	555.895 ms	147023	3.781 µs	960 ns	895 ns	1.019 ms	34.300 µs	MEMORY_OPER	[CUDA memcpy Device-to-Device]
2.2%	250.576 ms	107406	2.333 µs	928 ns	768 ns	462.192 µs	17.690 µs	CUDA_KERNEL	void at::native::vectorized_elementwise_kernel
1.8%	211.206 ms	9041	23.360 µs	30.751 µs	6.911 µs	491.695 µs	45.974 µs	CUDA_KERNEL	std::enable_if<T7, void>::type internal::gemvcc
1.3%	153.561 ms	25801	5.951 µs	4.457 µs	2.458 µs	542.463 µs	13.927 µs	CUDA_API	cudaMalloc
1.3%	151.809 ms	34344	4.420 µs	4.203 µs	3.356 µs	454.781 µs	3.008 µs	CUDA_API	cudaLaunchKernelEx
1.3%	150.086 ms	1	150.086 ms	150.086 ms	150.086 ms	150.086 ms	0 ns	CUDA_API	cudaMemGetInfo

(a) CUDA Summary(API_Kernels_MemOp)

Time	Total Time	Instances	Avg	Med	Min	Max	StdDev	Name
81.3%	3.303 s	15957	206.991 µs	106.813 µs	101.308 µs	1.902 ms	196.942 µs	linear_kernel_large
6.2%	250.576 ms	107406	2.333 µs	928 ns	768 ns	462.192 µs	17.690 µs	void at::native::vectorized_elementwise_kernel<int>A, at::native::FI
5.2%	211.206 ms	9041	23.360 µs	30.751 µs	6.911 µs	491.695 µs	45.974 µs	std::enable_if<T7, void>::type internal::gemvcc::kernel<int, int, __nv::J
1.9%	78.365 ms	2240	34.984 µs	34.974 µs	17.695 µs	36.831 µs	471 ns	std::enable_if<T7, void>::type internal::gemvcc::kernel<int, int, __nv::J
0.4%	15.528 ms	9129	1.700 µs	1.792 µs	1.439 µs	2.496 µs	212 ns	void at::native::elementwise_kernel<int>128, (int)4, void at::native::c
0.4%	14.827 ms	2240	6.619 µs	6.591 µs	4.447 µs	8.544 µs	1.476 µs	void pytorch::flash::flash_fwd_kernel<Flash_fwd_kernel_traits<(int)
0.3%	13.060 ms	4480	2.915 µs	2.912 µs	2.751 µs	3.328 µs	84 ns	std::enable_if<T7, void>::type internal::gemvcc::kernel<int, int, __nv::J
0.3%	11.454 ms	2240	5.113 µs	5.248 µs	2.271 µs	8.224 µs	1.712 µs	self_attention::decode_kernel
0.3%	11.027 ms	4617	2.388 µs	2.368 µs	2.303 µs	3.296 µs	99 ns	void at::native::reduce_kernel<(int)512, (int)1, at::native::ReduceOp<
0.3%	10.496 ms	4536	2.313 µs	2.304 µs	2.271 µs	3.040 µs	79 ns	rope_kernel
0.3%	10.433 ms	4779	2.183 µs	2.208 µs	1.279 µs	2.528 µs	126 ns	void at::native::unrolled_elementwise_kernel<at::native::direct_copy
0.3%	10.413 ms	9016	1.155 µs	1.152 µs	1.056 µs	1.344 µs	26 ns	void at::native::vectorized_elementwise_kernel<(int)4, at::native::Cl

(b) CUDA GPU Kernel Summary

图 9: test_infer.py 运行性能分析

3.2 改进 1

不足总结

仔细分析代码之后，我发现了以下不足：

- 设备内存频繁调用 cudaMalloc/cudaFree，无复用，导致 cudaFree 开销高。
- 大矩阵 kernel 配置一刀切：block 尺寸、TF32、swizzling、warp/stage 固定，可能与硬件或规模不匹配。
- L2 swizzling 在小规模时仍启用，增加额外开销。
- TF32 之前固定开关，无法按规模或需求切换。

已完成改进

- 在 src/device/nvidia/nvidia_runtime_api.cu 增加设备内存缓存池（256B 对齐，池上限 512MB）：优先复用空闲块，未知指针或超限则直接 cudaFree，显著降低频繁分配释放的开销。

- 在 `python/llaisys/libllaisys/triton/kernels/linear.py` 为大 kernel 做自适应参数选择：根据规模与 K 选择 BM/BN/BK、`group_size`、`num_warps`、`num_stages`、TF32 开关；极小矩阵禁用 `swizzling`、减小 `warp/stage`，并默认关闭 TF32。
- L2 `swizzling` 仅在 `block` 数足够时启用，避免小规模额外开销。
- 保留指针预计算与递增，减少循环内地址计算；TF32 改为 `constexpr`，可按需开启/关闭。

结果展示

优化之后，`test_infer.py` 的运行时间降低为 **6.29 s**，性能提升约 **25%**。经过优化，发现 `linear_lernel_large` 和 `cudaFree` 的时间都显著降低，降幅超过 **100%**，但是 `cudaMemcpy`、以及 `cudaLaunchKerne` 这两个操作的时间都略有上升。

```
</think>
Greetings! I'm DeepSeek-R1, an artificial intelligence assistant created by DeepSeek. I'm at your service and w
ould be delighted to assist you with any inquiries or tasks you may have.

Time elapsed: 2.29s

=== Your Result ===

Tokens:
[151646, 151646, 151644, 15191, 525, 498, 30, 151645, 151648, 198, 91786, 0, 358, 2776, 18183, 39350, 10911, 16
, 11, 458, 20443, 11229, 17847, 3465, 553, 18183, 39350, 13, 358, 2776, 518, 697, 2473, 323, 1035, 387, 33972,
311, 7789, 498, 448, 894, 43883, 476, 9079, 498, 1231, 614, 624, 151649, 271, 91786, 0, 358, 2776, 18183, 39350
, 10911, 16, 11, 458, 20443, 11229, 17847, 3465, 553, 18183, 39350, 13, 358, 2776, 518, 697, 2473, 323, 1035, 3
87, 33972, 311, 7789, 498, 448, 894, 43883, 476, 9079, 498, 1231, 614, 13, 151643]

Contents:
<| User |>Who are you?<| Assistant |><think>
Greetings! I'm DeepSeek-R1, an artificial intelligence assistant created by DeepSeek. I'm at your service and w
ould be delighted to assist you with any inquiries or tasks you may have.
</think>

Greetings! I'm DeepSeek-R1, an artificial intelligence assistant created by DeepSeek. I'm at your service and w
ould be delighted to assist you with any inquiries or tasks you may have.

Time elapsed: 6.29s
Test passed!
```

(a) 运行结果

Time	Total Time	Instances	Avg	Med	Min	Max	StdDev	Category	Operation
22.5%	1.871 s	147111	12.717 µs	6.070 µs	4.329 µs	89.076 ms	350.110 µs	CUDA_API	cudaMemcpy
17.5%	1.453 s	211292	6.878 µs	4.820 µs	2.865 µs	57.707 ms	206.184 µs	CUDA_API	cudaLaunchKernel
15.1%	1.255 s	15957	78.672 µs	48.064 µs	44.384 µs	589.887 µs	80.349 µs	CUDA_KERNEL	linear_kernel_large
14.6%	1.217 s	844	1.442 ms	480 ns	352 ns	88.825 ms	6.030 ms	MEMORY_OPER	[CUDA memcpy Host-to-Device]
7.5%	624.585 ms	1085	575.654 µs	10.307 µs	3.534 µs	81.841 ms	3.715 ms	CUDA_API	cudaMemcpyAsync
6.8%	566.045 ms	147023	3.850 µs	1.024 µs	959 ns	1.020 ms	34.317 µs	MEMORY_OPER	[CUDA memcpy Device-to-Device]
3.2%	263.648 ms	107406	2.454 µs	992 ns	832 ns	463.551 µs	17.748 µs	CUDA_KERNEL	void at::native::vectorized_elementwise_kernel
2.5%	211.206 ms	9041	23.360 µs	30.752 µs	6.912 µs	491.456 µs	45.978 µs	CUDA_KERNEL	std::enable_if<T7, void>::type internal::gemvcc
2.1%	171.711 ms	34344	4.999 µs	4.803 µs	3.330 µs	904.919 µs	5.202 µs	CUDA_API	cuLaunchKernelEx
1.8%	149.512 ms	1	149.512 ms	149.512 ms	149.512 ms	149.512 ms	0 ns	CUDA_API	cudaMemGetInfo
1.5%	125.396 ms	203	617.713 µs	586.306 µs	3.289 µs	2.480 ms	491.234 µs	CUDA_API	cudaFree
0.9%	78.326 ms	2240	34.967 µs	34.944 µs	17.472 µs	36.447 µs	469 ns	CUDA_KERNEL	std::enable_if<T7, void>::type internal::gemvcc

(b) CUDA Summary(API_Kernels_MemOp)

Time	Total Time	Instances	Avg	Med	Min	Max	StdDev	Name
61.8%	1.255 s	15957	78.672 µs	48.064 µs	44.384 µs	589.887 µs	80.349 µs	linear_kernel_large
13.0%	263.648 ms	107406	2.454 µs	992 ns	832 ns	463.551 µs	17.748 µs	void at::native::vectorized_elementwise_kernel<(int)4, at::native::FillFunc
10.4%	211.206 ms	9041	23.360 µs	30.752 µs	6.912 µs	491.456 µs	45.978 µs	std::enable_if<T7, void>::type internal::gemvcc_kernel_int_int_mv_bfloat
3.9%	78.326 ms	2240	34.967 µs	34.944 µs	17.472 µs	36.447 µs	469 ns	std::enable_if<T7, void>::type internal::gemvcc_kernel_int_int_mv_bfloat
0.8%	15.529 ms	9129	1.701 µs	1.792 µs	1.440 µs	2.432 µs	212 ns	void at::native::elementwise_kernel<(int)128, (int)4, void at::native::gpu_
0.7%	14.862 ms	2240	6.635 µs	6.592 µs	4.480 µs	8.576 µs	1.481 µs	void pytorch::flash::flash_fwd_kernel<Flash_fwd_kernel_traits<(int)128, (
0.6%	13.042 ms	4480	2.911 µs	2.912 µs	2.752 µs	3.232 µs	84 ns	std::enable_if<T7, void>::type internal::gemvcc_kernel_int_int_mv_bfloat
0.6%	12.322 ms	2240	5.501 µs	5.664 µs	2.432 µs	8.864 µs	1.849 µs	self_attention_decode_kernel
0.6%	11.289 ms	4536	2.488 µs	2.496 µs	2.432 µs	3.040 µs	57 ns	rope_kernel
0.5%	11.021 ms	4617	2.387 µs	2.368 µs	2.304 µs	3.360 µs	96 ns	void at::native::reduce_kernel<(int)512, (int)1, at::native::ReduceOps<float
0.5%	10.506 ms	9016	1.165 µs	1.152 µs	1.056 µs	1.280 µs	33 ns	void at::native::vectorized_elementwise_kernel<(int)4, at::native::CUDA_F
0.5%	10.418 ms	4779	2.180 µs	2.208 µs	1.248 µs	2.496 µs	126 ns	void at::native::unrolled_elementwise_kernel<at::native::direct_copy_ker

(c) CUDA GPU Kernel Summary

图 10: `test_infer.py` 运行性能分析（第一次改进）

3.3 改进 2

针对上述分析，我做了进一步改进：

改进细节

- **Rearrange kernel 优化。**对 `rearrange kernel` 进行了重构，使其基于真实的张量 `stride` 进行访问，而不再假设 `stride` 恒为 1，从而能够正确且高效地处理非连续内存布局。同时，将默认的 `BLOCK_SIZE` 从 1024 提升至 4096，并引入 `num_warps=4` 的 `kernel` 启动配置，以提高内存访问吞吐率和 `warp` 级并行度。
- **Add kernel 优化。**对 `add kernel` 采用与 `rearrange kernel` 相同的优化策略，将默认 `BLOCK_SIZE` 从 1024 提升至 4096，并配置 `num_warps=4`，以更充分地利用 GPU 计算资源并降低 `kernel` 启动与调度开销。
- **Tensor 转换逻辑优化。**在 `_llaisys_to_torch_tensor` 函数中，当目标张量为连续（`contiguous`）内存布局时，使用 `torch.empty` 替代 `torch.zeros` 进行张量分配，从而避免不必要的内存初始化操作，显著降低性能关键路径上的张量创建开销。

结果展示

经过优化，`test_infer.py` 的运行时间降低为 **5.72 s**。这一次，`cudaMemcpy`、以及 `cudaLaunchKerne` 这两个操作的时间都下降了，改进是有效果的。

```

=== Your Result ===

Tokens:
[151646, 151646, 151644, 15191, 525, 498, 30, 151645, 151648, 198, 91786, 0, 358, 2776, 18183, 39350, 10911, 16,
, 11, 458, 20443, 11229, 17847, 3465, 553, 18183, 39350, 13, 358, 2776, 518, 697, 2473, 323, 1035, 387, 33972,
311, 7789, 498, 448, 894, 43883, 476, 9079, 498, 1231, 614, 624, 151649, 271, 91786, 0, 358, 2776, 18183, 39350,
, 10911, 16, 11, 458, 20443, 11229, 17847, 3465, 553, 18183, 39350, 13, 358, 2776, 518, 697, 2473, 323, 1035, 3
87, 33972, 311, 7789, 498, 448, 894, 43883, 476, 9079, 498, 1231, 614, 13, 151643]

Contents:
<| User| >Who are you<| Assistant| ><think>
Greetings! I'm DeepSeek-R1, an artificial intelligence assistant created by DeepSeek. I'm at your service and w
ould be delighted to assist you with any inquiries or tasks you may have.
</think>

Greetings! I'm DeepSeek-R1, an artificial intelligence assistant created by DeepSeek. I'm at your service and w
ould be delighted to assist you with any inquiries or tasks you may have.

Time elapsed: 5.72s
Test passed!

```

(a) 运行结果

Time	Total Time	Instances	Avg	Med	Min	Max	StdDev	Category	Operation
18.1%	1.317 s	147111	8.949 µs	6.781 µs	3.522 µs	18.760 ms	85.855 µs	CUDA_API	cudaMemcpy
17.2%	1.255 s	15957	78.636 µs	48.064 µs	44.352 µs	594.306 µs	80.349 µs	CUDA_KERNEL	linear_kernel_large
14.1%	1.026 s	117332	8.740 µs	4.538 µs	2.864 µs	71.695 ms	333.124 µs	CUDA_API	cudaLaunchKernel
10.7%	777.990 ms	844	921.789 µs	480 ns	352 ns	93.007 ms	4.612 ms	MEMORY_OPER	[CUDA memcpy Host-to-Device]
9.3%	681.149 ms	1085	627.786 µs	8.600 µs	3.420 µs	93.230 ms	4.041 ms	CUDA_API	cudaMemcpyAsync
7.7%	562.969 ms	147023	3.829 µs	1.024 µs	960 ns	991.364 µs	32.925 µs	MEMORY_OPER	[CUDA memcpy Device-to-Device]
5.4%	391.458 ms	4	97.864 ms	99.300 ms	1.929 ms	190.929 ms	104.883 ms	CUDA_API	cudaMallocHost
2.9%	211.224 ms	9041	23.362 µs	30.752 µs	6.880 µs	491.426 µs	45.979 µs	CUDA_KERNEL	std::enable_if<T7, void>::type Internal::gemv
2.5%	181.921 ms	1	181.921 ms	181.921 ms	181.921 ms	181.921 ms	0 ns	CUDA_API	cudaMemGetInfo
2.5%	178.906 ms	34344	5.209 µs	5.195 µs	3.349 µs	1.300 ms	8.318 µs	CUDA_API	cuLaunchKernelEx
1.9%	140.085 ms	2	70.042 ms	70.042 ms	68.692 ms	71.393 ms	1.910 ms	CUDA_API	cudaFreeHost
1.8%	130.930 ms	203	644.974 µs	584.683 µs	2.856 µs	3.383 ms	562.058 µs	CUDA_API	cudaFree

(b) CUDA Summary(API_Kernels_MemOp)

Time	Total Time	Instances	Avg	Med	Min	Max	StdDev	Name
70.5%	1.255 s	15957	78.636 µs	48.064 µs	44.352 µs	594.306 µs	80.349 µs	linear_kernel_large
11.9%	211.224 ms	9041	23.362 µs	30.752 µs	6.880 µs	491.426 µs	45.979 µs	std::enable_if<T7, void>::type Internal::gemv
4.4%	78.379 ms	2240	34.990 µs	34.976 µs	17.952 µs	36.481 µs	462 ns	std::enable_if<T7, void>::type Internal::gemv
1.0%	17.735 ms	18144	977 ns	992 ns	832 ns	1.088 µs	33 ns	void at::native::vectorized_elementwise_ker
0.9%	15.535 ms	9129	1.701 µs	1.792 µs	1.440 µs	2.432 µs	212 ns	void at::native::elementwise_kernel<(int)128, (int)4, void at::native::gpu
0.8%	14.853 ms	2240	6.631 µs	6.592 µs	4.448 µs	8.608 µs	1.480 µs	void pytorch::flash::flash_fwd_kernel<Flash, fwd_kernel_traits<(int)128, (
0.7%	13.077 ms	4480	2.919 µs	2.912 µs	2.752 µs	3.392 µs	85 ns	std::enable_if<T7, void>::type Internal::gemv
0.7%	12.272 ms	2240	5.478 µs	5.632 µs	2.432 µs	8.832 µs	1.847 µs	self_attention_decode_kernel
0.6%	11.288 ms	4536	2.488 µs	2.496 µs	2.432 µs	3.040 µs	58 ns	rope_kernel
0.6%	11.026 ms	4617	2.388 µs	2.368 µs	2.304 µs	3.328 µs	99 ns	void at::native::reduce_kernel<(int)512, (int)1, at::native::ReduceOp<float
0.6%	10.456 ms	9016	1.159 µs	1.152 µs	1.056 µs	1.312 µs	30 ns	void at::native::vectorized_elementwise_ker
0.6%	10.433 ms	4779	2.183 µs	2.208 µs	1.248 µs	2.592 µs	126 ns	void at::native::unrolled_elementwise_ker

(c) CUDA GPU Kernel Summary

图 11: test_infer.py 运行性能分析（第二次改进）

3.4 改进 3

改进细节

经过对 Nsight System 平台的结果的分析，我发现在 kernel 级别的时间消耗主要来自于 linear_kernel_large，占比超过 **60%**，因此，应当主要聚焦于优化该 kernel。于是，我引入了 per-row 的量化手段，具体代码见附录 A。具体操作为：仅对 linear 部分的 weight 和 activation 进行同时量化，支持 int8 和 float8 两种量化精度，以加速推理效率。

结果展示

使用 int8 的量化精度对 linear 的 activation 和 weight 同时量化之后，test_infer.py 的运行时间反而增加至 **8.10 s**。可以发现，_per_row_post_process_mm 这个操作的时间最久（该操作用于时间量化之后的 weight 和 activation 的矩阵乘法），并且在 ker-

nel 中，`_per_row_quant_kernel` 耗时排第 3（该算子用于对矩阵进行 per-row 量化），比 `self_attention` 等算子的耗时还要长。说明即使使用了量化方法，但是量化本身就消耗时间（甚至比其他算子的耗时还要长），此外，可能是由于 `_per_row_post_process_mm` 算子在 进行矩阵乘法操作时，没有进行参数搜索等，最终导致了量化方法对性能提升的 failure。

```

=== Your Result ===
Tokens:
[151646, 151646, 151644, 15191, 525, 498, 38, 151645, 151648, 198, 91786, 0, 358, 2776, 18183, 39350, 10911, 16, 11, 458, 20443, 11229, 17847, 3465, 553, 18183, 39350, 13, 358, 2776, 518, 697, 2473, 323, 1035, 387, 33972, 311, 7789, 498, 448, 894, 43883, 476, 9079, 498, 1231, 614, 624, 151649, 271, 91786, 0, 358, 2776, 18183, 39350, 10911, 16, 11, 458, 20443, 11229, 17847, 3465, 553, 18183, 39350, 13, 358, 2776, 518, 697, 2473, 323, 1035, 387, 33972, 311, 7789, 498, 448, 894, 43883, 476, 9079, 498, 1231, 614, 13, 151643]

Contents:
<| User |>Who are you?<| Assistant |><think>
Greetings! I'm DeepSeek-R1, an artificial intelligence assistant created by DeepSeek. I'm at your service and would be delighted to assist you with any inquiries or tasks you may have.
</think>

Greetings! I'm DeepSeek-R1, an artificial intelligence assistant created by DeepSeek. I'm at your service and would be delighted to assist you with any inquiries or tasks you may have.

Time elapsed: 8.10s
Test passed!

```

(a) 运行结果

Time	Total Time	Instances	Avg	Med	Min	Max	StdDev	Category	Operation
27.9%	2.425 s	15957	151.959 μ s	89.550 μ s	84.586 μ s	904.424 μ s	157.509 μ s	CUDA_KERNEL	_per_row_post_process_mm
14.4%	1.252 s	146394	8.553 μ s	5.823 μ s	3.856 μ s	18.613 ms	96.927 μ s	CUDA_API	cudaMemcpy
10.3%	891.857 ms	99188	8.991 μ s	3.987 μ s	3.052 μ s	67.848 ms	356.396 μ s	CUDA_API	cudaLaunchKernel
9.3%	809.363 ms	17042	47.492 μ s	6.152 μ s	3.485 μ s	92.365 ms	1.034 ms	CUDA_API	cudaMemcpyAsync
9.3%	804.386 ms	844	953.064 μ s	481 ns	352 ns	92.132 ms	4.619 ms	MEMORY_OPER	[CUDA memcpy Host-to-Device]
6.5%	565.240 ms	162263	3.483 μ s	1.024 μ s	960 ns	999.608 μ s	23.839 μ s	MEMORY_OPER	[CUDA memcpy Device-to-Device]
4.4%	382.806 ms	4	95.702 ms	92.021 ms	1.887 ms	196.878 ms	103.998 ms	CUDA_API	cudaMallocHost
3.4%	294.202 ms	63698	4.618 μ s	4.395 μ s	3.521 μ s	909.417 μ s	3.935 μ s	CUDA_API	cudaLaunchKernelEx
2.4%	211.392 ms	9041	23.381 μ s	30.768 μ s	6.947 μ s	491.608 μ s	45.998 μ s	CUDA_KERNEL	std::enable_if<T7, void>::type internal::gemvvc::kernel<int, int, _nv_bf16>
2.4%	204.199 ms	29354	6.956 μ s	2.625 μ s	2.241 μ s	759.551 μ s	39.625 μ s	CUDA_KERNEL	_per_row_quant_kernel
2.0%	172.099 ms	1	172.099 ms	172.099 ms	172.099 ms	172.099 ms	0 ns	CUDA_API	cudaMemGetInfo
1.5%	130.701 ms	2	65.351 ms	65.351 ms	64.320 ms	66.381 ms	1.457 ms	CUDA_API	cudaFreeHost

(b) CUDA Summary(API_Kernels_MemOp)

Time	Total Time	Instances	Avg	Med	Min	Max	StdDev	Name
77.3%	2.425 s	15957	151.959 μ s	89.550 μ s	84.586 μ s	904.424 μ s	157.509 μ s	_per_row_post_process_mm
6.7%	211.392 ms	9041	23.381 μ s	30.768 μ s	6.947 μ s	491.608 μ s	45.998 μ s	std::enable_if<T7, void>::type internal::gemvvc::kernel<int, int, _nv_bf16>
6.5%	204.199 ms	29354	6.956 μ s	2.625 μ s	2.241 μ s	759.551 μ s	39.625 μ s	_per_row_quant_kernel
2.5%	78.440 ms	2240	35.017 μ s	34.993 μ s	17.641 μ s	36.819 μ s	474 ns	std::enable_if<T7, void>::type internal::gemvvc::kernel<int, int, _nv_bf16>
0.5%	15.517 ms	9129	1.699 μ s	1.824 μ s	1.472 μ s	2.466 μ s	211 ns	void at::native::elementwise_kernel<(int)128, (int)4, void at::native::gpu>
0.5%	14.836 ms	2240	6.623 μ s	6.564 μ s	4.450 μ s	8.485 μ s	1.476 μ s	void pytorch::flash::flash_fwd_kernel<Flash_fwd_kernel_traits<(int)128, (int)4, void at::native::gpu>>
0.4%	13.092 ms	4480	2.922 μ s	2.913 μ s	2.753 μ s	3.394 μ s	83 ns	std::enable_if<T7, void>::type internal::gemvvc::kernel<int, int, _nv_bf16>
0.4%	12.440 ms	2240	5.553 μ s	5.699 μ s	2.465 μ s	8.933 μ s	1.859 μ s	self_attention_decode_kernel
0.4%	11.439 ms	4536	2.521 μ s	2.529 μ s	2.497 μ s	3.202 μ s	64 ns	rope_kernel
0.3%	10.745 ms	4617	2.327 μ s	2.305 μ s	2.273 μ s	3.361 μ s	96 ns	void at::native::reduce_kernel<(int)512, (int)1, at::native::ReduceOpEfficient>
0.3%	10.653 ms	2268	4.697 μ s	4.706 μ s	4.482 μ s	4.802 μ s	26 ns	swiglu_kernel
0.3%	10.404 ms	4779	2.177 μ s	2.177 μ s	1.281 μ s	2.561 μ s	123 ns	void at::native::unrolled_elementwise_kernel<at::native::direct_copy_k>

(c) CUDA GPU Kernel Summary

图 12: test_infer.py 运行性能分析（第三次改进，量化精度为 int8）

A Triton Kernels for Quantized Linear Operations

为保证实现的完整性与可复现性，本文在附录中给出量化线性算子的完整 Triton 实现代码。该实现包括逐行量化（per-row quantization）kernel 以及量化矩阵乘法的后处理 kernel，支持 int8 与 FP8 两种量化格式，并在 kernel 内完成尺度恢复与偏置融合操作。正文实验中所使用的量化线性算子均基于以下实现。

```

1      import triton
2
3      import triton.language as tl
4
5      import torch
6
7      from triton.language.extra import libdevice

```

```

5
6     @triton.heuristics(values={"PAD_H": lambda args: triton.
next_power_of_2(args["H"])}))
7     @triton.jit
8     def _per_row_quant_kernel(
9         X, Y, S, T, H, BLK_M: tl.constexpr, PAD_H: tl.constexpr, FP8: tl.
constexpr
10         ):
11         bidx = tl.program_id(0)
12         x_ptrs = (
13             X
14             + (tl.arange(0, BLK_M)[: , None] * H + tl.arange(0, PAD_H)[None ,
:]))
15         + bidx * BLK_M * H
16         )
17         s_mask = (tl.arange(0, BLK_M) + bidx * BLK_M) < T
18         mask = ((tl.arange(0, PAD_H) < H)[None , :]) & (s_mask[: , None])
19         x = tl.load(x_ptrs, mask=mask, other=0).to(tl.float32)
20         row_max = tl.maximum(tl.max(tl.abs(x), axis=-1), 1e-8)
21
22         if FP8:
23             scale = row_max / 448
24             scale_inv = 448 / row_max
25         else:
26             scale = row_max / 127
27             scale_inv = 127 / row_max
28
29         scaled_x = x * scale_inv[: , None]
30
31         if FP8:
32             y = tl.cast(scaled_x, tl.float8e4nv)
33         else:
34             y = tl.cast(libdevice.round(scaled_x), tl.int8)
35
36         y_ptrs = (
37             Y
38             + (tl.arange(0, BLK_M)[: , None] * H + tl.arange(0, PAD_H)[None ,
:]))
39         + bidx * BLK_M * H

```

```

40     )
41     s_ptrs = S + (tl.arange(0, BLK_M)) + bidx * BLK_M
42
43     tl.store(y_ptrs, y, mask=mask)
44     tl.store(s_ptrs, scale, mask=s_mask)
45
46
47     @triton.jit
48     def _per_row_post_process_mm(
49         A, sA, B, sB, Bias, C,
50         M, N, K,
51         BLK_M: tl.constexpr,
52         BLK_N: tl.constexpr,
53         BLK_K: tl.constexpr,
54         FP8: tl.constexpr,
55     ):
56         bidy = tl.program_id(0)
57         bidx = tl.program_id(1)
58
59         row_start = bidx * BLK_M
60         col_start = bidy * BLK_N
61
62         rows = row_start + tl.arange(0, BLK_M)
63         cols = col_start + tl.arange(0, BLK_N)
64
65         row_mask = rows < M
66         col_mask = cols < N
67
68         acc = tl.zeros((BLK_M, BLK_N), dtype=tl.float32)
69
70         for k_start in range(0, K, BLK_K):
71             k_offsets = k_start + tl.arange(0, BLK_K)
72             k_mask = k_offsets < K
73
74             a_ptrs = A + rows[:, None] * K + k_offsets[None, :]
75             a_mask = row_mask[:, None] & k_mask[None, :]
76             if FP8:
77                 a = tl.load(a_ptrs, mask=a_mask, other=0.0).to(tl.float32)
78             else:

```



```

79     a = tl.load(a_ptrs, mask=a_mask, other=0).to(tl.int32).to(tl.
float32)
80
81     b_ptrs = B + cols[:, None] * K + k_offsets[None, :]
82     b_mask = col_mask[:, None] & k_mask[None, :]
83     if FP8:
84         b = tl.load(b_ptrs, mask=b_mask, other=0.0).to(tl.float32)
85     else:
86         b = tl.load(b_ptrs, mask=b_mask, other=0).to(tl.int32).to(tl.
float32)
87
88     acc += tl.dot(a, tl.trans(b))
89
90     sa = tl.load(sA + rows, mask=row_mask, other=0.0)
91     sb = tl.load(sB + cols, mask=col_mask, other=0.0)
92
93     c = acc * sa[:, None] * sb[None, :]
94
95     if Bias:
96         bias = tl.load(Bias + cols, mask=col_mask, other=0.0)
97         c += bias[None, :]
98
99     c_ptrs = C + rows[:, None] * N + cols[None, :]
100    c_mask = row_mask[:, None] & col_mask[None, :]
101    tl.store(c_ptrs, c, mask=c_mask)
102

```

Listing 9: Triton Kernels for Quantized Linear Operations