

Night Compute

triton文件生成并使用ui界面分析 .ncu-rep 流程:

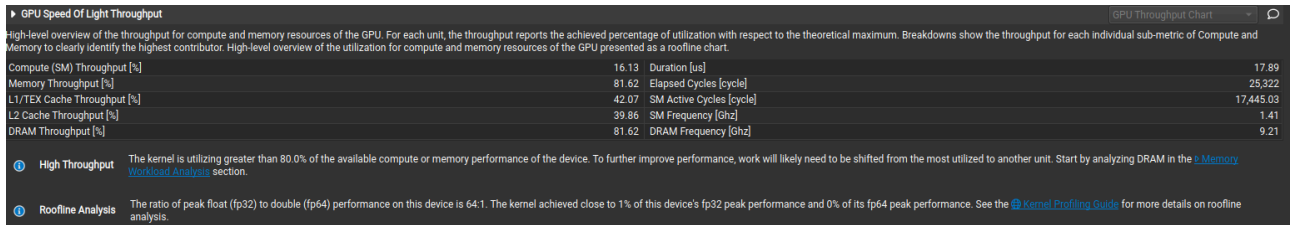
```
yst@yst-X299-WU8 ~/文/y/C/C/C/lec1_profile_CUDA_kernels (main)> /usr/local/NVIDIA-Nsight-Compute-2024.3/ncu --set full -o matrix_squre python triton_sample.py
==PROF== Connected to process 2203444 (/home/yst/miniconda3/envs/yst_pytorch/bin/python3.8)
==PROF== Profiling "distribution_elementwise_grid..." - 0: 0%...50%...100% - 37 passes
==PROF== Profiling "square_kernel_0d1d234" - 1: 0%...50%...100% - 37 passes
==PROF== Profiling "vectorized_elementwise_kernel" - 2: 0%...50%...100% - 38 passes
==PROF== Profiling "vectorized_elementwise_kernel" - 3: 0%...50%...100% - 37 passes
==PROF== Profiling "vectorized_elementwise_kernel" - 4: 0%...50%...100% - 37 passes
==PROF== Profiling "vectorized_elementwise_kernel" - 5: 0%...50%...100% - 37 passes
==PROF== Profiling "vectorized_elementwise_kernel" - 6: 0%...50%...100% - 37 passes
==PROF== Profiling "vectorized_elementwise_kernel" - 7: 0%...50%...100% - 37 passes
==PROF== Profiling "vectorized_elementwise_kernel" - 8: 0%...50%...100% - 37 passes
==PROF== Profiling "vectorized_elementwise_kernel" - 9: 0%...50%...100% - 37 passes
==PROF== Profiling "vectorized_elementwise_kernel" - 10: 0%...50%...100% - 37 passes
==PROF== Profiling "vectorized_elementwise_kernel" - 11: 0%...50%...100% - 37 passes
==PROF== Profiling "vectorized_elementwise_kernel" - 12: 0%...50%...100% - 37 passes
==PROF== Profiling "vectorized_elementwise_kernel" - 13: 0%...50%...100% - 37 passes
==PROF== Profiling "vectorized_elementwise_kernel" - 14: 0%...50%...100% - 37 passes
==PROF== Profiling "vectorized_elementwise_kernel" - 15: 0%...50%...100% - 37 passes
==PROF== Profiling "reduce_kernel" - 16: 0%...50%...100% - 37 passes
passed
==PROF== Disconnected from process 2203444
==PROF== Report: /home/yst/文档/yst/CUDA/cuda-mode/CUDA-Optimization/lec1_profile_CUDA_kernels/matrix_squre.ncu-rep
yst@yst-X299-WU8 ~/文/y/C/C/C/lec1_profile_CUDA_kernels (main)> /usr/local/NVIDIA-Nsight-Compute-2024.3/ncu-ut matrix_squre.ncu-rep
```

源码:

```
6 # 装饰器: 使用triton的jit编译器来编译这个函数
7 # BLOCK_SIZE: thread数量, 用于处理的元素数量, 使用 tl.constexpr 指定为编译时常量
8 @triton.jit
9 def square_kernel(output_ptr, input_ptr, output_row_stride, input_row_stride, n_cols, BLOCK_SIZE: tl.constexpr):
10     row_idx = tl.program_id(0) # 获取当前程序id, block_id
11     row_start_ptr = input_ptr + row_idx * input_row_stride;
12     col_offsets = tl.arange(0, BLOCK_SIZE) # thread_id
13     input_ptrs = row_start_ptr + col_offsets
14
15     # 从GPU中加载数据
16     # mask: BLOCK_SIZE = triton.next_power_of_2(n_cols) 可能大于 n_cols
17     row = tl.load(input_ptrs, mask=col_offsets < n_cols, other=-float('inf'))
18
19     # 计算加载元素的平方
20     square_output = row * row
21
22     output_row_start_ptr = output_ptr + row_idx * output_row_stride
23     output_ptrs = output_row_start_ptr + col_offsets
24     tl.store(output_ptrs, square_output, mask=col_offsets < n_cols)
25
26 def square(x):
27     n_rows, n_cols = x.shape
28     BLOCK_SIZE = triton.next_power_of_2(n_cols) # 找到最接近的大于所给参数的2^n eg: n_cols: 20 -> 返回32
29     num_warps = 4
30     if BLOCK_SIZE >= 2048:
31         num_warps = 8
32     if BLOCK_SIZE >= 4096:
33         num_warps = 16
34
35     y = torch.empty_like(x)
36
37     # [:]: grid_size
38     # (n_rows, ): 一维grid
39     square_kernel[(n_rows, )](y, x, y.stride(0), x.stride(0), n_cols, BLOCK_SIZE)
40
41     return y
```

Night Compute分析:

1. GPU Speed of Light Throughput

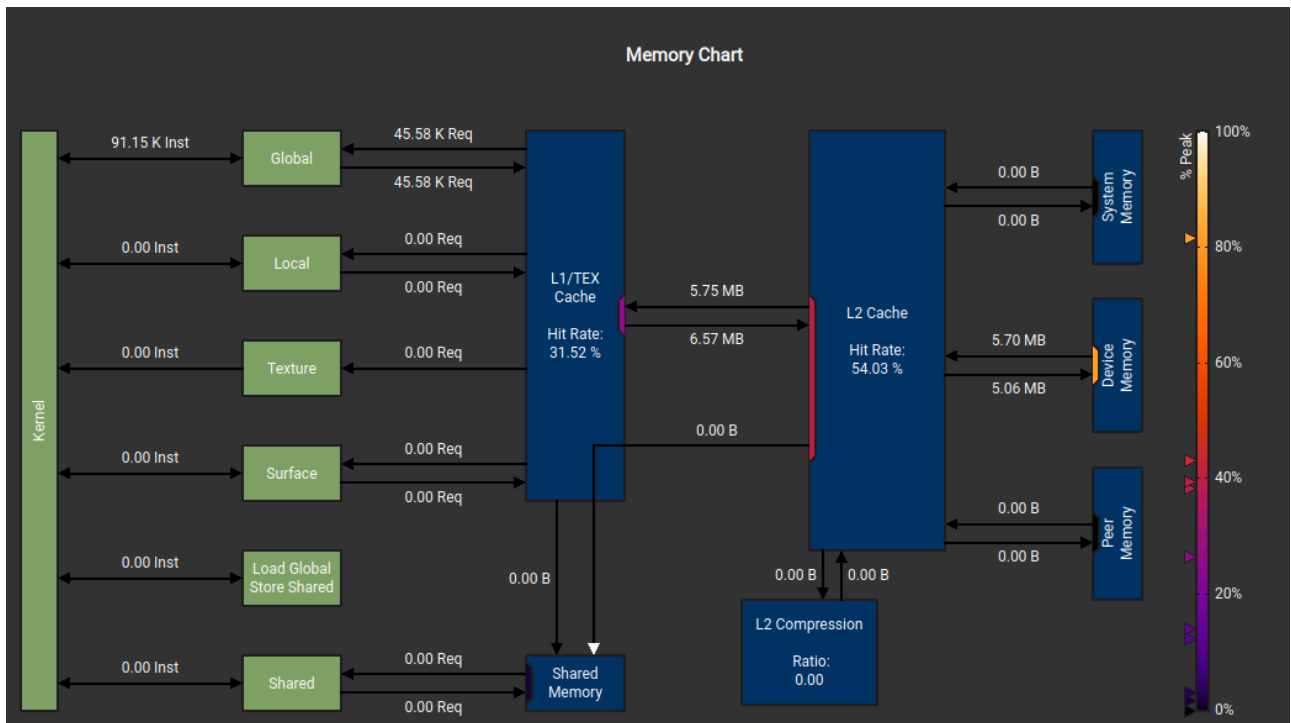
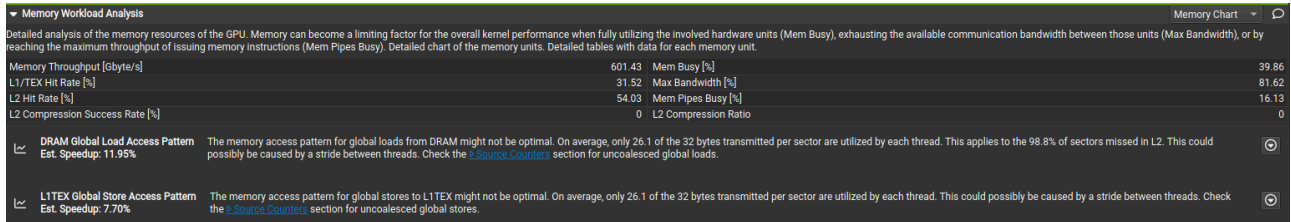


可以看出, L1 和 L2 利用率不高, 并且Memory > Compute, 说明是访存密集型算子

eg: x: [H, W]

计算次数: $H \times W$ 访存次数: $H \times W \times 2$

2. Memory Workload Analysis

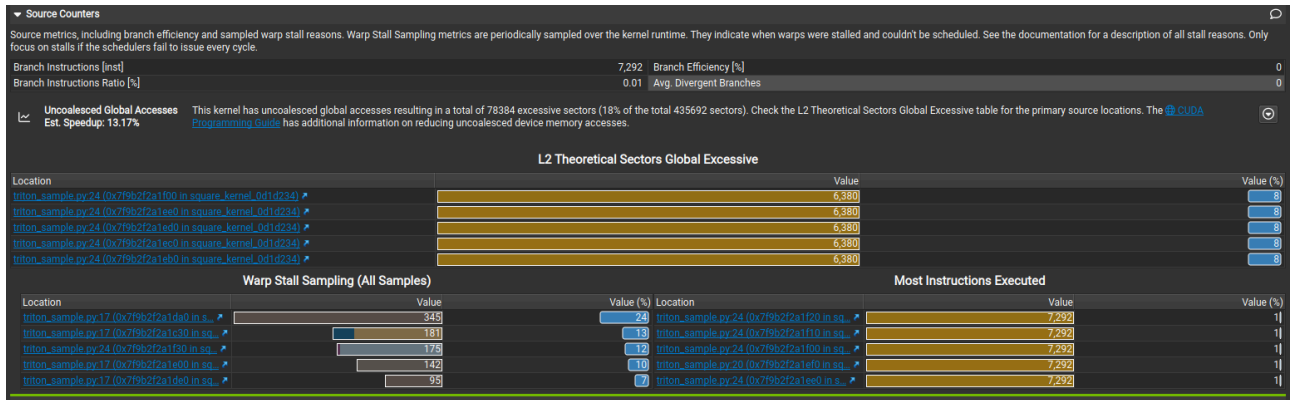


L1 和 L2 Hit Rate 不高, Max BandWidth很低(3080: 790GB/s)

eg: x: [1823, 781] FP32

数据量: $1823 \times 781 \times 32 / 8 / 1024 / 1024 = 5.43\text{MB}$

3. Source Counters



分支指令:

一共7292个分支, 分支效率和平均发散分支都是0 -> 分支预测做得很好(因为确实没有分支)

未合并的全局访问:

有未合并的全局访问, 导致78384个多余的sector

L2 未有效利用:

问题出在24行

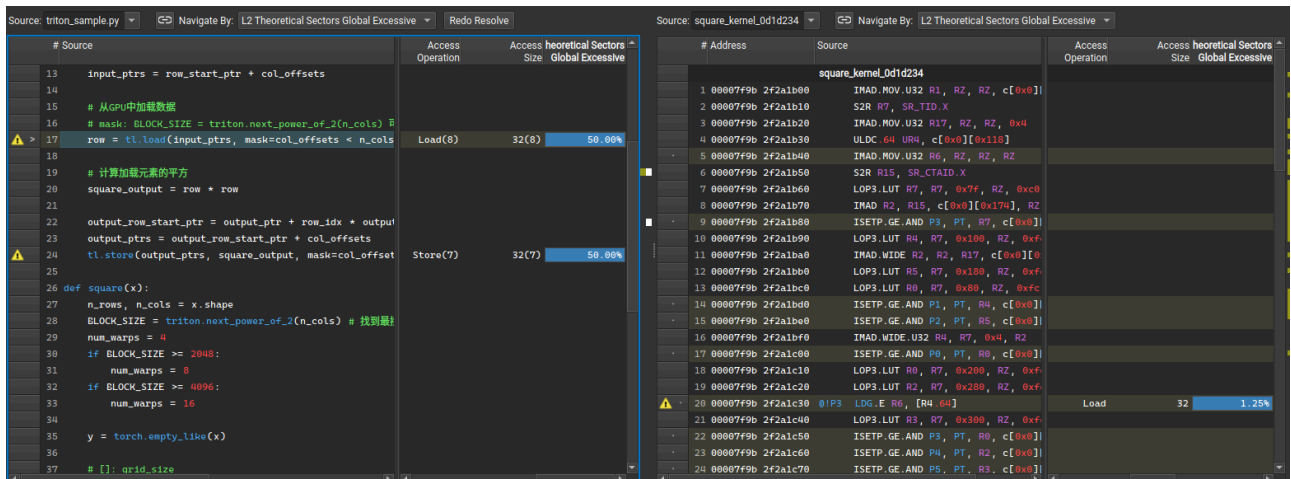
`tl.store(output_ptrs, square_output, mask=col_offsets < n_cols)`

Warp停滞

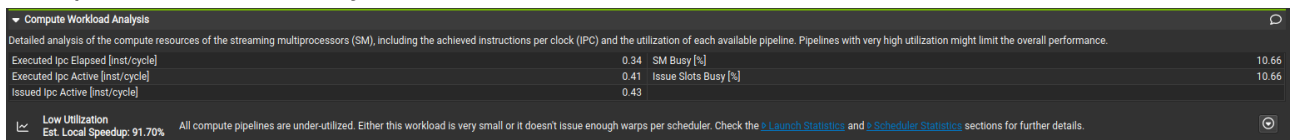
问题处在17和24行(一个读一个存)

`row = tl.load(input_ptrs, mask=col_offsets < n_cols, other=float('inf'))`

`tl.store(output_ptrs, square_output, mask=col_offsets < n_cols)`



4. Compute Workland Analysis



计算pipeline没有被充分利用(因为是访存密集)

5. Launch Statistics

Launch Statistics

Summary of the configuration used to launch the kernel. The launch configuration defines the size of the kernel grid, the division of the grid into blocks, and the GPU resources needed to execute the kernel. Choosing an efficient launch configuration maximizes device utilization.

Grid Size	1,823	Function Cache Configuration	CachePreferNone
Registers Per Thread [register/thread]	20	Static Shared Memory Per Block [byte/block]	0
Block Size	128	Dynamic Shared Memory Per Block [byte/block]	0
Threads [thread]	233,344	Driver Shared Memory Per Block [kbyte/block]	1.02
Waves Per SM	2.23	Shared Memory Configuration Size [kbyte]	16.38
Uses Green Context	0	# SMs [SM]	68

Tail Effect

Est. Speedup: 33.33%

A wave of thread blocks is defined as the maximum number of blocks that can be executed in parallel on the target GPU. The number of blocks in a wave depends on the number of multiprocessors and the theoretical occupancy of the kernel. This kernel launch results in 2 full waves and a partial wave of 190 thread blocks. Under the assumption of a uniform execution duration of all thread blocks, the partial wave may account for up to 33.3% of the total kernel runtime with a lower occupancy of 25.3%. Try launching a grid with no partial wave. The overall impact of this tail effect also lessens with the number of full waves executed for a grid. See the [Hardware Mode](#) description for more details on launch configurations.