

MSRA SH Triton Study Group

4. Nsight Compute

2025/08/08

What is NVIDIA Nsight Compute

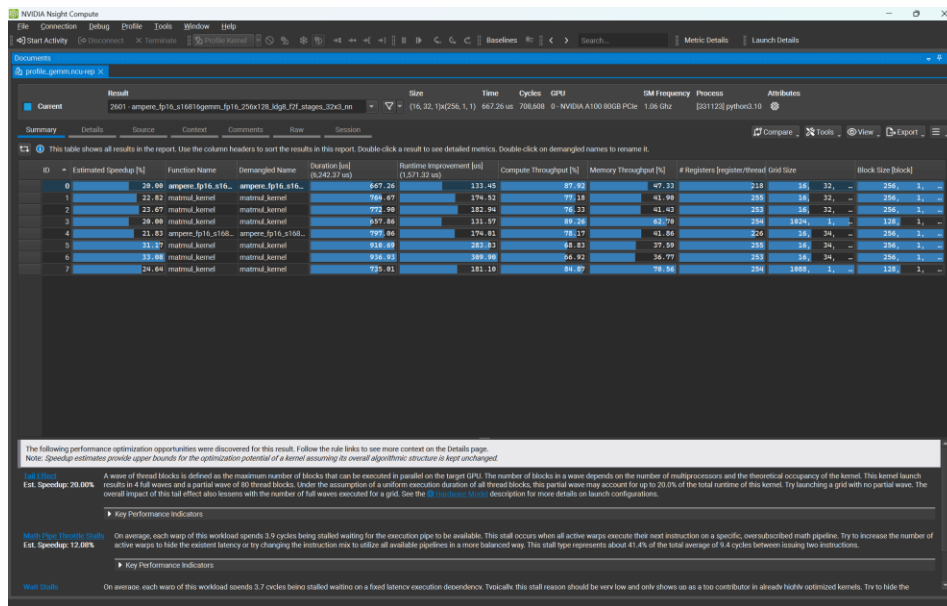
- NVIDIA Nsight Compute is an interactive profiler for CUDA and NVIDIA OptiX that provides detailed performance metrics and API debugging.
- Download: [Getting Started with Nsight Compute | NVIDIA Developer](#)
 - Host: on (Windows) Laptop
 - Target: on Linux (already installed on GCR A100 Sandbox: /usr/local/cuda/bin/ncu)

How to use NVIDIA Nsight Compute

- On target device: ncu + any program that uses GPUs

```
ncu="/usr/local/cuda/bin/ncu"
python="$ (which python)"
sudo $ncu --set full -o profile_gemm -f $python ./call.py
```

- Download the log file (xxx.ncu-rep) and open with GUI



Example: Triton GeMM Kernel

- The triton GeMM kernel in the last class is not optimized

```
Problem size: (4096, 4096, 4096)  
Torch GeMM Latency: 0.632 ms  
Triton (Naive) GeMM Latency: 0.679 ms
```

Prepare for NCU profiling

- Get Autotune result

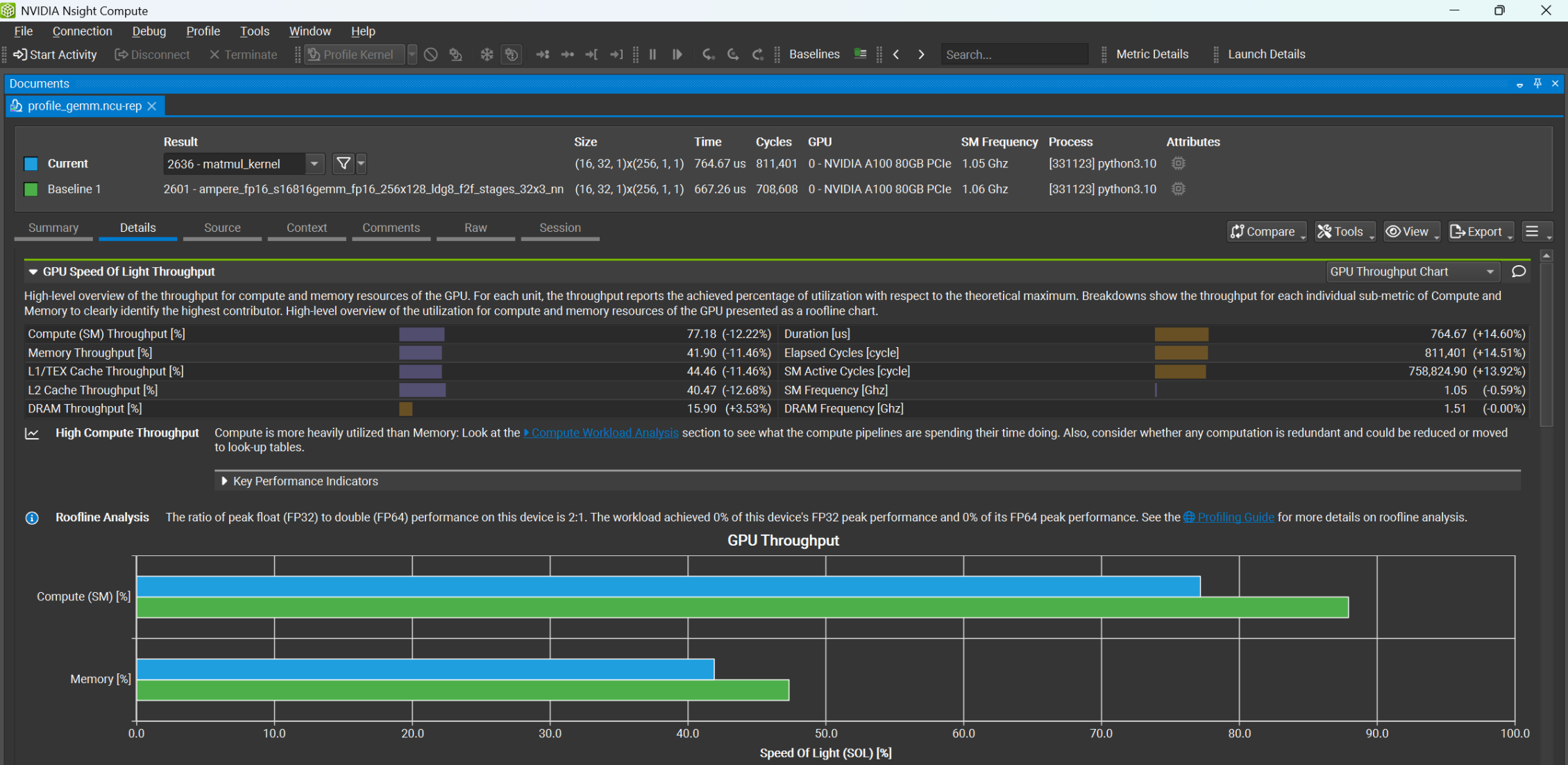
```
(moe) chengzhang@microsoft.com@GCRAZGDL1536:~/TritonStudyGroup/3_Triton_GeMM$  
export TRITON_PRINT_AUTOTUNING=1  
(moe) chengzhang@microsoft.com@GCRAZGDL1536:~/TritonStudyGroup/3_Triton_GeMM$  
python test.py  
Triton autotuning for function matmul_kernel finished after 1.67s; best config  
selected: BLOCK_SIZE_M: 128, BLOCK_SIZE_N: 256, BLOCK_SIZE_K: 64, num_warps:  
8, num_ctas: 1, num_stages: 3, maxnreg: None;
```

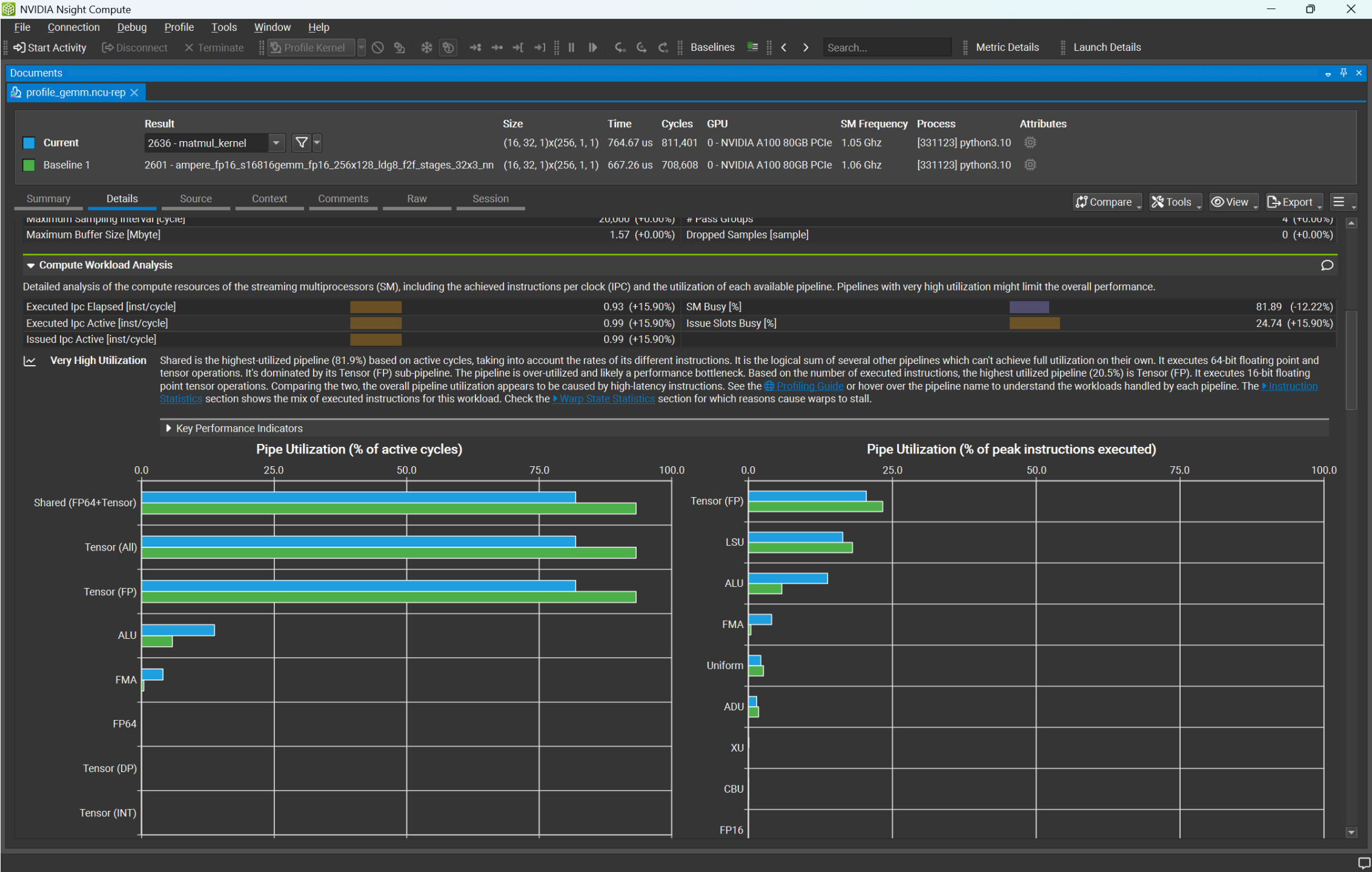
- Rewrite triton kernel launch

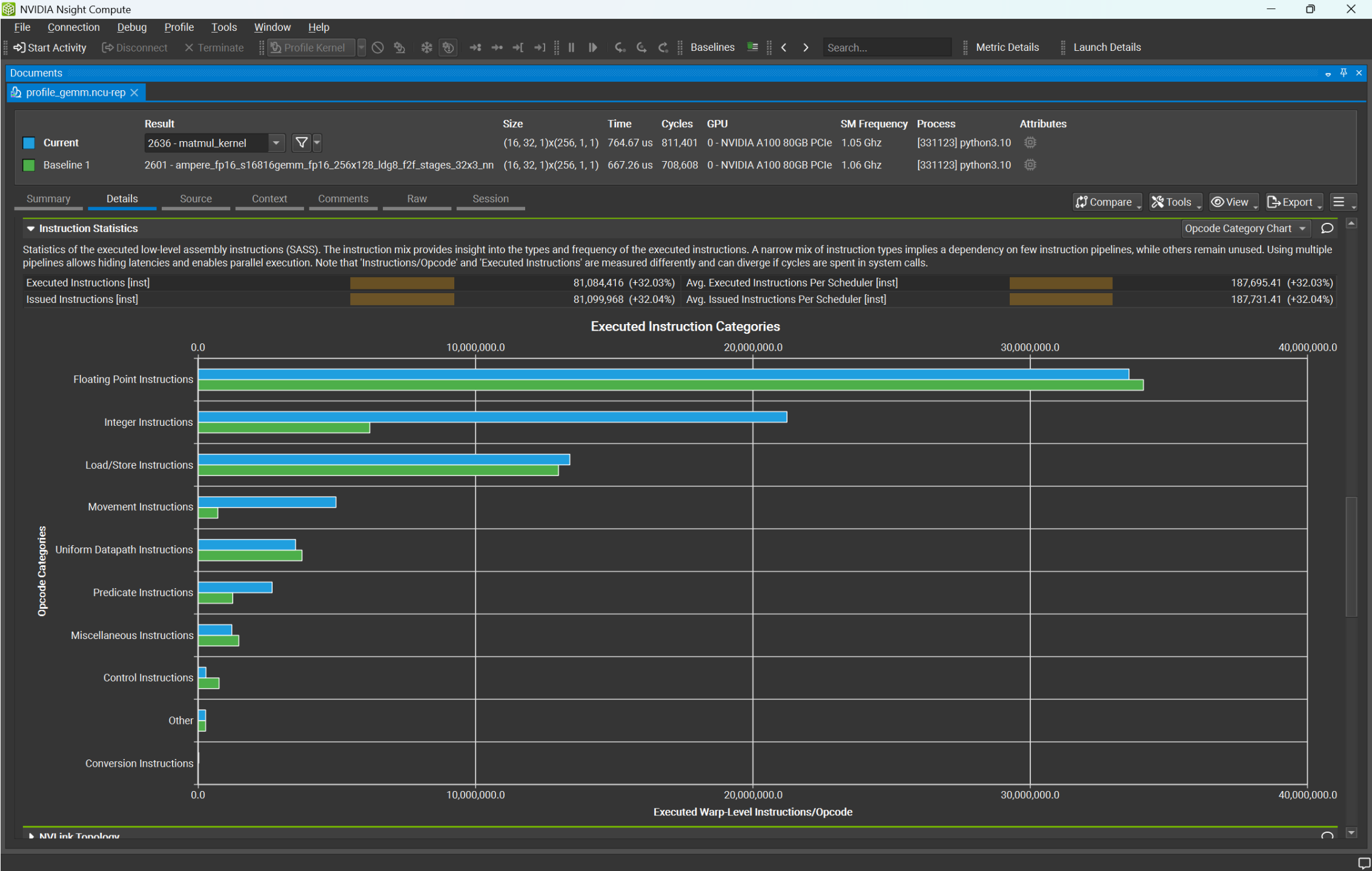
```
def triton_matmul(a: torch.Tensor, b: torch.Tensor) -> torch.Tensor:  
    M, K = a.shape  
    K2, N = b.shape  
    assert K == K2, "Inner dimensions must match for matrix multiplication"  
    c = torch.empty((M, N), dtype=a.dtype, device=a.device)  
    BLOCK_SIZE_M = 128  
    BLOCK_SIZE_N = 256  
    BLOCK_SIZE_K = 64  
    grid = (triton.cdiv(N, BLOCK_SIZE_N), triton.cdiv(M, BLOCK_SIZE_M), )  
    matmul_kernel[grid](  
        a, b, c, M, N, K,  
        a.stride(0), a.stride(1),  
        b.stride(0), b.stride(1),  
        c.stride(0), c.stride(1),  
        num_stages=3, num_warps=8,  
        BLOCK_SIZE_M=BLOCK_SIZE_M, BLOCK_SIZE_N=BLOCK_SIZE_N, BLOCK_SIZE_K=BLOCK_SIZE_K,  
    )  
    return c
```

Run NCU profile and open replay

```
==PROF== Connected to process 21919 (/home/chengzhang/anaconda3/envs/moe/bin/python3.10)
==PROF== Profiling "ampere_fp16_s16816gemm_fp16_2..." - 0: 0%....50%....100% - 49 passes
==PROF== Profiling "matmul_kernel" - 1: 0%....50%....100% - 50 passes
==PROF== Profiling "matmul_kernel" - 2: 0%....50%....100% - 49 passes
==PROF== Profiling "matmul_kernel" - 3: 0%....50%....100% - 50 passes
==PROF== Profiling "ampere_fp16_s16816gemm_fp16_2..." - 4: 0%....50%....100% - 50 passes
==PROF== Profiling "matmul_kernel" - 5: 0%....50%....100% - 49 passes
==PROF== Profiling "matmul_kernel" - 6: 0%....50%....100% - 49 passes
==PROF== Profiling "matmul_kernel" - 7: 0%....50%....100% - 49 passes
==PROF== Disconnected from process 21919
```







Reduce ALU & Integer Instructions

```
# Create pointers for the first blocks of A and B
offs_m = (start_m + t1.arange(0, BLOCK_SIZE_M)) % M
offs_n = (start_n + t1.arange(0, BLOCK_SIZE_N)) % N
offs_k = t1.arange(0, BLOCK_SIZE_K)
a_ptrs = a_ptr + offs_m[:, None] * stride_am + offs_k[None, :] * stride_ak
b_ptrs = b_ptr + offs_k[:, None] * stride_bk + offs_n[None, :] * stride_bn
```



```
# Create pointers for the first blocks of A and B
offs_m = start_m + t1.arange(0, BLOCK_SIZE_M)
offs_n = start_n + t1.arange(0, BLOCK_SIZE_N)
offs_k = t1.arange(0, BLOCK_SIZE_K)
offs_m = t1.max_contiguous(t1.multiple_of(offs_m % M, BLOCK_SIZE_M), BLOCK_SIZE_M)
offs_n = t1.max_contiguous(t1.multiple_of(offs_n % N, BLOCK_SIZE_N), BLOCK_SIZE_N)
a_ptrs = a_ptr + offs_m[:, None] * stride_am + offs_k[None, :] * stride_ak
b_ptrs = b_ptr + offs_k[:, None] * stride_bk + offs_n[None, :] * stride_bn
```

```

for k in range(0, K, BLOCK_SIZE_K):
    mask_k = k + offs_k < K # Mask for valid K indices
    a = tl.load(a_ptrs, mask=mask_k[None, :], other=0.) # Load a block of A to shared memory
    b = tl.load(b_ptrs, mask=mask_k[:, None], other=0.) # Load a block of B to shared memory
    accumulator = tl.dot(a, b, accumulator) # Call tensor core
    a_ptrs += BLOCK_SIZE_K * stride_ak # Move to the next block of A
    b_ptrs += BLOCK_SIZE_K * stride_bk # Move to the next block of B

```



```

mid = (K // BLOCK_SIZE_K) * BLOCK_SIZE_K
for k in range(0, mid, BLOCK_SIZE_K):
    a = tl.load(a_ptrs) # Load a block of A to shared memory
    b = tl.load(b_ptrs) # Load a block of B to shared memory
    accumulator = tl.dot(a, b, accumulator) # Call tensor core
    a_ptrs += BLOCK_SIZE_K * stride_ak # Move to the next block of A
    b_ptrs += BLOCK_SIZE_K * stride_bk # Move to the next block of B
if mid < K:
    mask_k = mid + offs_k < K # Mask for valid K indices
    a = tl.load(a_ptrs, mask=mask_k[None, :], other=0.) # Load a block of A to shared memory
    b = tl.load(b_ptrs, mask=mask_k[:, None], other=0.) # Load a block of B to shared memory
    accumulator = tl.dot(a, b, accumulator) # Call tensor core

```

Reduce ALU & Integer Instructions

- Not ideal

```
Problem size: (4096, 4096, 4096)
Torch GeMM Latency: 0.632 ms
Triton (Naive) GeMM Latency: 0.679 ms
Triton (Good) GeMM Latency: 0.684 ms
```

- Reason: Limited Registers

Function Name	Demangled Name	Duration [us] (6,242.37 us)	Runtime Improvement [us] (1,571.32 us)	Compute Throughput [%]	Memory Throughput [%]	# Registers [register/thread]
ampere_fp16_s168...	ampere_fp16_s168...	667.26 (-12.74%)	133.45	87.92 (+13.92%)	47.33 (+12.95%)	218 (-14.51%)
matmul_kernel	matmul_kernel	764.67	174.52	77.18	41.90	255
matmul_kernel	matmul_kernel	772.90 (+1.08%)	182.94	76.33 (-1.10%)	41.43 (-1.13%)	253 (-0.78%)

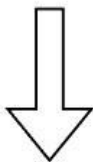
Threadblock Swizzle

CTA-0	CTA-1	CTA-2	CTA-3	CTA-4	CTA-5	CTA-6	CTA-7	CTA-8	CTA-9	CTA-10	CTA-11
CTA-12	CTA-13	CTA-14	CTA-15	CTA-16	CTA-17	CTA-18	CTA-19				
							CTA-31	CTA-32	CTA-33	CTA-34	CTA-35
											CTA-47

C Matrix

CTA-0	CTA-1	CTA-2	CTA-3	CTA-4	CTA-5	CTA-6	CTA-7	CTA-8	CTA-9	CTA-10	CTA-11
CTA-12	CTA-13	CTA-14	CTA-15	CTA-16	CTA-17	CTA-18	CTA-19				
							CTA-31	CTA-32	CTA-33	CTA-34	CTA-35
											CTA-47

C Matrix



TILE-0				TILE-1				TILE-2			
CTA-0	CTA-1	CTA-2	CTA-3	CTA-16	CTA-17	CTA-18	CTA-19	CTA-32	CTA-33	CTA-34	CTA-35
CTA-4	CTA-5	CTA-6	CTA-7								
CTA-8	CTA-9	CTA-10	CTA-11								
CTA-12	CTA-13	CTA-14	CTA-15				CTA-31				CTA-47

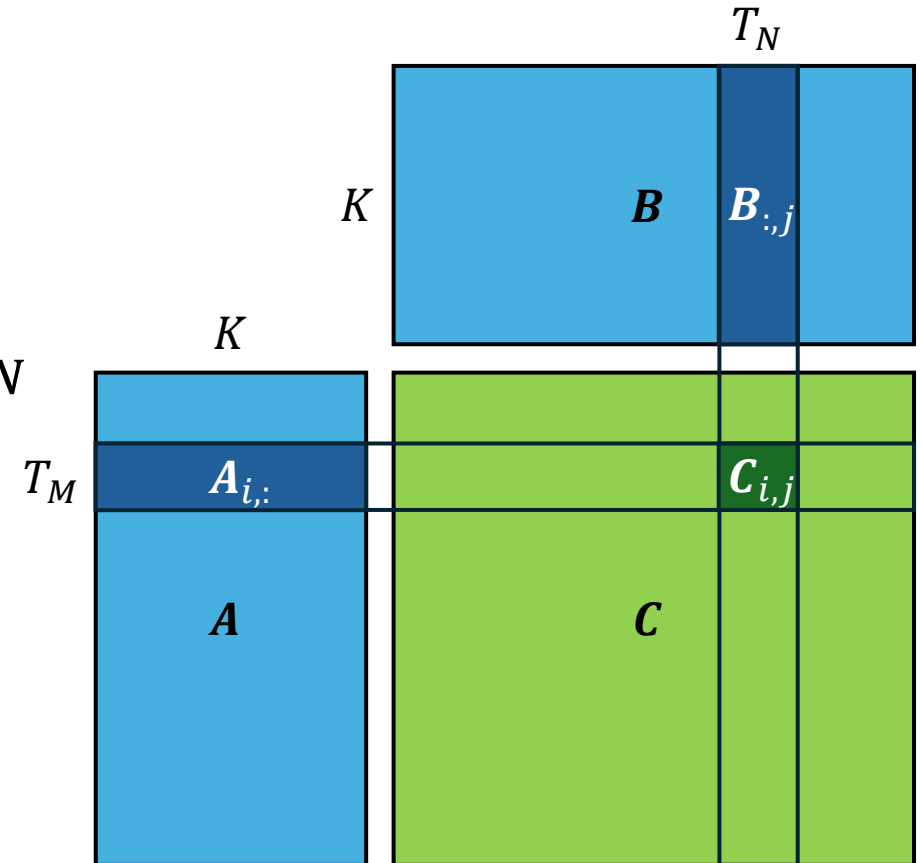
C Matrix

Arithmetic Intensity of GeMM

- Block size $[T_M, T_N] \ll K$
- $\mathbf{C}_{ij} = \sum_{k=0}^K \mathbf{A}_{ik} \mathbf{B}_{kj}$
- Block HBM Access: $T_M K + T_N K + T_M T_N$
- Block Calculation: $T_M T_N K$
- Arithmetic Intensity :

$$\frac{T_M T_N K}{T_M K + T_N K + T_M T_N} \approx \frac{T_M T_N}{T_M + T_N}$$

- $\frac{128 \times 256}{128 + 256} > \frac{128 \times 128}{128 + 128}$



Change Block Size & Threadblock Swizzle

```
BLOCK_SIZE_M = 128
BLOCK_SIZE_N = 256
BLOCK_SIZE_K = 64
grid = (triton.cdiv(N, BLOCK_SIZE_N), triton.cdiv(M, BLOCK_SIZE_M), )
matmul_kernel[grid](
    a, b, c, M, N, K,
    a.stride(0), a.stride(1),
    b.stride(0), b.stride(1),
    c.stride(0), c.stride(1),
    num_stages=3, num_warps=8,
```

(Compute bound)



(Balanced)

```
BLOCK_SIZE_M = 128
BLOCK_SIZE_N = 128
BLOCK_SIZE_K = 32
grid = (triton.cdiv(N, BLOCK_SIZE_N) * triton.cdiv(M, BLOCK_SIZE_M), )
matmul_kernel[grid](
    a, b, c, M, N, K,
    a.stride(0), a.stride(1),
    b.stride(0), b.stride(1),
    c.stride(0), c.stride(1),
    num_stages=4, num_warps=4,
```



```
# Thread block index
pid_n = tl.program_id(axis=0)
pid_m = tl.program_id(axis=1)
start_m = pid_m * BLOCK_SIZE_M
start_n = pid_n * BLOCK_SIZE_N
```



```
# Thread block index
GROUP_SIZE_M = 8
pid = tl.program_id(axis=0)
num_pid_m = tl.cdiv(M, BLOCK_SIZE_M)
num_pid_n = tl.cdiv(N, BLOCK_SIZE_N)
num_pid_in_group = GROUP_SIZE_M * num_pid_n
group_id = pid // num_pid_in_group
first_pid_m = group_id * GROUP_SIZE_M
group_size_m = min(num_pid_m - first_pid_m, GROUP_SIZE_M)
pid_m = first_pid_m + ((pid % num_pid_in_group) % group_size_m)
pid_n = (pid % num_pid_in_group) // group_size_m
start_m = pid_m * BLOCK_SIZE_M
start_n = pid_n * BLOCK_SIZE_N
```

Change Block Size & Threadblock Swizzle

- Good

```
Problem size: (4096, 4096, 4096)
Torch GeMM Latency: 0.632 ms
Triton (Naive) GeMM Latency: 0.679 ms
Triton (Good) GeMM Latency: 0.684 ms
Triton (Better) GeMM Latency: 0.665 ms
Problem size: (4321, 4080, 4080)
Torch GeMM Latency: 0.729 ms
Triton (Naive) GeMM Latency: 0.792 ms
Triton (Good) GeMM Latency: 0.819 ms
Triton (Better) GeMM Latency: 0.733 ms
```

Reading Materials: Nsight Compute

- [4. Nsight Compute CLI — NsightCompute 13.0 documentation](#)
- [NVIDIA Development Tools Solutions - ERR_NVGPUCTRPERM: Permission issue with Performance Counters | NVIDIA Developer](#)

Homework

- Play with [TritonStudyGroup/4_Nsight_Compute at main · Starmys/TritonStudyGroup](#)

Reading Materials: Triton Static Hints

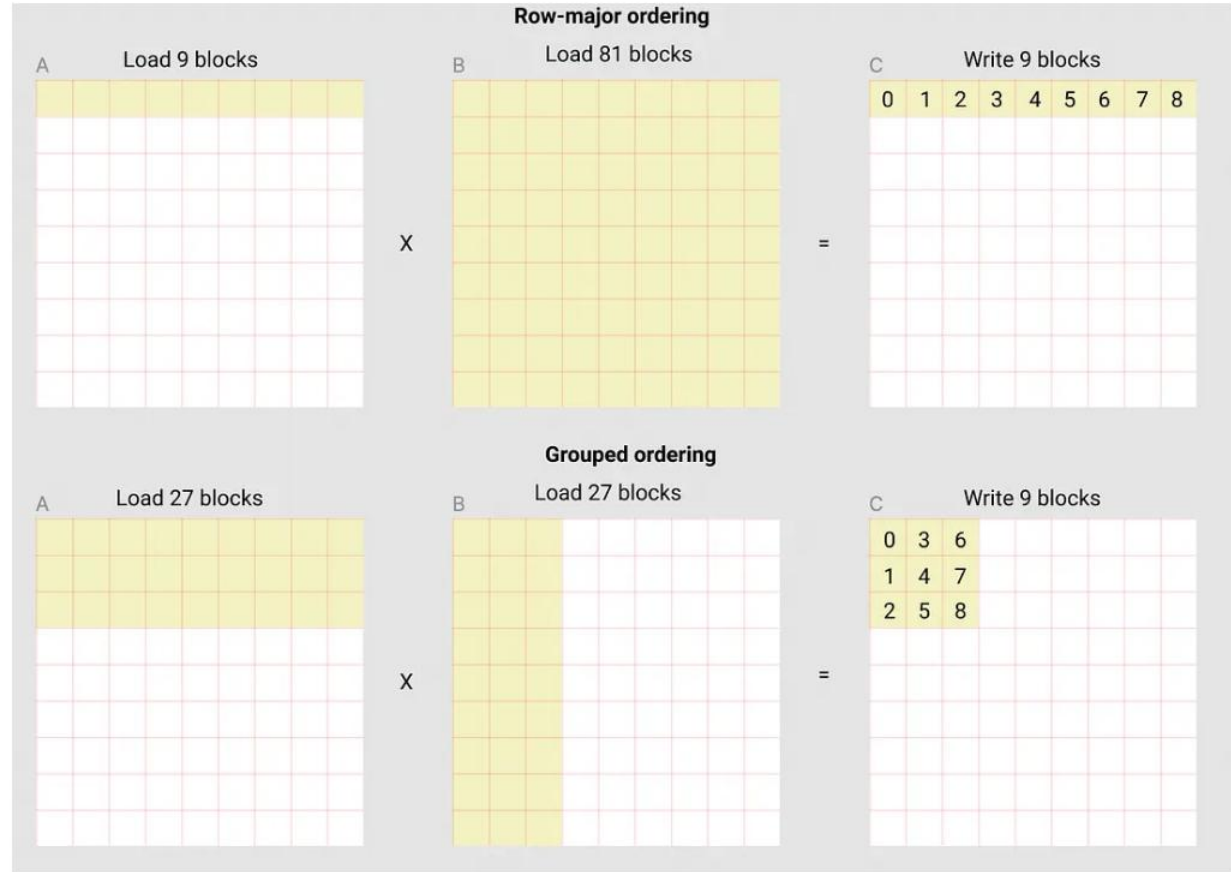
- [Confusion of arange, multiple_of, max_contiguous · triton-lang/triton · Discussion #5939](#)
- [triton/include/triton/Analysis/AxisInfo.h at 65d9862f5a9029827a7ae04439f07d7e7eadf859 · triton-lang/triton](#)

Reading Materials

- [Understanding the Triton Tutorials Part 1 | by Isamu Isozaki | Medium](#)
- [Deep Dive on CUTLASS Ping-Pong GEMM Kernel – PyTorch](#)

L2 Cache Optimizations

One way the tutorial suggests we do this is very simple, just increase the number of rows loaded at once!



As you see, if we load one row and load all the columns, we can write 9 blocks in the output. But if we load 3 rows and 3 columns, we can write 9 blocks but without loading $81 - 54 = 27$ rows! In practice, this can save around 10% of computing.