

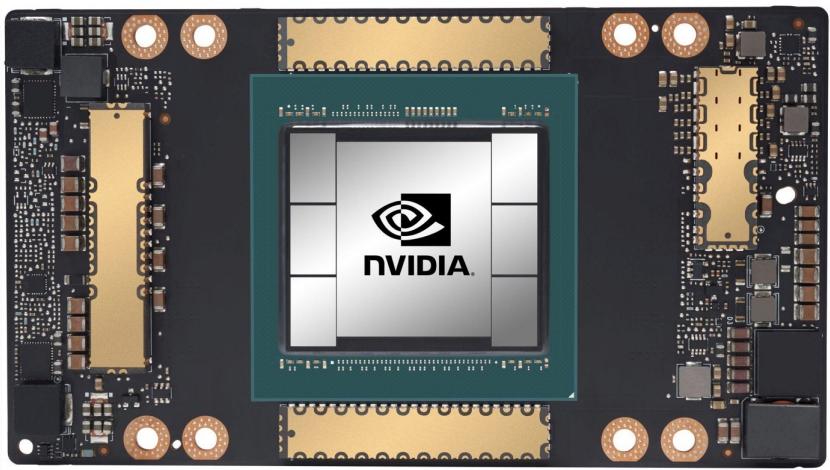
CUDA实践下的 GPU并行计算原理与机制介绍

博士后： 纪焘

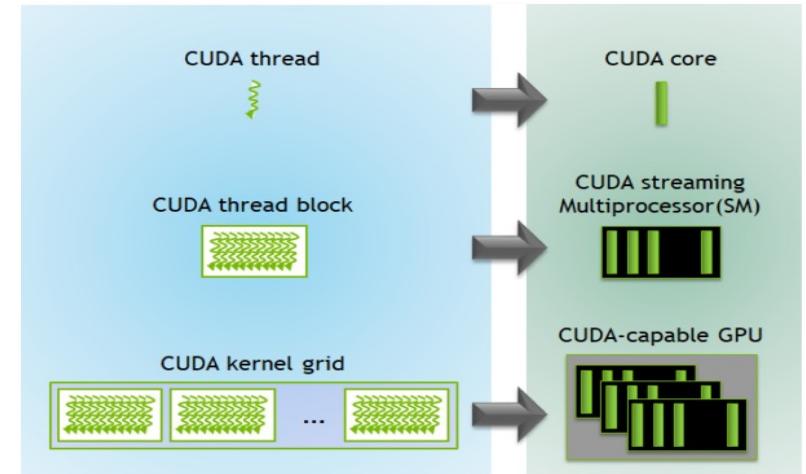
助教： 熊静飞、郭彬

taoji@fudan.edu.cn, jfxiong24@m.fudan.edu.cn, binguo@stu.ecnu.edu.cn

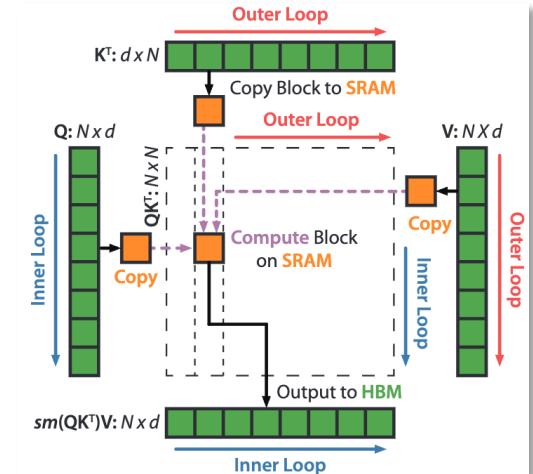
AI-Infra Seminar



GPU



CUDA、Triton、Cutlass



FlashAttention

Q&A time

- ▶ 17:15 ~ 18:00 (250715~250717)
- ▶ A5029, 江湾叉二
- ▶ #腾讯会议: [354-8238-6983](#)

Lecture

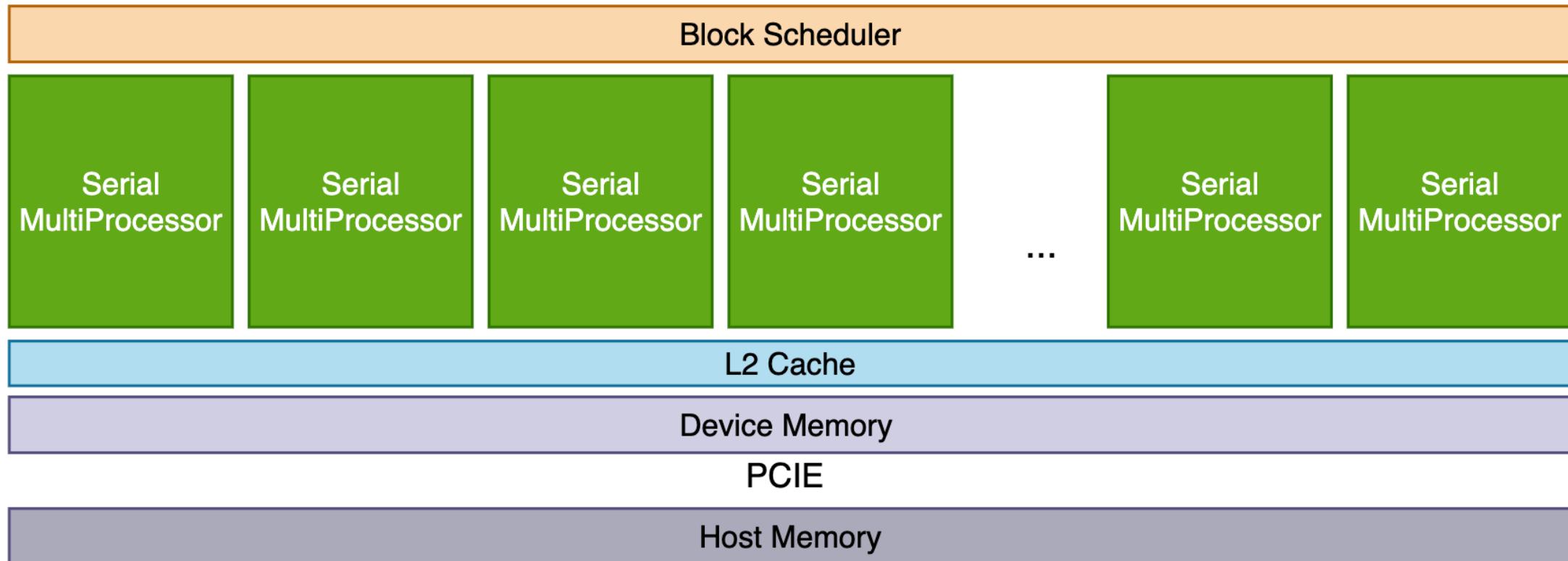
<https://github.com/JT-Ushio/AI-Infra-Seminar>

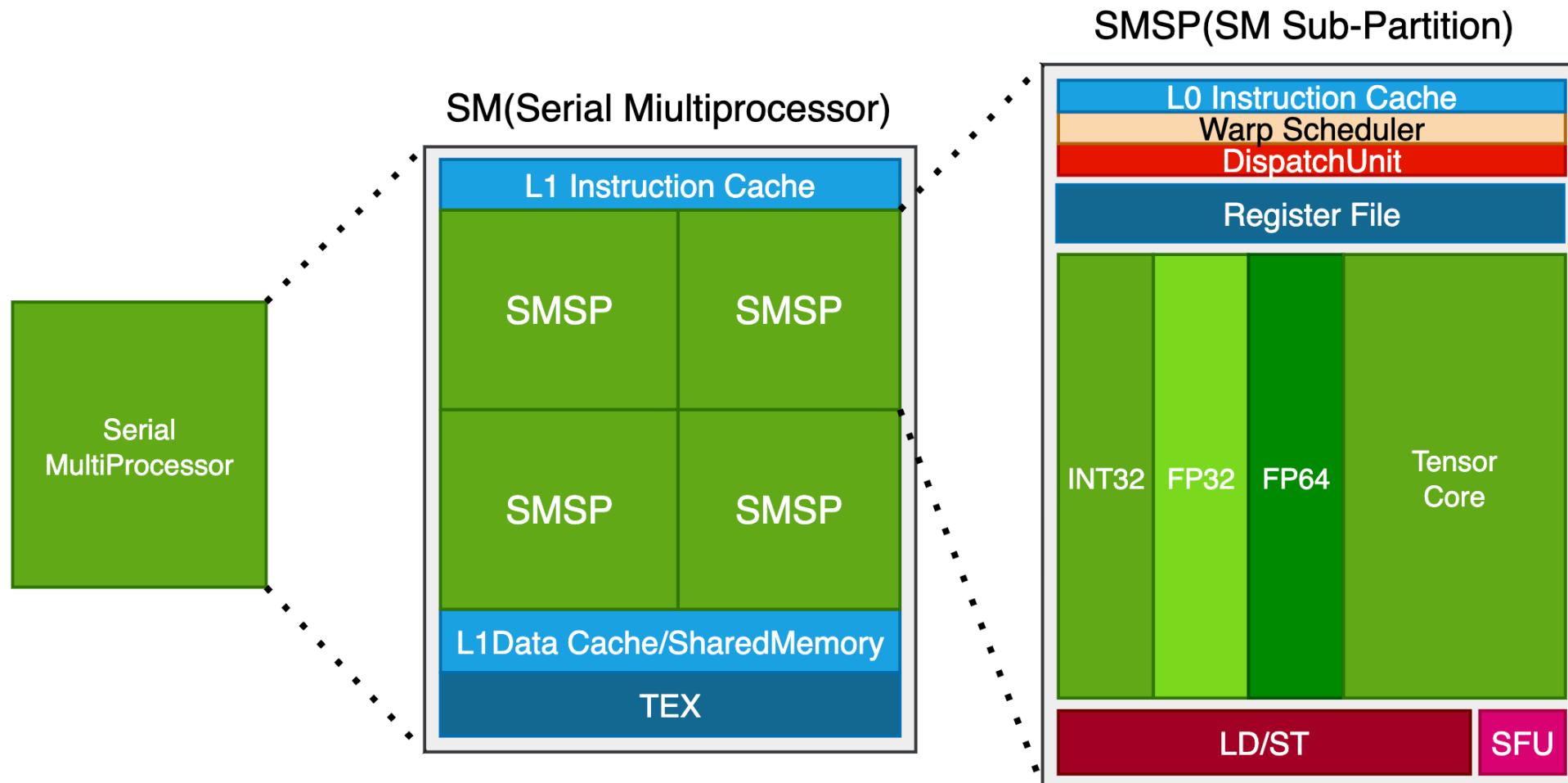


群聊: 并行计算&CUDA讨论班

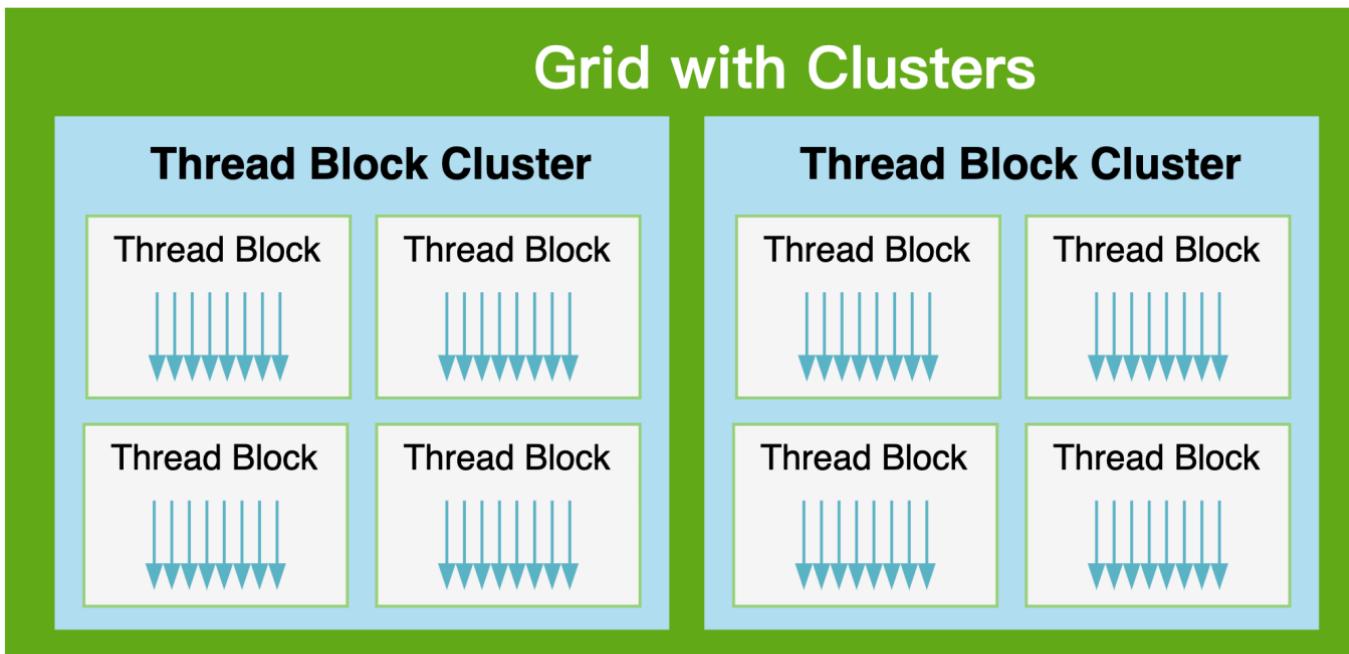


GPU





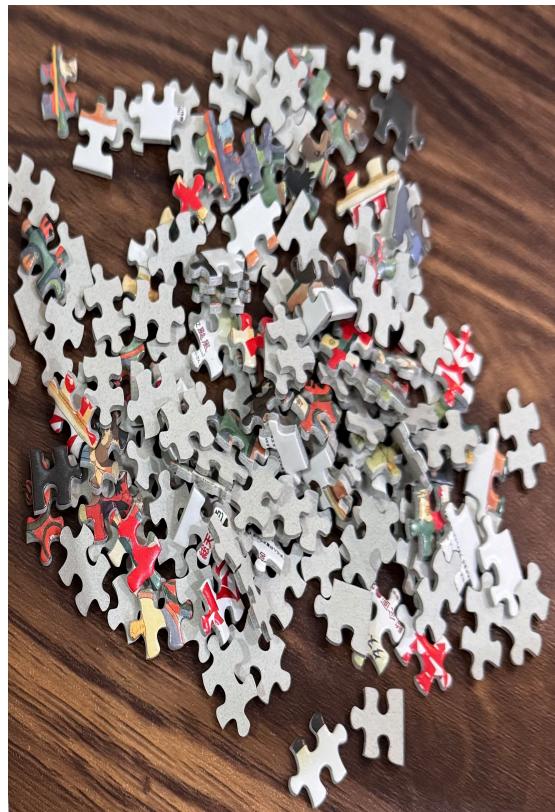
CUDA Programming



software concept	hardware concept
grid(Kernel) 108/114/132SMs	device
thread block cluster (TBC) 12~14SMs	GPU Processing Cluster (GPC)
thread block (CTA) <1024	Serial MultiProcessor(SM) <2048
warp 32*1	Serial MultiProcessor Sub-Partitiom (SMSP)
thread 1	cuda core

- ▶ Thread-level Programming
 - ▶ Single instruction, multiple threads (SIMT)

CUDA Programming: SIMD

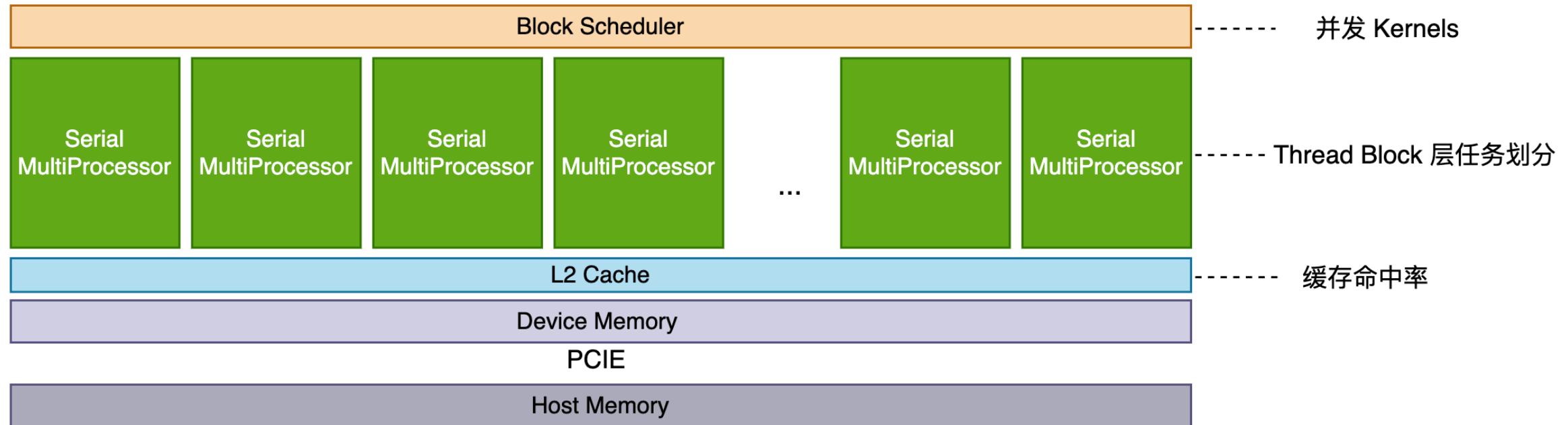


•
•
•



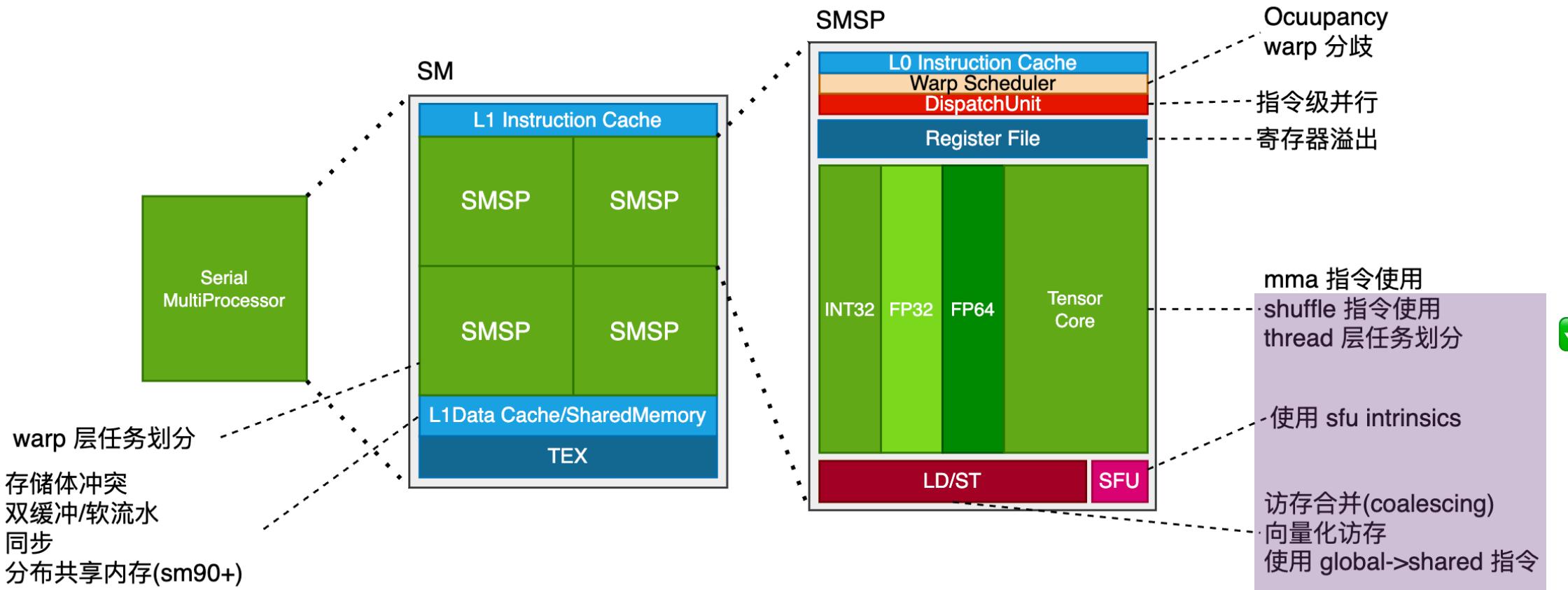
- ▶ ID
- ▶ $R, C = ID // \text{Col}, ID \% \text{Col}$
- ▶

CUDA Programming

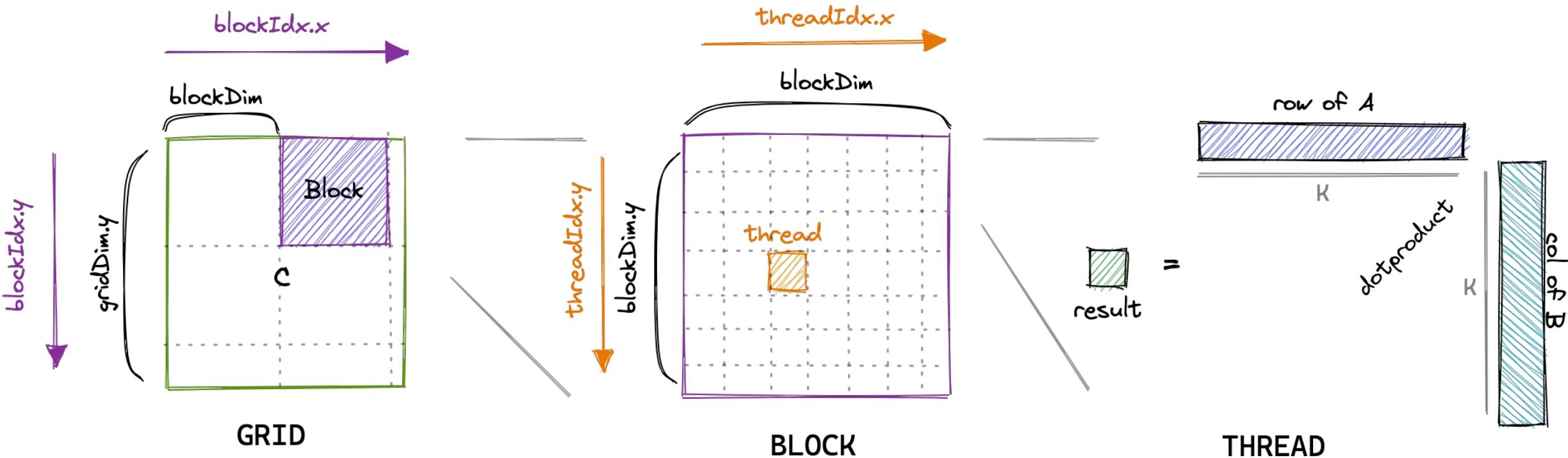


outside an SM

Conclusion



Matmul v0



Matmul v0



0_naive_matmul.py

```
1 module = load_inline(
2     name="matmul_naive_mod",
3     cpp_sources=cpp_source,
4     cuda_sources=cuda_source,
5     functions=["matmul_naive"],
6     verbose=True,
7 )
8
9 M, K, N = 1024, 1024, 1024
10 A = torch.randn(M, K, device="cuda").contiguous()
11 B = torch.randn(K, N, device="cuda").contiguous()
12
13 with profile(
14     activities=[ProfilerActivity.CUDA],
15     schedule=torch.profiler.schedule(wait=1, warmup=0, active=2),
16     on_trace_ready=torch.profiler.tensorboard_trace_handler("./log/matmul_naive")
17 , record_shapes=False,
18     with_stack=False,
19 ) as prof:
20     for _ in range(3):
21         with record_function("matmul_naive"), torch.no_grad():
22             C1 = module.matmul_naive(A, B)
23         with record_function("torch.matmul"), torch.no_grad():
24             C2 = torch.matmul(A, B)
25         prof.step()
26
27 assert torch.allclose(C1, C2, atol=1e-3, rtol=1e-3)
28 print(prof.key_averages().table(sort_by="cuda_time_total", row_limit=10))
29
```

Matmul v0

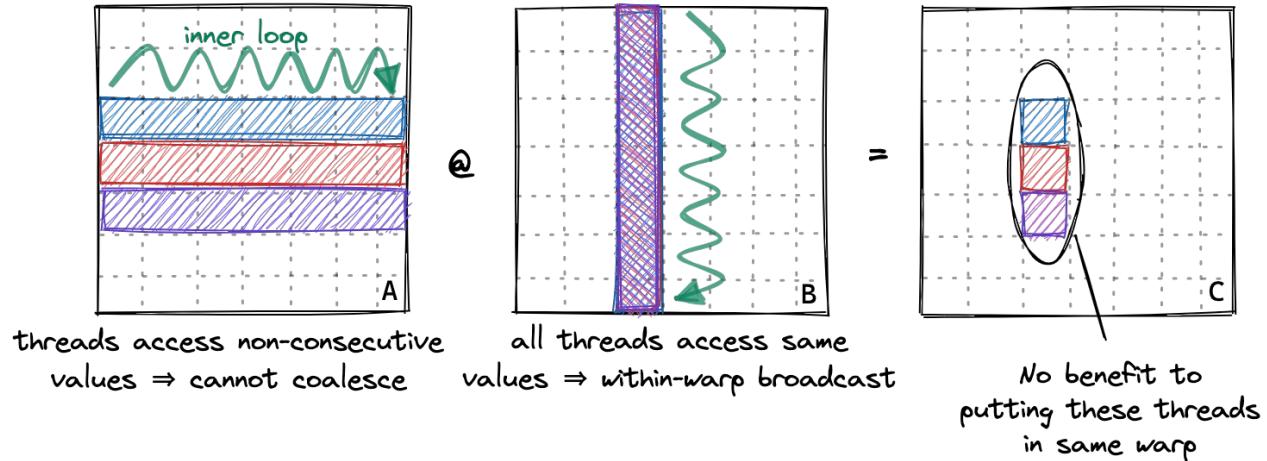


0_naive_matmul.py

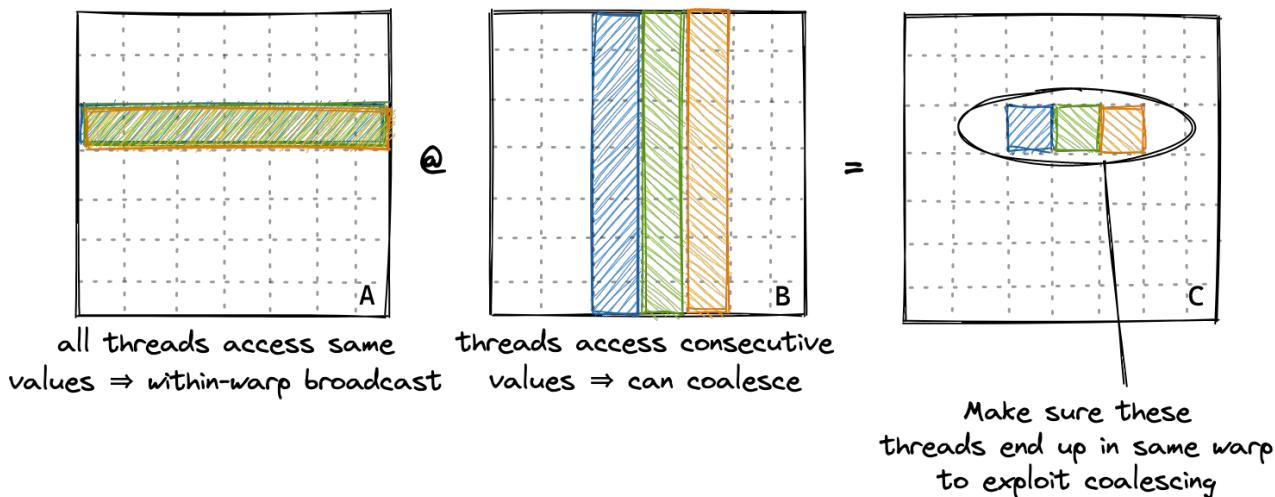
```
1  cuda_source = """  
2  __global__ void matmul_naive_kernel(const float* A, const float* B, float* C, int M, int N, int K) {  
3      int row = blockIdx.y * blockDim.y + threadIdx.y;  
4      int col = blockIdx.x * blockDim.x + threadIdx.x;  
5  
6      if (row < M && col < N) {  
7          float val = 0.0f;  
8          for (int k = 0; k < K; ++k) {  
9              val += A[row * K + k] * B[k * N + col];  
10         }  
11         C[row * N + col] = val;  
12     }  
13 }  
14  
15 torch::Tensor matmul_naive(torch::Tensor A, torch::Tensor B) {  
16     TORCH_CHECK(A.size(1) == B.size(0), "Incompatible matrix sizes");  
17     int M = A.size(0), K = A.size(1), N = B.size(1);  
18  
19     auto C = torch::empty({M, N}, A.options());  
20  
21     dim3 threads(32, 32);  
22     dim3 blocks((N + threads.x - 1) / threads.x, (M + threads.y - 1) / threads.y);  
23  
24     matmul_naive_kernel<<<blocks, threads>>>(br/>25         A.data_ptr<float>(), B.data_ptr<float>(), C.data_ptr<float>(), M, N, K  
26     );  
27  
28     return C;  
29 }  
30 """
```

Matmul v0

Naive kernel:



Coalescing kernel:

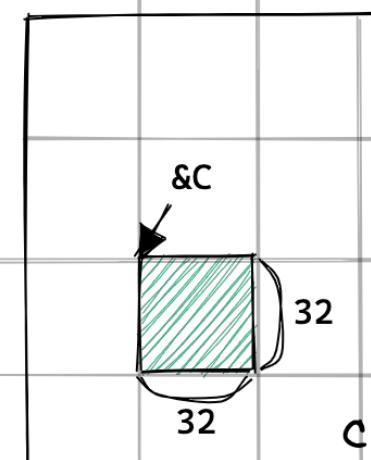
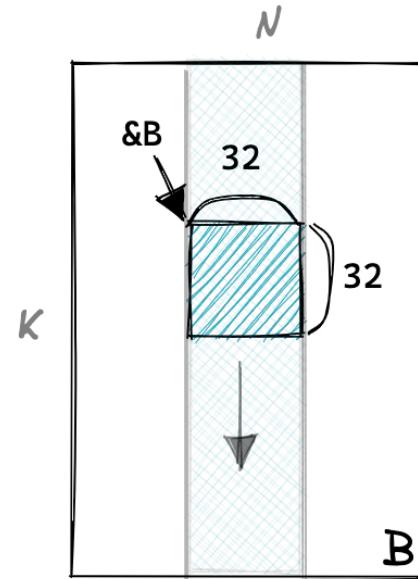
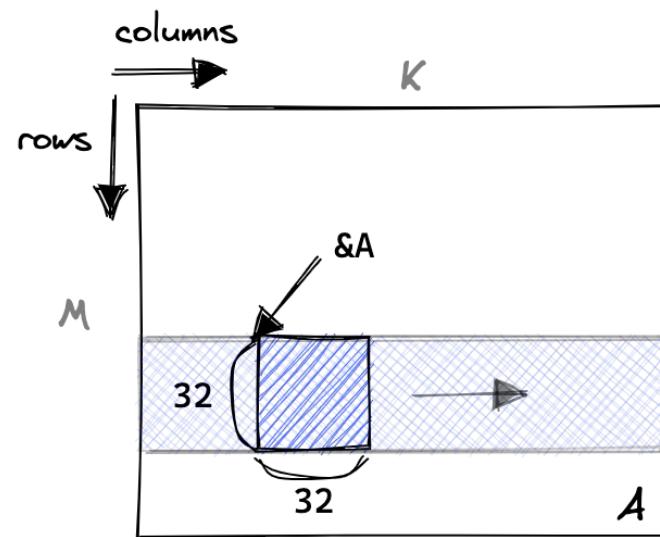


Matmul v0

Name	Self CPU %	Self CPU	CPU total %	CPU total	CPU time avg	Self CUDA	Self CUDA %	CUDA total	CUDA time avg
matmul_naive_kernel(float const*, float const*, flo...	0.00%	0.000us	0.00%	0.000us	0.000us	1.480ms	82.40%	1.480ms	740.228us
ampere_sgemm_128x64_nn	0.00%	0.000us	0.00%	0.000us	0.000us	313.440us	17.45%	313.440us	156.720us
Memset (Device)	0.00%	0.000us	0.00%	0.000us	0.000us	2.816us	0.16%	2.816us	1.408us
cudaLaunchKernel	4.29%	60.915us	4.29%	60.915us	15.229us	0.000us	0.00%	0.000us	0.000us
cudaMemsetAsync	1.35%	19.184us	1.35%	19.184us	9.592us	0.000us	0.00%	0.000us	0.000us
cudaOccupancyMaxActiveBlocksPerMultiprocessor	0.70%	10.011us	0.70%	10.011us	5.005us	0.000us	0.00%	0.000us	0.000us
cudaDeviceSynchronize	93.66%	1.331ms	93.66%	1.331ms	1.331ms	0.000us	0.00%	0.000us	0.000us

Matmul v1

1024*K*2



cCol=1

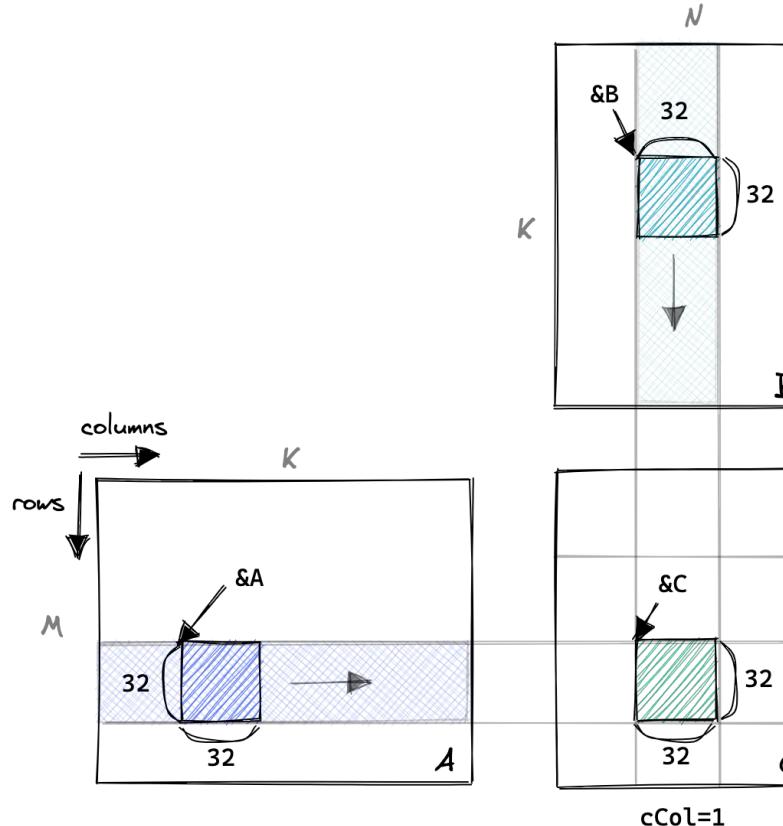
cRow=2

Outer loop:

Advance $\&A, \&B$ by size of cacheblock ($=32 \times 32$) until C is fully calculated

32*K*2
1024*K*2

Matmul v1

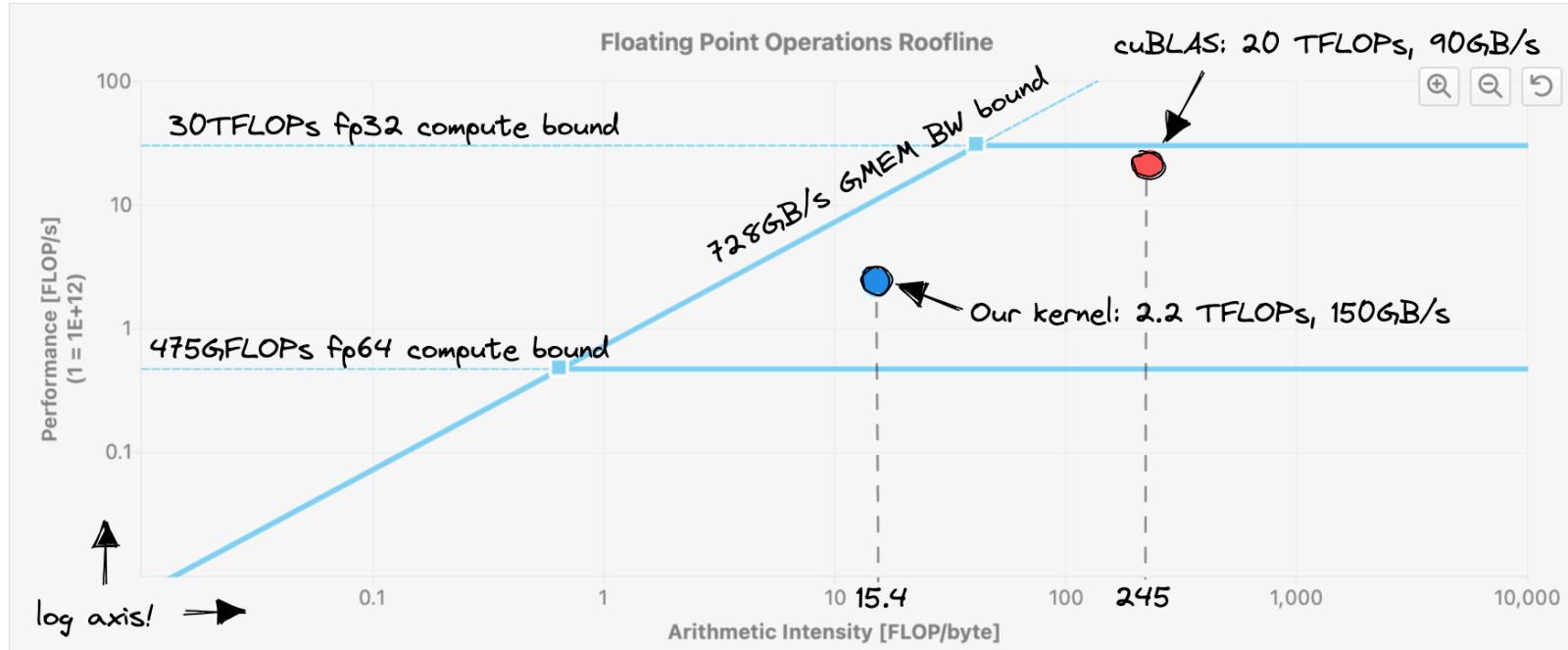


```
1_smem_matmul.py
1  cuda_source = """
2  __global__ void matmul_shared_kernel(const float* A, const float* B, float* C, int M, int N, int K) {
3      __shared__ float As[32][32];
4      __shared__ float Bs[32][32];
5
6      int row = blockIdx.y * blockDim.y + threadIdx.y;
7      int col = blockIdx.x * blockDim.x + threadIdx.x;
8
9      float acc = 0.0f;
10
11     for (int t = 0; t < (K + blockDim.x - 1) / blockDim.x; ++t) {
12         if (row < M && t * blockDim.x + threadIdx.x < K)
13             As[threadIdx.y][threadIdx.x] = A[row * K + t * blockDim.x + threadIdx.x];
14         else
15             As[threadIdx.y][threadIdx.x] = 0.0f;
16
17         if (col < N && t * blockDim.x + threadIdx.y < K)
18             Bs[threadIdx.y][threadIdx.x] = B[(t * blockDim.x + threadIdx.y) * N + col];
19         else
20             Bs[threadIdx.y][threadIdx.x] = 0.0f;
21
22         __syncthreads();
23
24         for (int k = 0; k < blockDim.x; ++k)
25             acc += As[threadIdx.y][k] * Bs[k][threadIdx.x];
26
27         __syncthreads();
28     }
29
30     if (row < M && col < N)
31         C[row * N + col] = acc;
32 }
33
34 torch::Tensor matmul_shared(torch::Tensor A, torch::Tensor B) {
35     TORCH_CHECK(A.size(1) == B.size(0), "Incompatible matrix sizes");
36     int M = A.size(0), K = A.size(1), N = B.size(1);
37
38     auto C = torch::empty({M, N}, A.options());
39
40     dim3 threads(32, 32);
41     dim3 blocks((N + 32 - 1) / 32, (M + 32 - 1) / 32);
42
43     matmul_shared_kernel<<<blocks, threads>>>(
44         A.data_ptr<float>(), B.data_ptr<float>(), C.data_ptr<float>(), M, N, K
45     );
46
47     return C;
48 }
49 """
```

Matmul v1

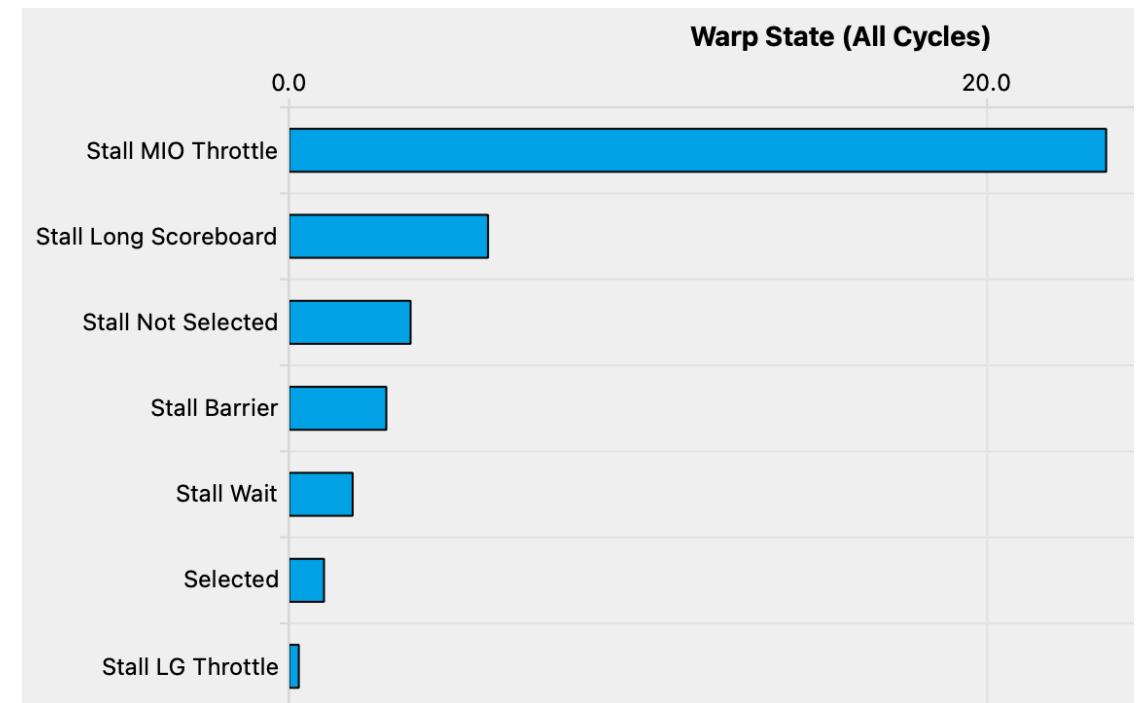
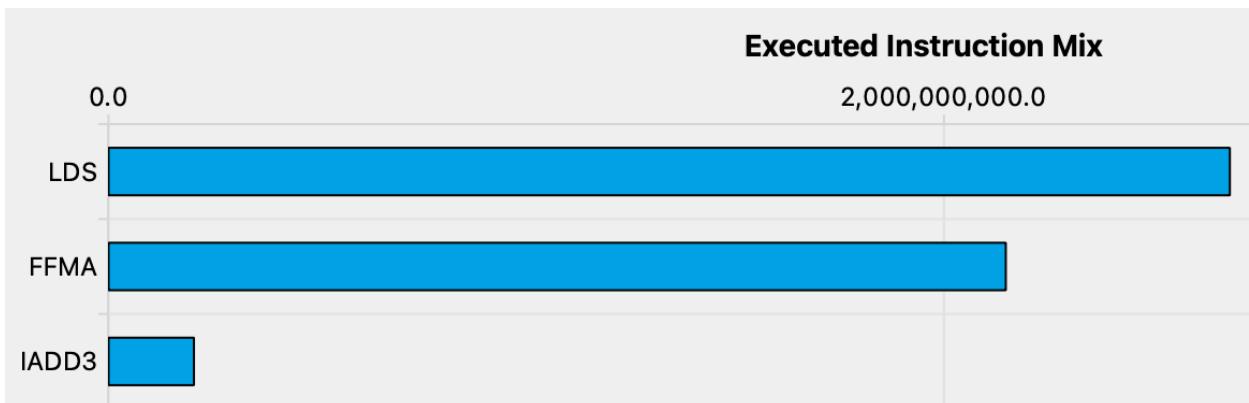
	Name	Self CPU %	Self CPU	CPU total %	CPU total	CPU time avg	Self CUDA	Self CUDA %	CUDA total	CUDA time avg	# of Calls
matmul_shared_kernel(float const*, float const*, flo...		0.00%	0.000us	0.00%	0.000us	0.000us	1.128ms	78.12%	1.128ms	563.923us	2
ampere_sgemm_128x64_nn		0.00%	0.000us	0.00%	0.000us	0.000us	313.025us	21.68%	313.025us	156.512us	2
Memset (Device)		0.00%	0.000us	0.00%	0.000us	0.000us	2.816us	0.20%	2.816us	1.408us	2
cudaLaunchKernel		6.03%	63.850us	6.03%	63.850us	15.962us	0.000us	0.00%	0.000us	0.000us	4
cudaMemsetAsync		1.80%	19.038us	1.80%	19.038us	9.519us	0.000us	0.00%	0.000us	0.000us	2
cudaOccupancyMaxActiveBlocksPerMultiprocessor		0.78%	8.212us	0.78%	8.212us	4.106us	0.000us	0.00%	0.000us	0.000us	2
cudaDeviceSynchronize		91.40%	968.317us	91.40%	968.317us	968.317us	0.000us	0.00%	0.000us	0.000us	1

Matmul v2



Matmul v2

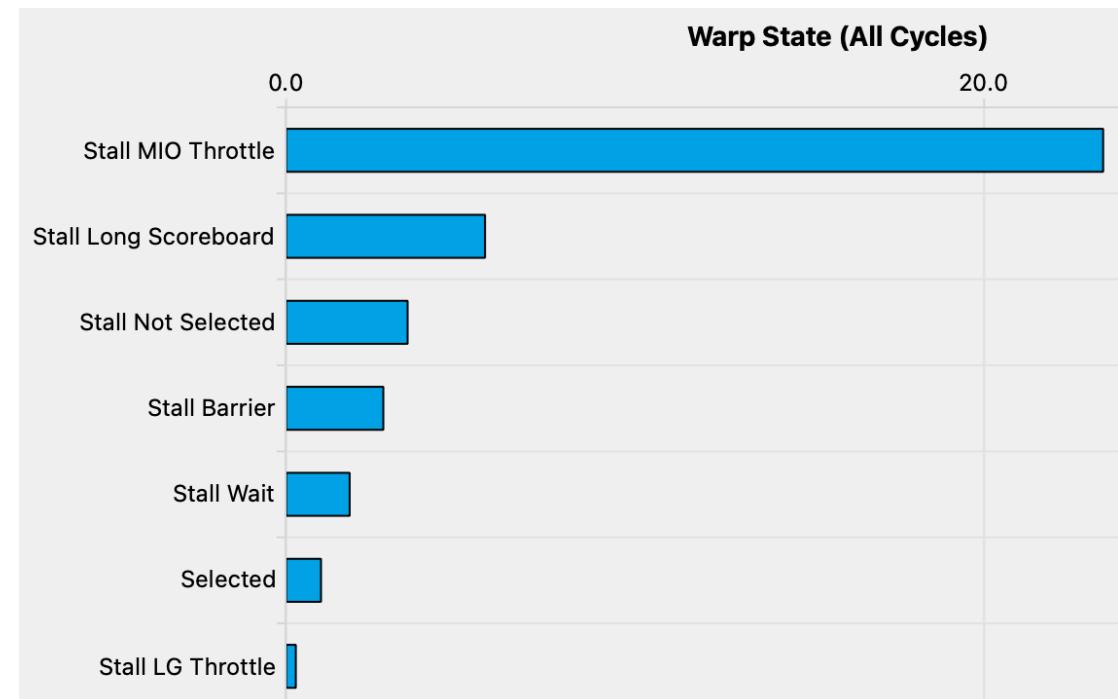
[NVIDIA Nsight Compute](#)



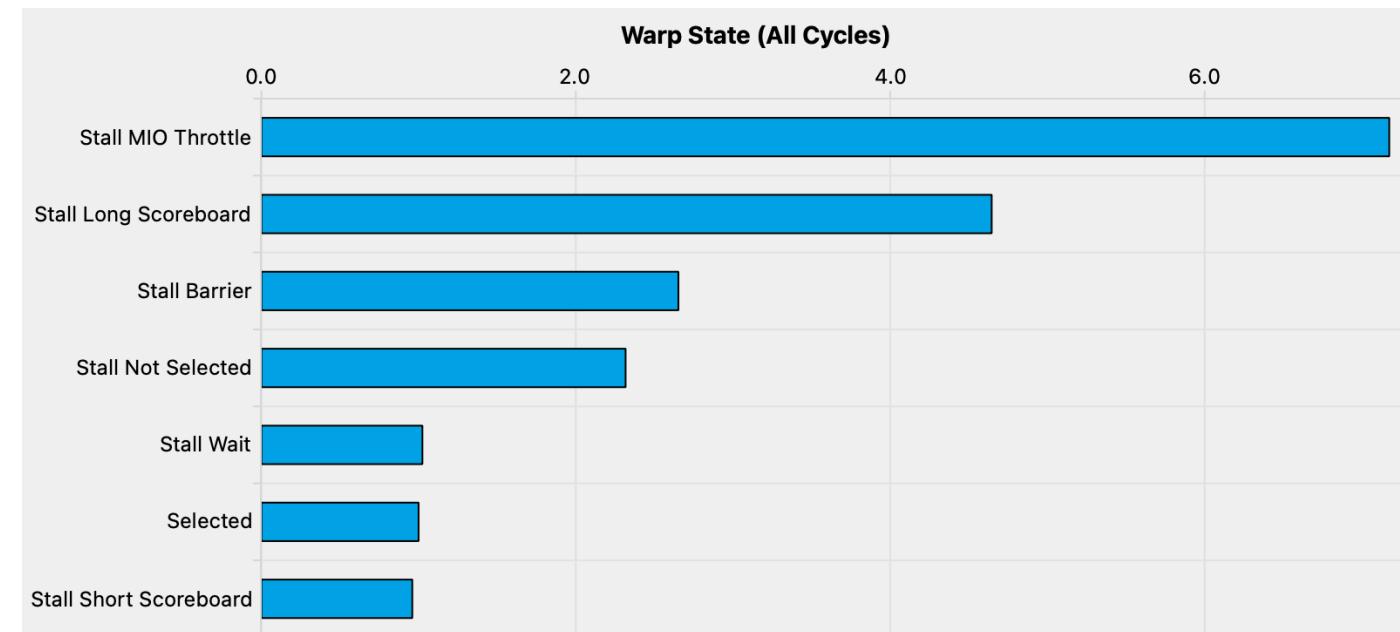
Stall Pipe Busy
Stall Long Scoreboard

Matmul v2

Thread任务：单线程计算单个C元素



Thread任务：单线程计算多个C元素

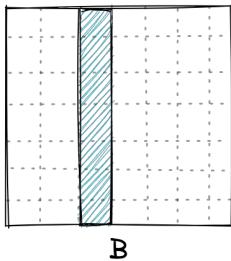
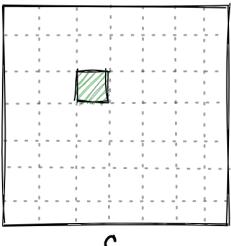
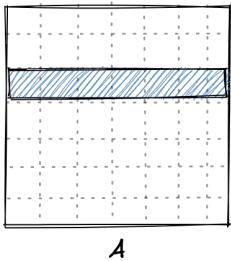


Matmul v2

calculating 1 result per thread requires:

- 7 loads from A
- 7 loads from B
- 1 load & 1 store to C

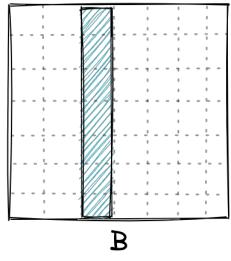
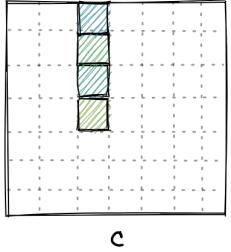
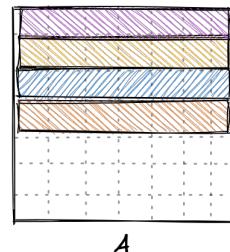
⇒ 15 loads & 1 store per result



calculating 4 results per thread requires:

- 28 loads from A
- 7 loads from B
- 4 loads & 4 stores to C

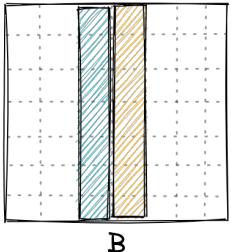
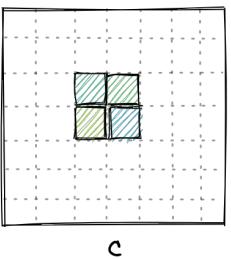
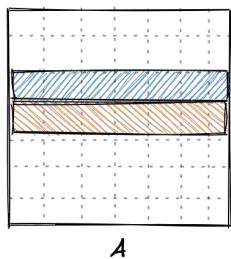
⇒ 9.75 loads & 1 store per result



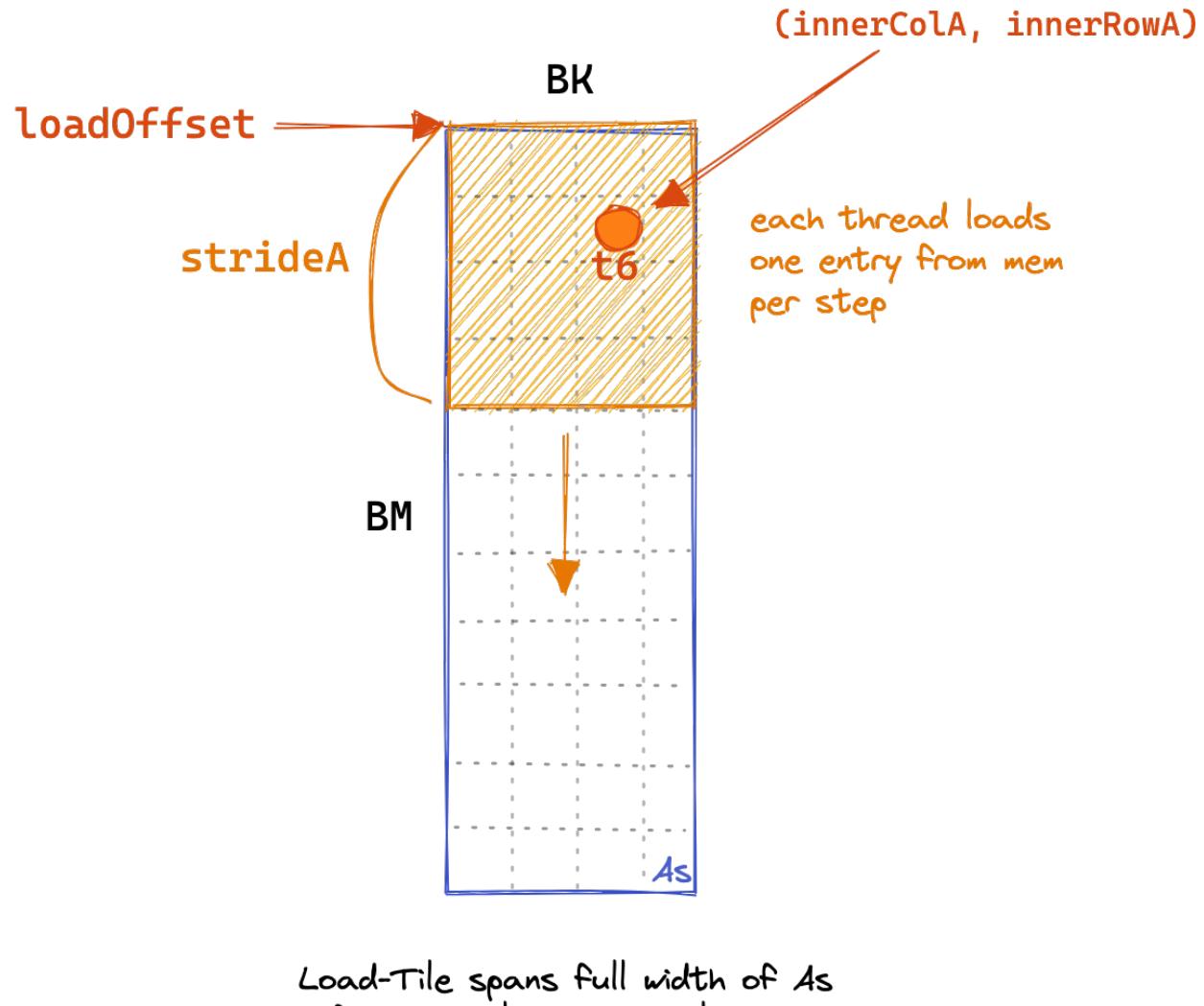
calculating 4 results per thread requires:

- 14 loads from A
- 14 loads from B
- 4 loads & 4 stores to C

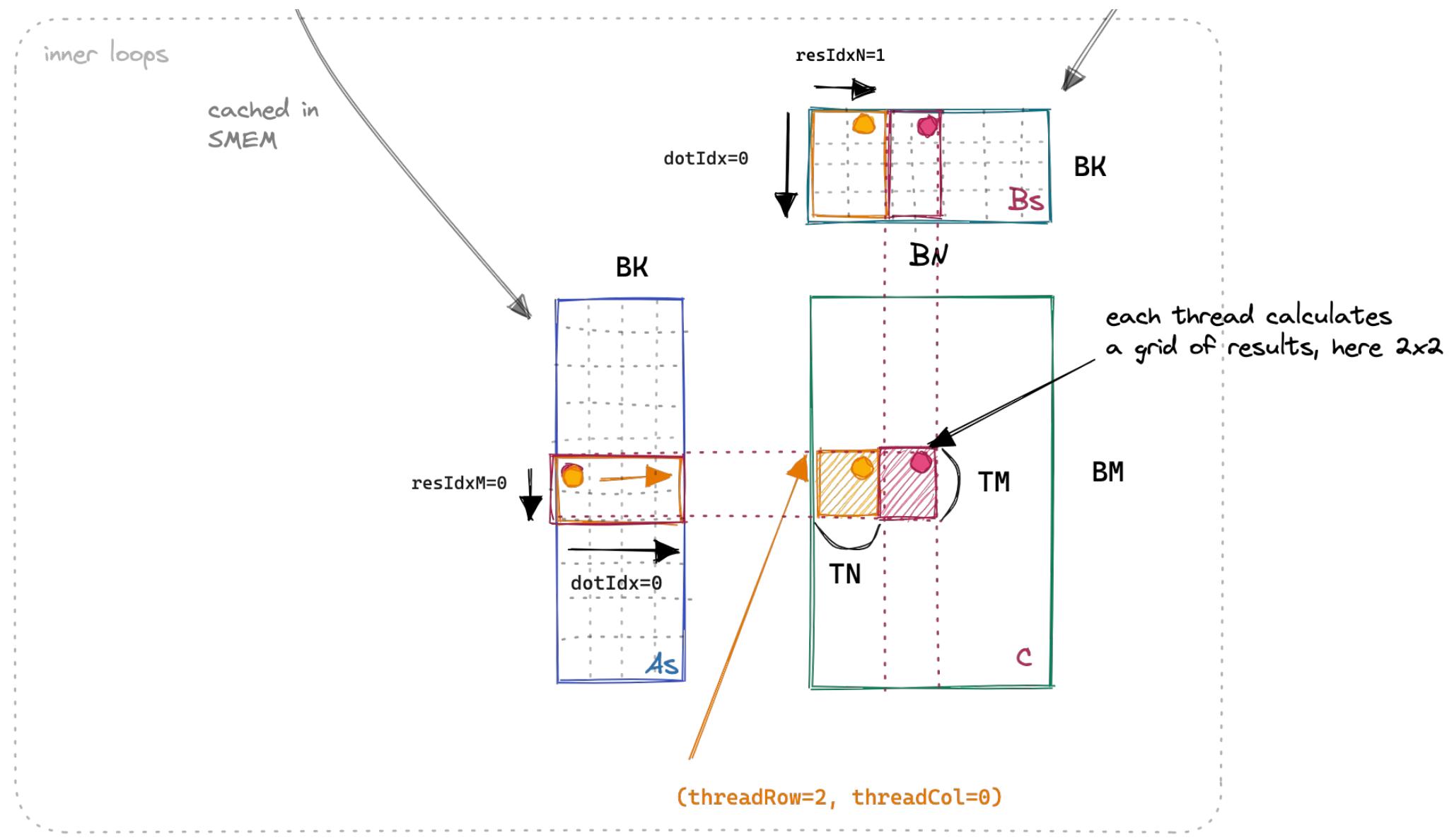
⇒ 8 loads & 1 store per result



Matmul v2



Matmul v2



Matmul v3



Matmul v3

GPU Features	NVIDIA Tesla P100	NVIDIA Tesla V100	NVIDIA A100
GPU Codename	GP100	GV100	GA100
GPU Architecture	NVIDIA Pascal	NVIDIA Volta	NVIDIA Ampere
Compute Capability	6.0	7.0	8.0
Threads / Warp	32	32	32
Max Warps / SM	64	64	64
Max Threads / SM	2048	2048	2048
Max Thread Blocks / SM	32	32	32
Max 32-bit Registers / SM	65536	65536	65536
Max Registers / Block	65536	65536	65536
Max Registers / Thread	255	255	255
Max Thread Block Size	1024	1024	1024
FP32 Cores / SM	64	64	64
Ratio of SM Registers to FP32 Cores	1024	1024	1024
Shared Memory Size / SM	64 KB	Configurable up to 96 KB	Configurable up to 164 KB

$$2048/64=32$$

$$65536/(72*64) = 14$$

$$164K/16K = 10$$

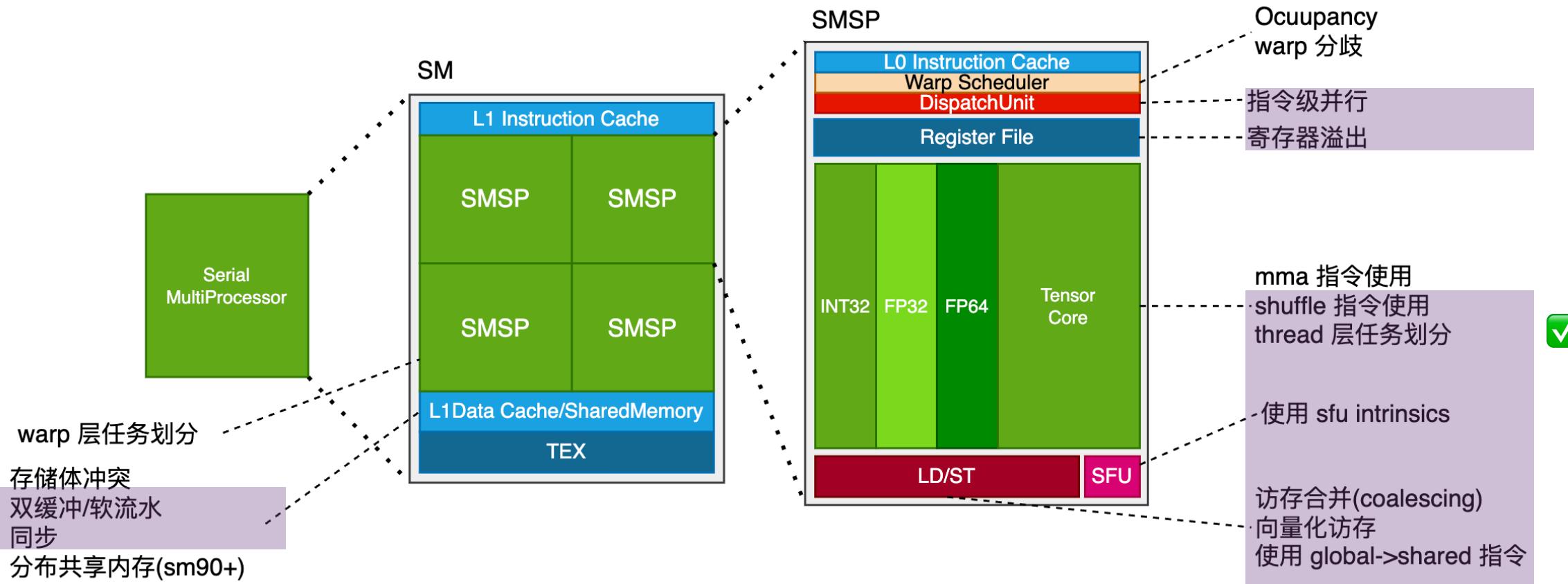
Matmul v4

类型	对应位宽
float2	64-bit
float4	128-bit
int2	64-bit
int4	128-bit
half2	32-bit
__half2	32-bit
bfloat162	32-bit
bfloat168 *	128-bit
char4	32-bit
uchar4	32-bit
ushort2	32-bit
自定义结构体	任意

Matmul

作业三：试试看大模型能生成什么速度的CUDA Matmul

Conclusion

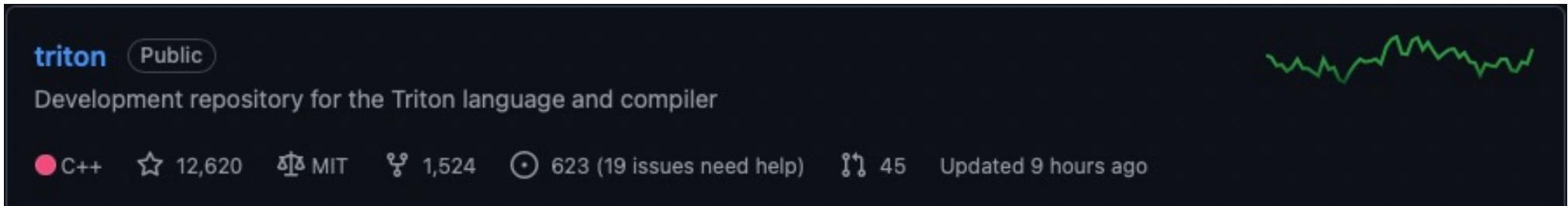


Triton: An Intermediate Language and Compiler for Tiled Neural Network Computations

Philippe Tillet
Harvard University
USA

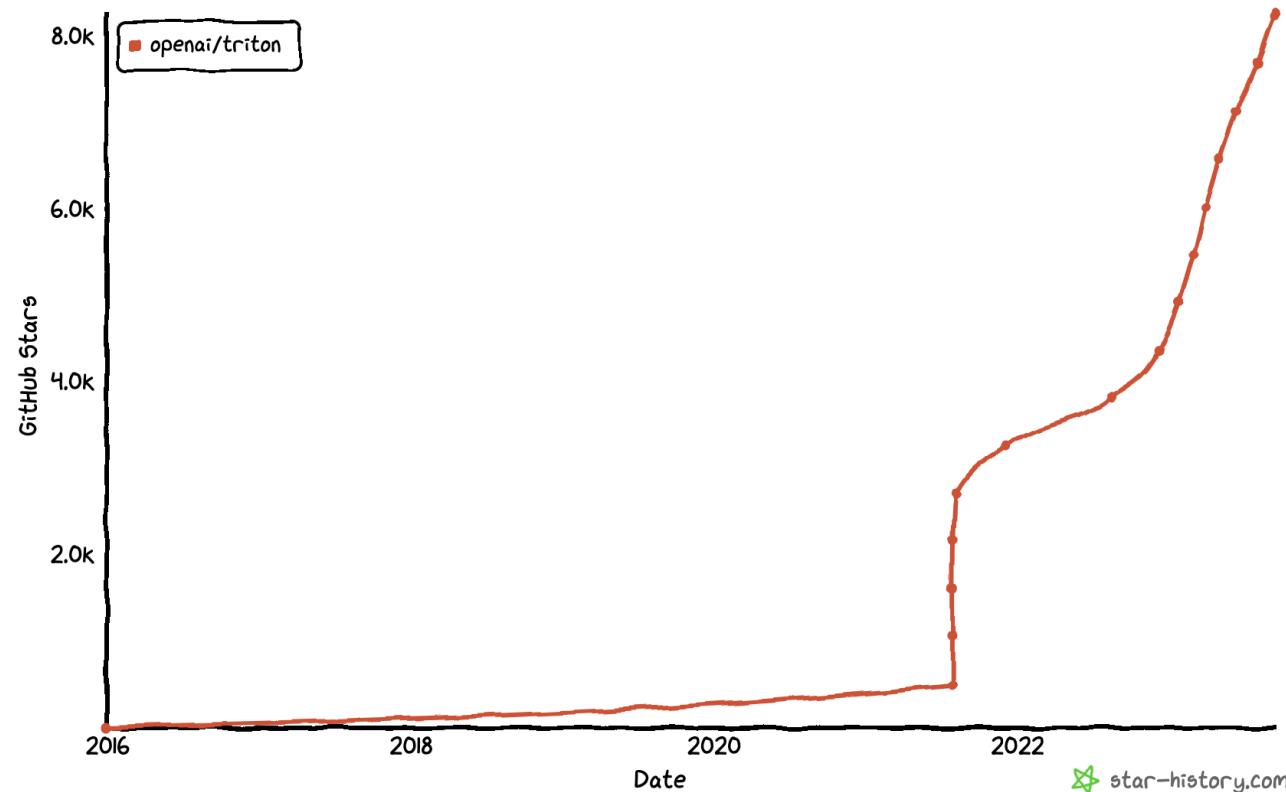
H. T. Kung
Harvard University
USA

David Cox
Harvard University, IBM
USA



Triton

- ▶ Triton是一个用于编写高效自定义深度学习算子的语言和编译器。相比CUDA，Triton生产力更高、更灵活。
- ▶ 2022年开始，Triton的热度呈爆发式增长，Triton生态逐渐形成规模。



Triton

- ▶ Triton 是一门 kernel 编程语言，主要用于编写 GPU kernels，也可以作为 AI 编译器的 Kernel 代码生成目标语言。和 CUDA, OpenCL 等属于并列的关系。

Model	Pytorch/Tensorflow/JAX
Graph	XLA/HLO TVM/Relay Pytorch/fx
Kernel	CUDA OpenCL ROCM HIP SYCL Triton
Device	GPUs/CPUs/...



Triton vs. CUDA



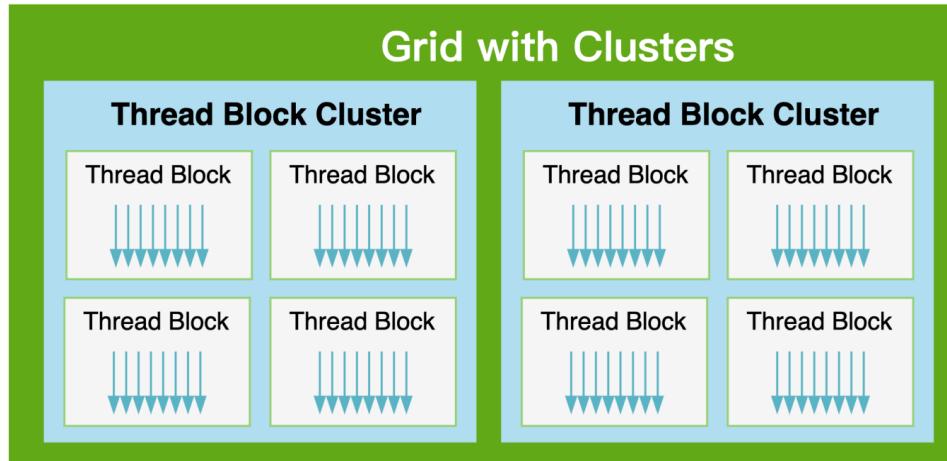
CUDA



Triton

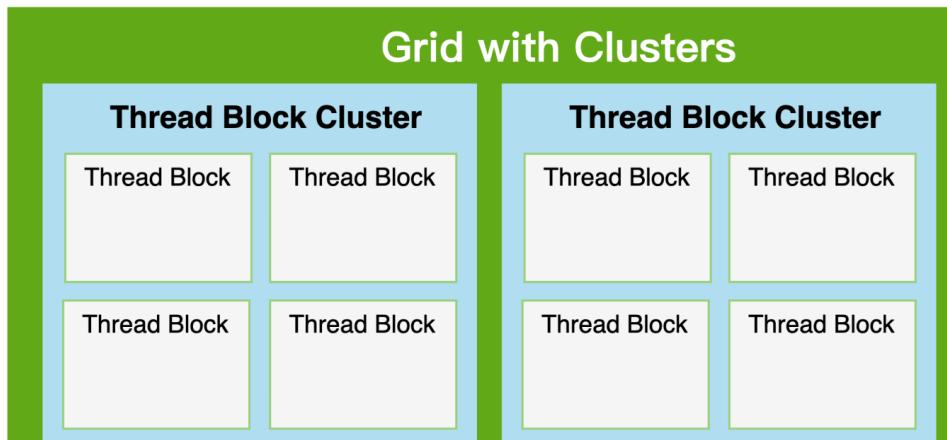
Triton vs. CUDA

► Thread-level Programming, Single instruction, multiple threads (SIMT)



software concept	hardware concept
grid(Kernel)	device
thread block cluster (TBC)	GPU Processing Cluster (GPC)
thread block (CTA)	Serial MultiProcessor(SM)
warp	Serial MultiProcessor Sub-Partitiom (SMSP)
thread	cuda core

► Block-level Programming, Single instruction, multiple data (SIMD)



software concept	hardware concept
Kernel	device
cluster	GPU Processing Cluster (GPC)
program(CTA)	Serial MultiProcessor(SM)

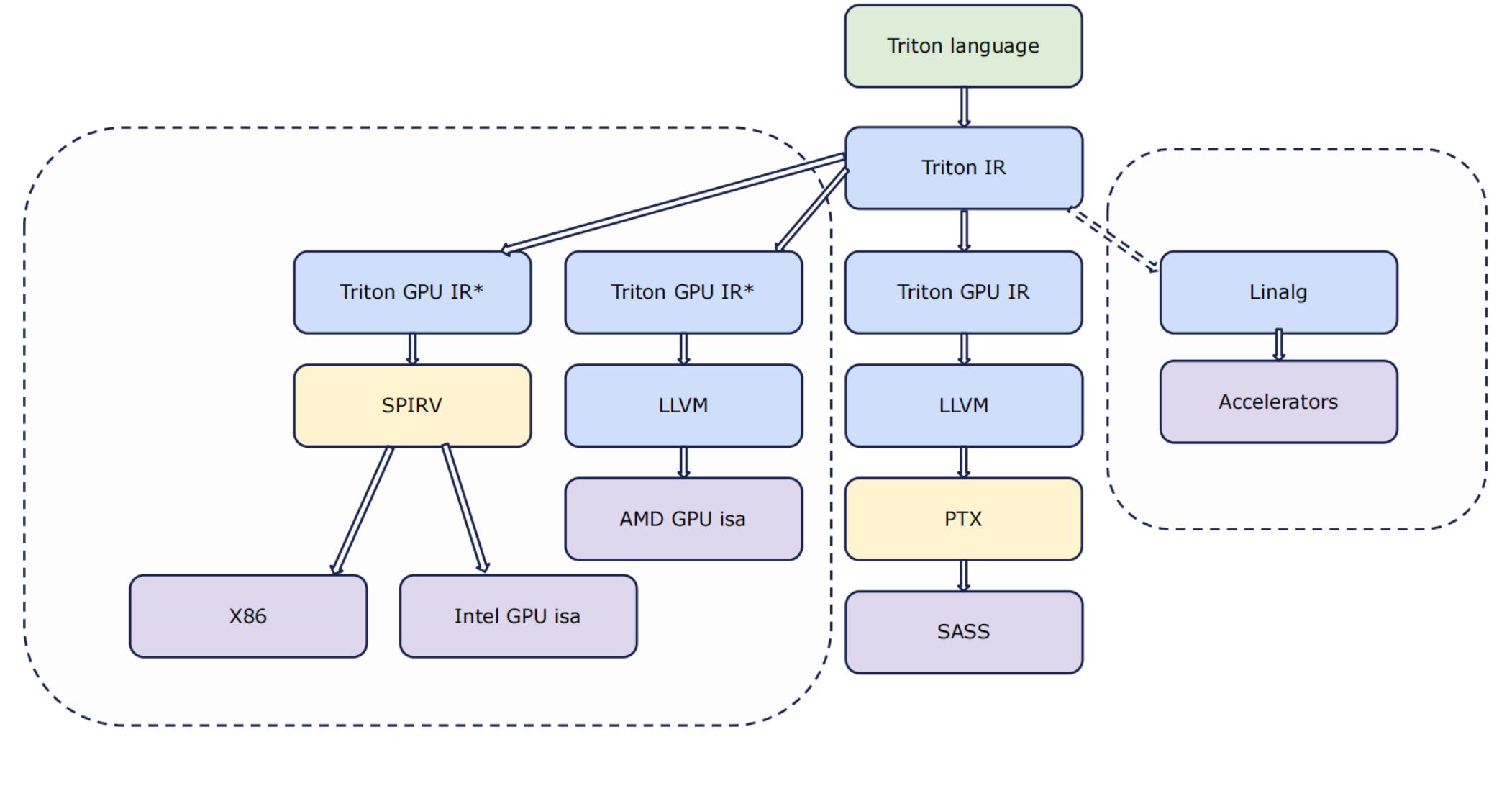
Triton vs. CUDA

	CUDA	Triton	Torch Op
Algorithm	User	User	Compiler
Shared Memory	User	Compiler	Compiler
Barriers	User	Compiler	Compiler
Distribution to blocks	User	User	Compiler
Grid size	User	User	Compiler
Distribution to Warps/threads	User	Compiler	Compiler
Tensor Core usage	User	Compiler	Compiler
Coalescing	User	Compiler	Compiler
Intermediate data layout	User	Compiler	Compiler
Workgroup size	User	User	Compiler

Why Triton?

- ▶ 更高抽象层级的编程模型，更低的心智负担。
- ▶ Block 级别编程，SPMD 编程模型，更简单的任务划分；
- ▶ 定义在 SRAM (寄存器和 shared memory)上的 tensor；
- ▶ 用于 tensor 的 primitive.(reduce, scan, sort, dot, ...);
- ▶ triton compiler 的优化使非 GPU/CUDA 编程专家也能写出性能不俗的 kernel.
- ▶ Triton 编译器代码开源，对多元 AI 芯片有更好的可移植性。

Why Triton?



Triton vs. CUDA

block_id	thread_id	offs	offs < n	x_value	y_value	z_value
0	0	0	True	1	0	1
0	1	1	True	2	1	3
0	2	2	True	3	0	3
0	3	3	True	4	1	5
1	0	4	True	5	0	5
1	1	5	True	6	1	7
1	2	6	False	undefined	undefined	undefined
1	3	7	False	undefined	undefined	undefined

block_id	offs	offs < n	x_value	y_value	z_value
0	[0, 1, 2, 3]	[True, True, True, True]	[1, 2, 3, 4]	[0, 1, 0, 1]	[1, 3, 3, 5]
1	[4, 5, 6, 7]	[True, True, False, False]	[5, 6, 'undefined', 'undefined']	[0, 1, 'undefined', 'undefined']	[5, 7, 'undefined', 'undefined']

Triton Index

index	name	tags	link
1	Triton RMSNorm + Residual	norm	link
2	Flash-Decoding / Split-KV Attention	attention, decoding	link
3	Cross entropy loss	loss	link
4	Geglu	Activation	link
5	RMSNorm	Normalization	link
6	Rope Embedding	Embedding	link
7	Swiglu	Activation	link
8	Gemm Int4	Quant, Dequant, Matmul, Gemm	link
9	Gemm Int8	Quant, Dequant, Matmul, Gemm	link
10	Batched Conv2D	Conv	link
11	Batched Addition	Add	link
12	Batched Patchng	Rearranging	link
13	Block-Sparse Matrix Multiplication	Sparsity, Matmul	link
14	Fused Linear Cross Entropy	Linear Cross Entropy	link
n+1	your contribution	its tags	

算子开发示例: Vector Add (1 tile/CTA)

1 tile/CTA

CTA 0	CTA 1	CTA 2	CTA 3	CTA 4
TILE_SIZE				

```
@triton.jit
def add_kernel(a_ptr, b_ptr, o_ptr, N, TILE_SIZE: tl.constexpr):

    pid = tl.program_id(0)
    offsets = pid * TILE_SIZE + tl.arange(0, TILE_SIZE)
    mask = offsets < N

    a = tl.load(a_ptr + offsets, mask=mask)
    b = tl.load(b_ptr + offsets, mask=mask)
    o = a + b

    tl.store(o_ptr + offsets, o, mask=mask)
```

```
def add(a, b):
    o = torch.empty_like(a)
    N = a.numel()

    TILE_SIZE = 128
    grid = (triton.cdiv(N, TILE_SIZE), 1, 1)
    add_kernel[grid](a, b, o, N, TILE_SIZE, num_warps=4)

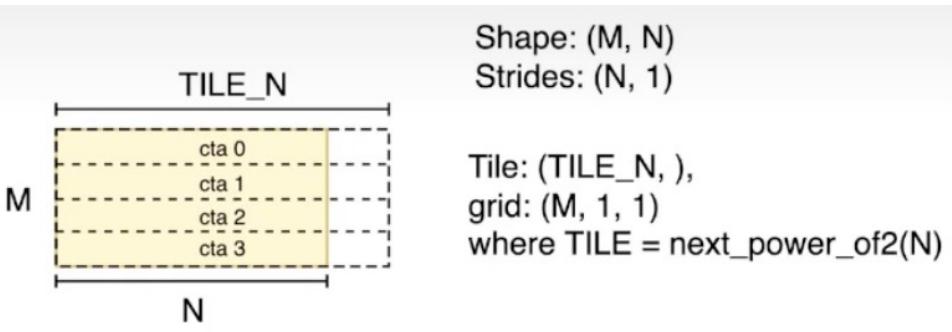
    return o
```

- ▶ Kernel: 描述每个CTA的运算
- ▶ 通过program_id区分CTA
- ▶ 通过range生成tensor

- ▶ Wrapper: 指定grid, 启动kernel

算子开发示例: Softmax

$$\text{Softmax}(x_i) = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$

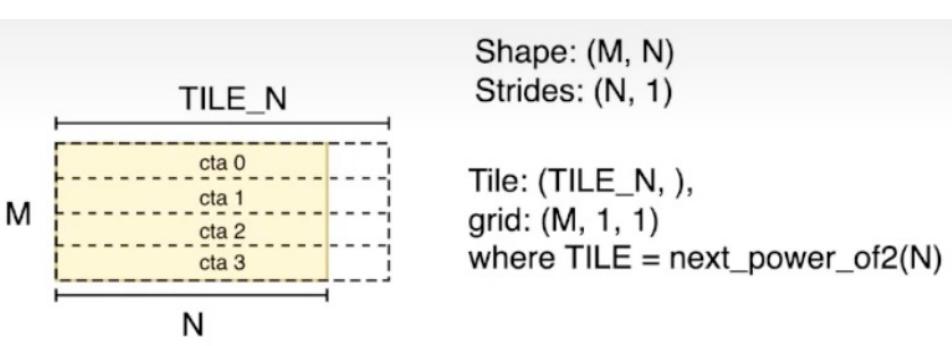


0_naive_softmax.py

```
1 x = torch.randn((4096, 32768), device="cuda")
2 # x = torch.randn((4096, 8192), device="cuda")
3 # x = torch.randn((4096, 4096), device="cuda")
4
5 with profile(
6     activities=[ProfilerActivity.CUDA],
7     schedule=torch.profiler.schedule(wait=1, warmup=1, active=3),
8     on_trace_ready=torch.profiler.tensorboard_trace_handler("./log/naive_softmax"),
9     record_shapes=False,
10    with_stack=False,
11 ) as prof:
12     for i in range(5):
13         with torch.no_grad():
14             out1 = softmax(x)
15         with torch.no_grad():
16             out0 = torch.softmax(x, dim=-1)
17         prof.step()
18
19     out1 = softmax(x)
20     out0 = torch.softmax(x, dim=-1)
21     assert torch.allclose(out1, out0, atol=1e-3)
22     print(prof.key_averages().table(sort_by="cuda_time_total", row_limit=10))
```

算子开发示例: Softmax

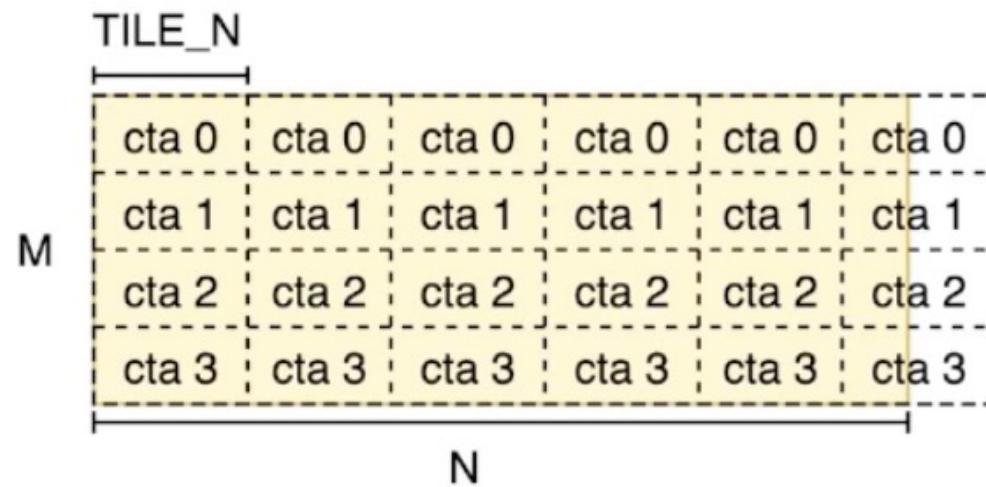
$$\text{Softmax}(x_i) = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$



0_naive_softmax.py

```
1  @triton.jit
2  def softmax_kernel(output_ptr, input_ptr, M, N, TILE_N: tl.constexpr):
3      pid_m = tl.program_id(0) # (0...M-1, 0, 0)
4      n_offsets = tl.arange(0, TILE_N) # [0, ..., TILE_N-1]
5      offset = pid_m * N + n_offsets # (pid_m, [0, ..., TILE_N-1])
6
7      x = tl.load(input_ptr + offset, mask=n_offsets < N, other=-float("inf")
8      ) m = tl.max(x)
9      e = tl.exp(x - m)
10     z = tl.sum(e)
11     out = e / z
12     tl.store(output_ptr + offset, out, mask=n_offsets < N)
13
14
15 def softmax(x):
16     M, N = x.shape
17     out = torch.empty_like(x)
18     TILE_N = triton.next_power_of_2(N)
19     grid = (M, 1, 1)
20     softmax_kernel[grid](out, x, M, N, TILE_N)
21     return out
```

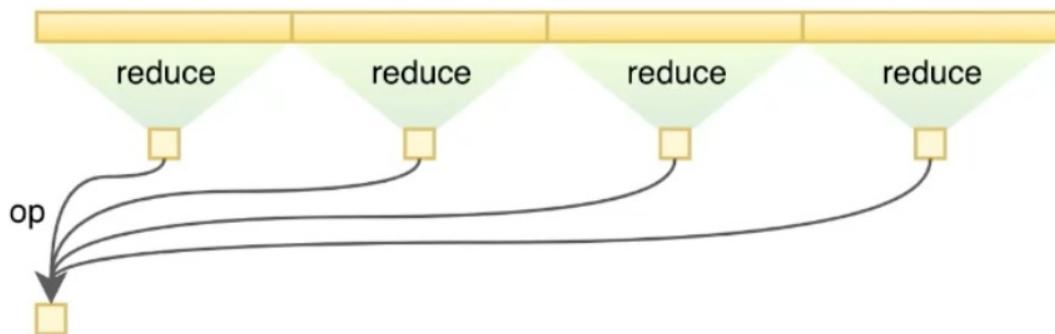
算子开发示例: Softmax



1_multi_cta_softmax.py

```
1 def softmax(x):
2     M, N = x.shape
3     out = torch.empty_like(x)
+ 4     TILE_N = min(4096, triton.next_power_of_2(N))
5     grid = (M, 1, 1)
6     softmax_kernel_loop_v1[grid](out, x, M, N, TILE_N)
7 ) return out
```

算子开发示例: Softmax



1_multi_cta_softmax.py

```
1  @triton.jit
2  def softmax_kernel_loop_v1(output_ptr, input_ptr, M, N, TILE_N: tl.constexpr):
3      pid_m = tl.program_id(0) # (0..M-1, 0, 0)
4
5      m = tl.full(() , value=-float("inf") , dtype=output_ptr.dtype.element_ty)
+ 6      for start_cta in range(0, N, TILE_N): # num of CTA
7          offsets_cta = start_cta + tl.arange(0, TILE_N)
8          offsets_x = pid_m * N + offsets_cta
9          x = tl.load(input_ptr + offsets_x, mask=offsets_cta < N, other=-float("inf"))
10         m = tl.maximum(m, tl.max(x))
11
12     z = tl.full(() , value=0, dtype=output_ptr.dtype.element_ty)
+ 13    for start_cta in range(0, N, TILE_N): # num of CTA
14        offsets_cta = start_cta + tl.arange(0, TILE_N)
15        offsets_x = pid_m * N + offsets_cta
16        x = tl.load(input_ptr + offsets_x, mask=offsets_cta < N, other=-float("inf"))
17        e = tl.exp(x - m)
18        z += tl.sum(e)
19
+ 20    for start_cta in range(0, N, TILE_N): # num of CTA
21        offsets_cta = start_cta + tl.arange(0, TILE_N)
22        offsets_x = pid_m * N + offsets_cta
23        x = tl.load(input_ptr + offsets_x, mask=offsets_cta < N, other=-float("inf"))
24        e = tl.exp(x - m)
25        out = e / z
26        tl.store(output_ptr + offsets_x, out, mask=offsets_cta < N)
27
```

算子开发示例: Softmax

$x: [3, 2, 5, 1]$

$m0=-\inf$

$z0=0$

$m1=3$

$z1=e(3-5)$

$m2=3$

$z2=z1+e(2-5)$

$m3=5$

$z3=z2+e(5-5)$

$m4=5$

$z4=z3+e(1-5)$

算子开发示例: Softmax

CTA i

$m = \max_{CTA_0_i}$
 $z = \sum_{CTA_0_i} \exp(x)$

CTA i+1

$m^* = \max(m, \max_{CTA_i+1})$
 $z^* = z \times \exp(m - m^*) + z_{CTA_i+1}$



2_online_softmax.py

```
1 @triton.jit
2 def softmax_kernel_online_v1(output_ptr, input_ptr, M, N, TILE_N: tl.constexpr):
3     pid_m = tl.program_id(0)
4     m = tl.full(() , value=-float("inf") , dtype=output_ptr.dtype.element_ty)
5     z = tl.full(() , value=0 , dtype=output_ptr.dtype.element_ty)
6     for start_cta in range(0, N, TILE_N):
7         offsets_cta = start_cta + tl.arange(0, TILE_N)
8         offsets_x = pid_m * N + offsets_cta
9         x = tl.load(input_ptr + offsets_x, mask=offsets_cta < N, other=-float("inf"))
+ 10        new_m = tl.maximum(m, tl.max(x))
+ 11        new_z = tl.exp(m - new_m) * z + tl.sum(tl.exp(x - new_m))
12        m = new_m
13        z = new_z
14
15    for start_cta in range(0, N, TILE_N):
16        offsets_cta = start_cta + tl.arange(0, TILE_N)
17        offsets_x = pid_m * N + offsets_cta
18        x = tl.load(input_ptr + offsets_x, mask=offsets_cta < N, other=-float("inf"))
19        e = tl.exp(x - m)
20        out = e / z
21        tl.store(output_ptr + offsets_x, out, mask=offsets_cta < N)
```

算子开发示例: Softmax

```
● ● ● 3_online_softmax_v2.py  
+ 1  @triton.jit  
+ 2  def prev_multiple_of(a, b):  
+ 3      return tl.cdiv(a, b) * b - b  
+ 4  
+ 5  @triton.jit  
+ 6  def softmax_kernel_online_v1(output_ptr, input_ptr, M, N, TILE_N: tl.constexpr):  
+ 7      pid_m = tl.program_id(0)  
+ 8      m = tl.full(() , value=-float("inf"), dtype=output_ptr.dtype.element_ty)  
+ 9      z = tl.full(() , value=0, dtype=output_ptr.dtype.element_ty)  
+10     before_last = prev_multiple_of(N, TILE_N)  
+11     for start_cta in range(0, before_last, TILE_N):  
+12         offsets_cta = start_cta + tl.arange(0, TILE_N)  
+13         offsets_x = pid_m * N + offsets_cta  
+14         x = tl.load(input_ptr + offsets_x)  
-15         # x = tl.load(input_ptr + offsets_x, mask=offsets_cta < N, other=-float("inf"))  
+16         new_m = tl.maximum(m, tl.max(x))  
+17         new_z = tl.exp(m - new_m) * z + tl.sum(tl.exp(x - new_m))  
+18         m = new_m  
+19         z = new_z  
+20  
+21     for start_cta in range(before_last, N, TILE_N):  
+22         offsets_cta = start_cta + tl.arange(0, TILE_N)  
+23         offsets_x = pid_m * N + offsets_cta  
+24         x = tl.load(input_ptr + offsets_x, mask=offsets_cta < N, other=-float("inf"))  
+25         new_m = tl.maximum(m, tl.max(x))  
+26         new_z = tl.exp(m - new_m) * z + tl.sum(tl.exp(x - new_m))  
+27         m = new_m  
+28         z = new_z  
+29  
+30     for start_cta in range(0, before_last, TILE_N):  
+31         offsets_cta = start_cta + tl.arange(0, TILE_N)  
+32         offsets_x = pid_m * N + offsets_cta  
+33         x = tl.load(input_ptr + offsets_x)  
+34         e = tl.exp(x - m)  
+35         out = e / z  
+36         tl.store(output_ptr + offsets_x, out)  
+37  
+38     for start_cta in range(before_last, N, TILE_N):  
+39         offsets_cta = start_cta + tl.arange(0, TILE_N)  
+40         offsets_x = pid_m * N + offsets_cta  
+41         x = tl.load(input_ptr + offsets_x, mask=offsets_cta < N, other=-float("inf"))  
+42         e = tl.exp(x - m)  
+43         out = e / z  
+44         tl.store(output_ptr + offsets_x, out, mask=offsets_cta < N)
```

算子开发示例: Softmax

```
● ● ● 4_online_softmax_v3.py

1  @triton.jit
2  def softmax_kernel_online_v1(output_ptr, input_ptr, M, N, TILE_N: tl.constexpr):
3      pid_m = tl.program_id(0)
4      m = tl.full(() , value=-float("inf") , dtype=output_ptr.dtype.element_ty)
5      z = tl.full(() , value=0 , dtype=output_ptr.dtype.element_ty)
6      before_last = prev_multiple_of(N, TILE_N)
7      for start_cta in range(0, before_last, TILE_N):
8          offsets_cta = start_cta + tl.arange(0, TILE_N)
9          offsets_x = pid_m * N + offsets_cta
10         x = tl.load(input_ptr + offsets_x)
11         # x = tl.load(input_ptr + offsets_x, mask=offsets_cta < N, other=-float("inf"))
12         new_m = tl.maximum(m, tl.max(x))
13         new_z = tl.exp(m - new_m) * z + tl.sum(tl.exp(x - new_m))
14         m = new_m
15         z = new_z
16
17     for start_cta in range(before_last, N, TILE_N):
18         offsets_cta = start_cta + tl.arange(0, TILE_N)
19         offsets_x = pid_m * N + offsets_cta
20         x = tl.load(input_ptr + offsets_x, mask=offsets_cta < N, other=-float("inf"))
21         new_m = tl.maximum(m, tl.max(x))
22         new_z = tl.exp(m - new_m) * z + tl.sum(tl.exp(x - new_m))
23         m = new_m
24         z = new_z
25
+ 26     for start_cta in range(0, TILE_N, TILE_N):
+ 27         offsets_cta = before_last + tl.arange(0, TILE_N)
28         offsets_x = pid_m * N + offsets_cta
29         x = tl.load(input_ptr + offsets_x, mask=offsets_cta < N, other=-float("inf"))
30         e = tl.exp(x - m)
31         out = e / z
32         tl.store(output_ptr + offsets_x, out, mask=offsets_cta < N)
33
+ 34     for start_cta in range(TILE_N, N, TILE_N):
+ 35         offsets_cta = (before_last - start_cta) + tl.arange(0, TILE_N)
36         offsets_x = pid_m * N + offsets_cta
37         x = tl.load(input_ptr + offsets_x)
38         e = tl.exp(x - m)
39         out = e / z
40         tl.store(output_ptr + offsets_x, out)
```

算子开发示例: Softmax



5_online_softmax_v4.py

```
1  for start_cta in range(0, TILE_N, TILE_N):
2      offsets_cta = before_last + tl.arange(0, TILE_N)
3      offsets_x = pid_m * N + offsets_cta
4      x = tl.load(
5          input_ptr + offsets_x,
+ 6          eviction_policy="evict_first",
7          mask=offsets_cta < N,
8          other=-float("inf"),
9      )
10     e = tl.exp(x - m)
11     out = e / z
12     tl.store(output_ptr + offsets_x, out, mask=offsets_cta < N)
13
14     for start_cta in range(TILE_N, N, TILE_N):
15         offsets_cta = (before_last - start_cta) + tl.arange(0, TILE_N)
16         offsets_x = pid_m * N + offsets_cta
+ 17         x = tl.load(input_ptr + offsets_x, eviction_policy="evict_first"
18     )     e = tl.exp(x - m)
19     out = e / z
20     tl.store(output_ptr + offsets_x, out)
```

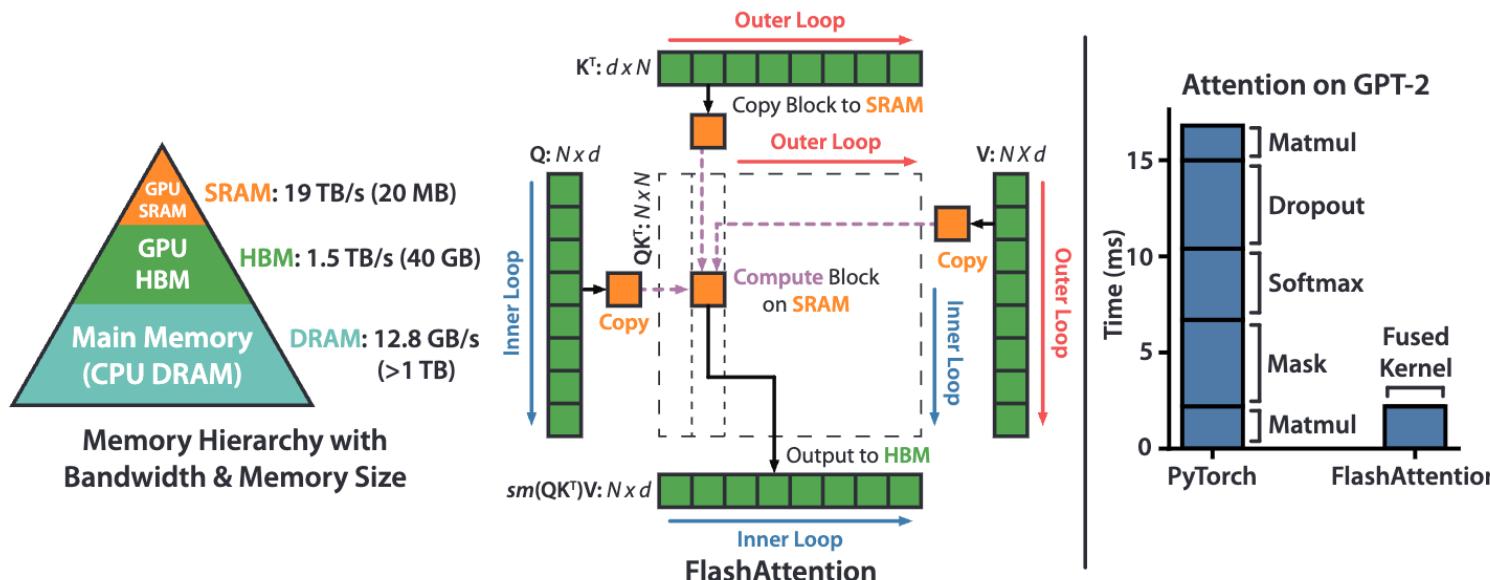
FlashAttention

FLASHATTENTION: Fast and Memory-Efficient Exact Attention with IO-Awareness

Tri Dao[†], Daniel Y. Fu[†], Stefano Ermon[†], Atri Rudra[‡], and Christopher Ré[†]

[†]Department of Computer Science, Stanford University

[‡]Department of Computer Science and Engineering, University at Buffalo, SUNY



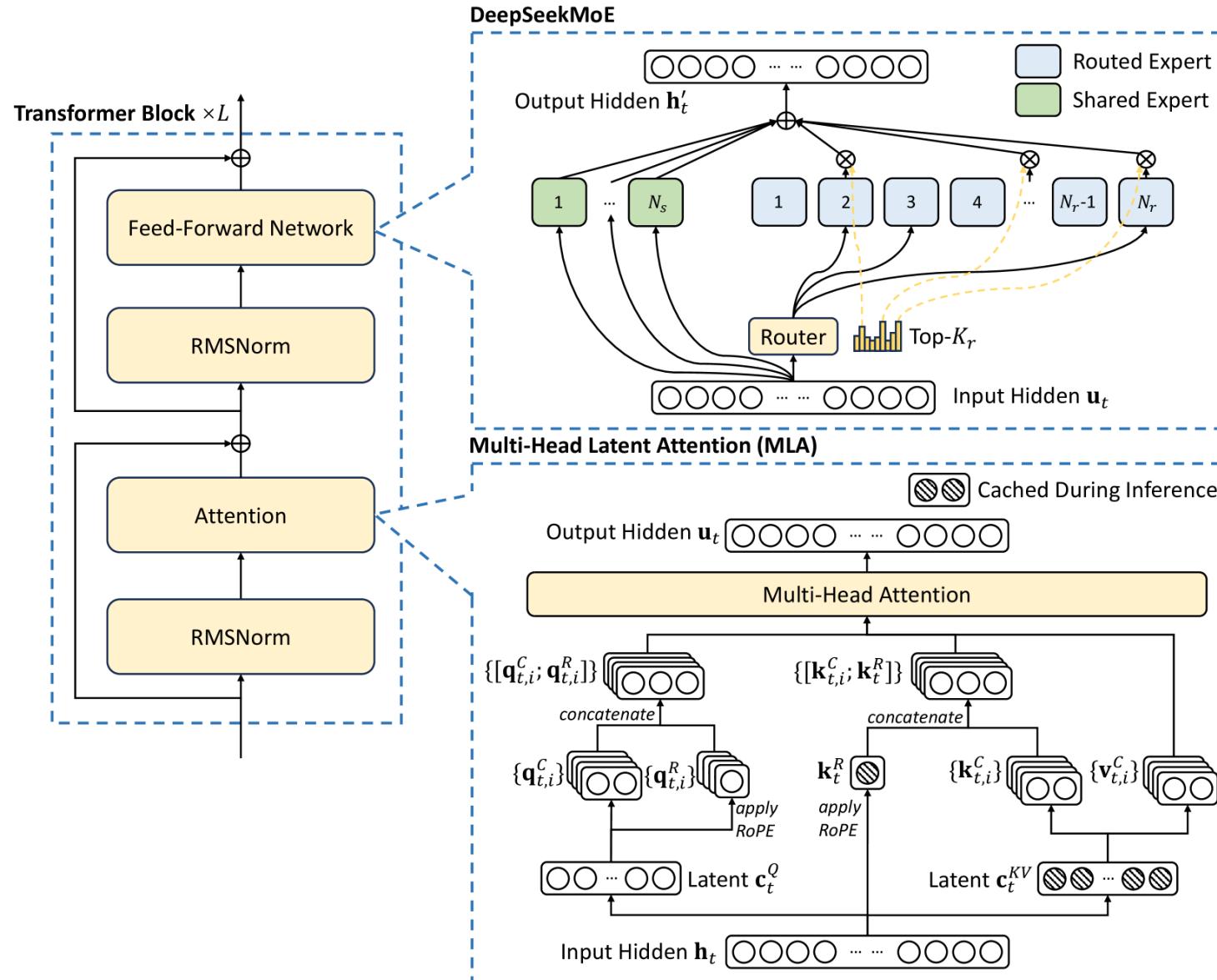
FlashAttention

Algorithm 1 FLASHATTENTION

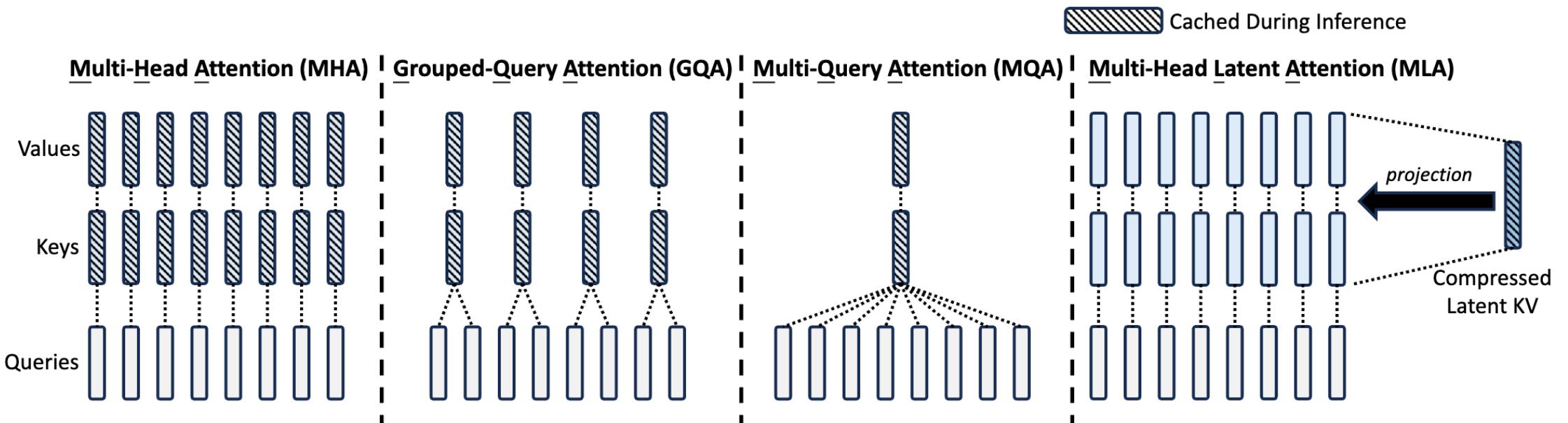
Require: Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ in HBM, on-chip SRAM of size M .

- 1: Set block sizes $B_c = \lceil \frac{M}{4d} \rceil, B_r = \min(\lceil \frac{M}{4d} \rceil, d)$.
 - 2: Initialize $\mathbf{O} = (0)_{N \times d} \in \mathbb{R}^{N \times d}, \ell = (0)_N \in \mathbb{R}^N, m = (-\infty)_N \in \mathbb{R}^N$ in HBM.
 - 3: Divide \mathbf{Q} into $T_r = \left\lceil \frac{N}{B_r} \right\rceil$ blocks $\mathbf{Q}_1, \dots, \mathbf{Q}_{T_r}$ of size $B_r \times d$ each, and divide \mathbf{K}, \mathbf{V} in to $T_c = \left\lceil \frac{N}{B_c} \right\rceil$ blocks $\mathbf{K}_1, \dots, \mathbf{K}_{T_c}$ and $\mathbf{V}_1, \dots, \mathbf{V}_{T_c}$, of size $B_c \times d$ each.
 - 4: Divide \mathbf{O} into T_r blocks $\mathbf{O}_i, \dots, \mathbf{O}_{T_r}$ of size $B_r \times d$ each, divide ℓ into T_r blocks $\ell_i, \dots, \ell_{T_r}$ of size B_r each, divide m into T_r blocks m_1, \dots, m_{T_r} of size B_r each.
 - 5: **for** $1 \leq j \leq T_c$ **do**
 - 6: Load $\mathbf{K}_j, \mathbf{V}_j$ from HBM to on-chip SRAM.
 - 7: **for** $1 \leq i \leq T_r$ **do**
 - 8: Load $\mathbf{Q}_i, \mathbf{O}_i, \ell_i, m_i$ from HBM to on-chip SRAM.
 - 9: On chip, compute $\mathbf{S}_{ij} = \mathbf{Q}_i \mathbf{K}_j^T \in \mathbb{R}^{B_r \times B_c}$.
 - 10: On chip, compute $\tilde{m}_{ij} = \text{rowmax}(\mathbf{S}_{ij}) \in \mathbb{R}^{B_r}, \tilde{\mathbf{P}}_{ij} = \exp(\mathbf{S}_{ij} - \tilde{m}_{ij}) \in \mathbb{R}^{B_r \times B_c}$ (pointwise), $\tilde{\ell}_{ij} = \text{rowsum}(\tilde{\mathbf{P}}_{ij}) \in \mathbb{R}^{B_r}$.
 - 11: On chip, compute $m_i^{\text{new}} = \max(m_i, \tilde{m}_{ij}) \in \mathbb{R}^{B_r}, \ell_i^{\text{new}} = e^{m_i - m_i^{\text{new}}} \ell_i + e^{\tilde{m}_{ij} - m_i^{\text{new}}} \tilde{\ell}_{ij} \in \mathbb{R}^{B_r}$.
 - 12: Write $\mathbf{O}_i \leftarrow \text{diag}(\ell_i^{\text{new}})^{-1} (\text{diag}(\ell_i) e^{m_i - m_i^{\text{new}}} \mathbf{O}_i + e^{\tilde{m}_{ij} - m_i^{\text{new}}} \tilde{\mathbf{P}}_{ij} \mathbf{V}_j)$ to HBM.
 - 13: Write $\ell_i \leftarrow \ell_i^{\text{new}}, m_i \leftarrow m_i^{\text{new}}$ to HBM.
 - 14: **end for**
 - 15: **end for**
 - 16: Return \mathbf{O} .
-

Multi-head Latent Attention

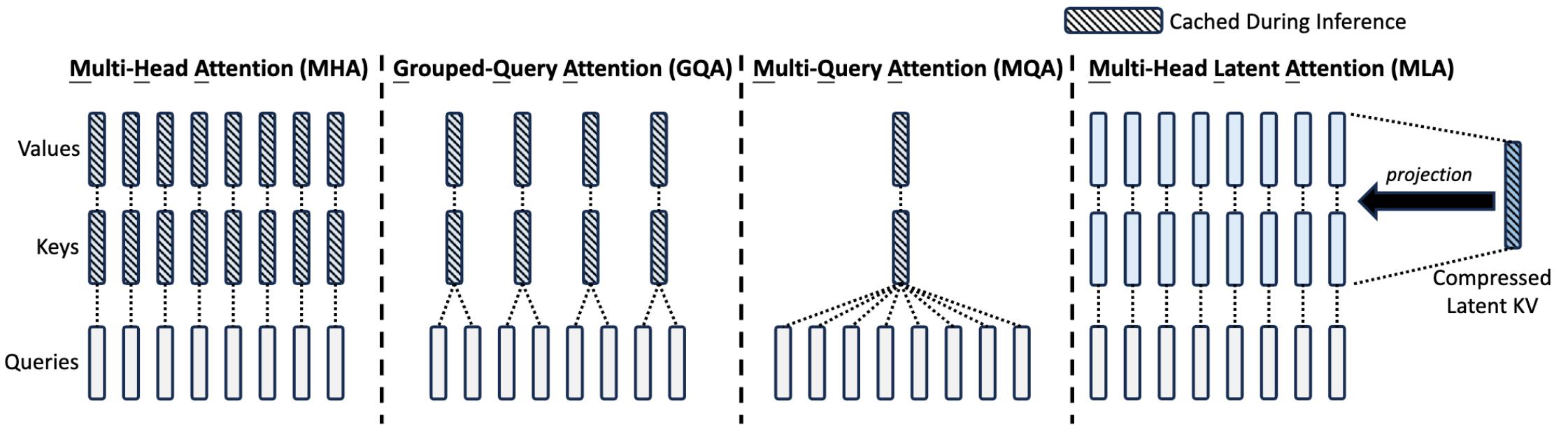


MHA vs. MLA



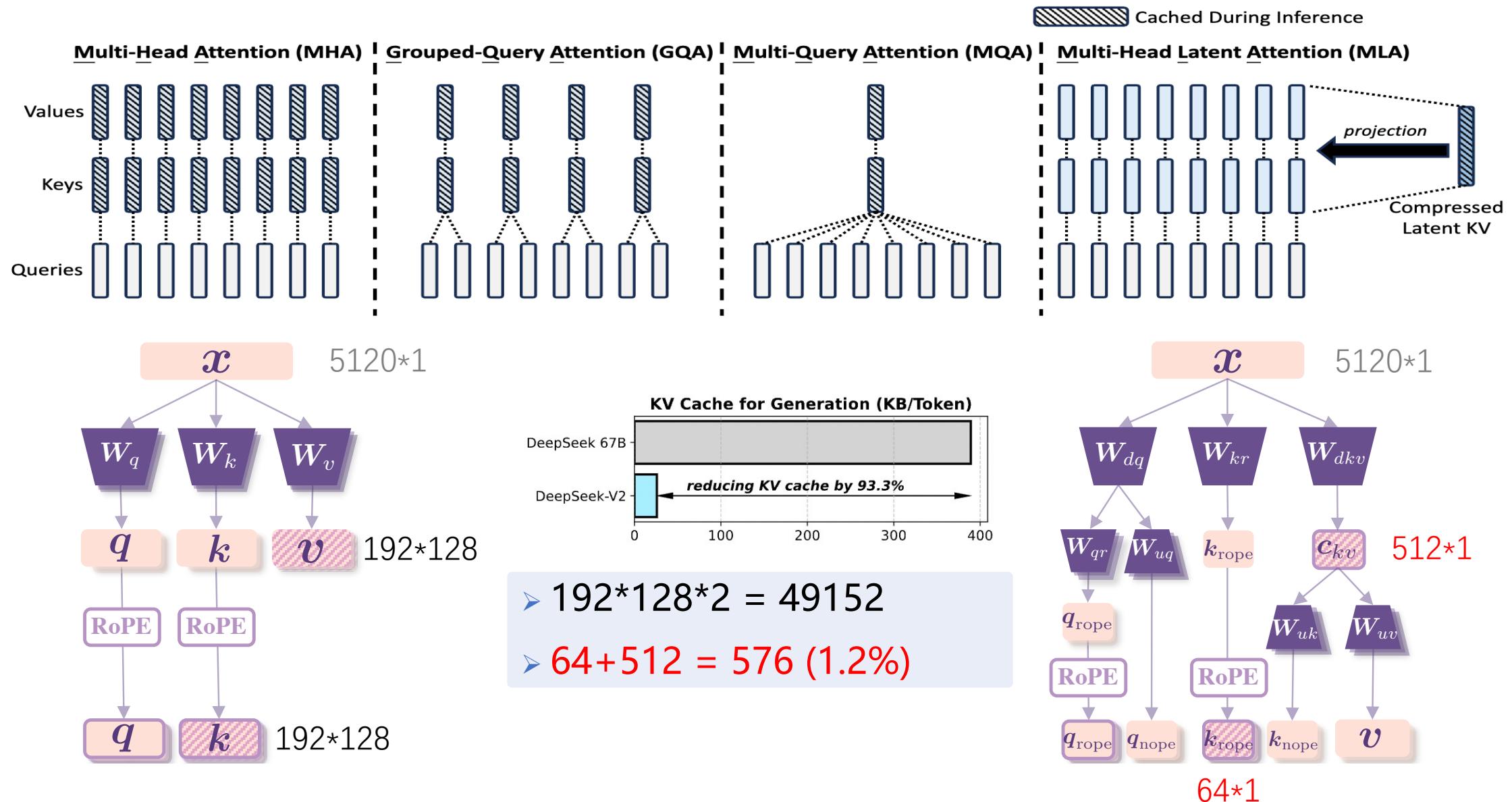
- 注意力存在**访存瓶颈** (KV缓存)
- GQA/MQA减少N倍缓存, 但**性能损失**严重
- DeepSeek提出**MLA**

MHA vs. MLA



- 注意力存在**访存瓶颈** (KV缓存)
- GQA/MQA减少N倍缓存, 但**性能损失**严重
- DeepSeek提出**MLA**

MHA vs. MLA



MHA vs. MLA

Benchmark (Metric)	# Shots	Dense 7B w/ MQA	Dense 7B w/ GQA (8 Groups)	Dense 7B w/ MHA
# Params	-	7.1B	6.9B	6.9B
BBH (EM)	3-shot	33.2	35.6	37.0
MMLU (Acc.)	5-shot	37.9	41.2	45.2
C-Eval (Acc.)	5-shot	30.0	37.7	42.9
CMMLU (Acc.)	5-shot	34.6	38.4	43.5
		-8.2	-3.9	

➤ GQA/MQA即使参数量对齐，**性能损失**仍然严重

MHA vs. MLA

Benchmark (Metric)	# Shots	Small MoE w/ MHA	Small MoE w/ MLA	Large MoE w/ MHA	Large MoE w/ MLA
# Activated Params	-	2.5B	2.4B	25.0B	21.5B
# Total Params	-	15.8B	15.7B	250.8B	247.4B
KV Cache per Token (# Element)	-	110.6K	15.6K	860.2K	34.6K
BBH (EM)	3-shot	37.9	39.0	46.6	50.7
MMLU (Acc.)	5-shot	48.7	50.0	57.5	59.0
C-Eval (Acc.)	5-shot	51.6	50.9	57.9	59.2
CMMLU (Acc.)	5-shot	52.3	53.4	60.7	62.5

+0.7

+2.2

➤ MLA减少90%以上缓存，**性能持平或更优**

MHA vs. MLA

$$\begin{aligned}\mathbf{q}_i^{(h)} &= \mathbf{x}_i \mathbf{W}_q^{(\bar{h})} \\ \mathbf{k}_i^{(h)} &= \mathbf{x}_i \mathbf{W}_k^{(h)} \\ \mathbf{v}_i^{(h)} &= \mathbf{x}_i \mathbf{W}_v^{(h)} \\ \mathbf{q}_{i,\text{rope}}^{(h)} &= \text{RoPE}(\mathbf{q}_i^{(h)}) \\ \mathbf{k}_{i,\text{rope}}^{(h)} &= \text{RoPE}(\mathbf{k}_i^{(h)})\end{aligned}$$

$$\begin{aligned}\mathbf{o}_i^{(h)} &= \text{Softmax} \left(\mathbf{q}_{i,\text{rope}}^{(h)} \mathbf{k}_{\leq i, \text{rope}}^{(h)\top} \right) \mathbf{v}_{\leq i}^{(h)}, \\ \text{MHA}(\mathbf{x}_i) &= \left[\mathbf{o}_i^{(1)}, \dots, \mathbf{o}_i^{(n_h)} \right] \mathbf{W}_o,\end{aligned}$$

$$\begin{aligned}\mathbf{q}_{i,\text{rope}}^{(h)}, \mathbf{k}_{i,\text{rope}} &= \text{RoPE} \left(\mathbf{x}_i \mathbf{W}_{dq} \mathbf{W}_{qr}^{(h)}, \mathbf{x}_i \mathbf{W}_{kr} \right), \\ \mathbf{q}_{i,\text{nope}}^{(h)} &= \mathbf{x}_i \mathbf{W}_{dq} \mathbf{W}_{qc}^{(h)}, \\ \mathbf{c}_{i,kv} &= \mathbf{x}_i \mathbf{W}_{dkv}, \\ \mathbf{k}_{i,\text{nope}}^{(h)}, \mathbf{v}_i^{(h)} &= \mathbf{c}_{i,kv} \mathbf{W}_{uk}^{(h)}, \mathbf{c}_{i,kv} \mathbf{W}_{uv}^{(h)}, \\ \mathbf{o}_i^{(h)} &= \text{Softmax} \left(\mathbf{q}_{i,\text{rope}}^{(h)} \mathbf{k}_{\leq i, \text{rope}}^{(h)\top} + \mathbf{q}_{i,\text{nope}} \mathbf{k}_{\leq i, \text{nope}}^{(h)\top} \right) \\ &\quad \cdot \mathbf{v}_{\leq i}^{(h)} \\ \text{MLA}(\mathbf{x}_i) &= \left[\mathbf{o}_i^{(1)}, \dots, \mathbf{o}_i^{(n_h)} \right] \cdot \mathbf{W}_o.\end{aligned}\tag{3}$$

MHA vs. MLA

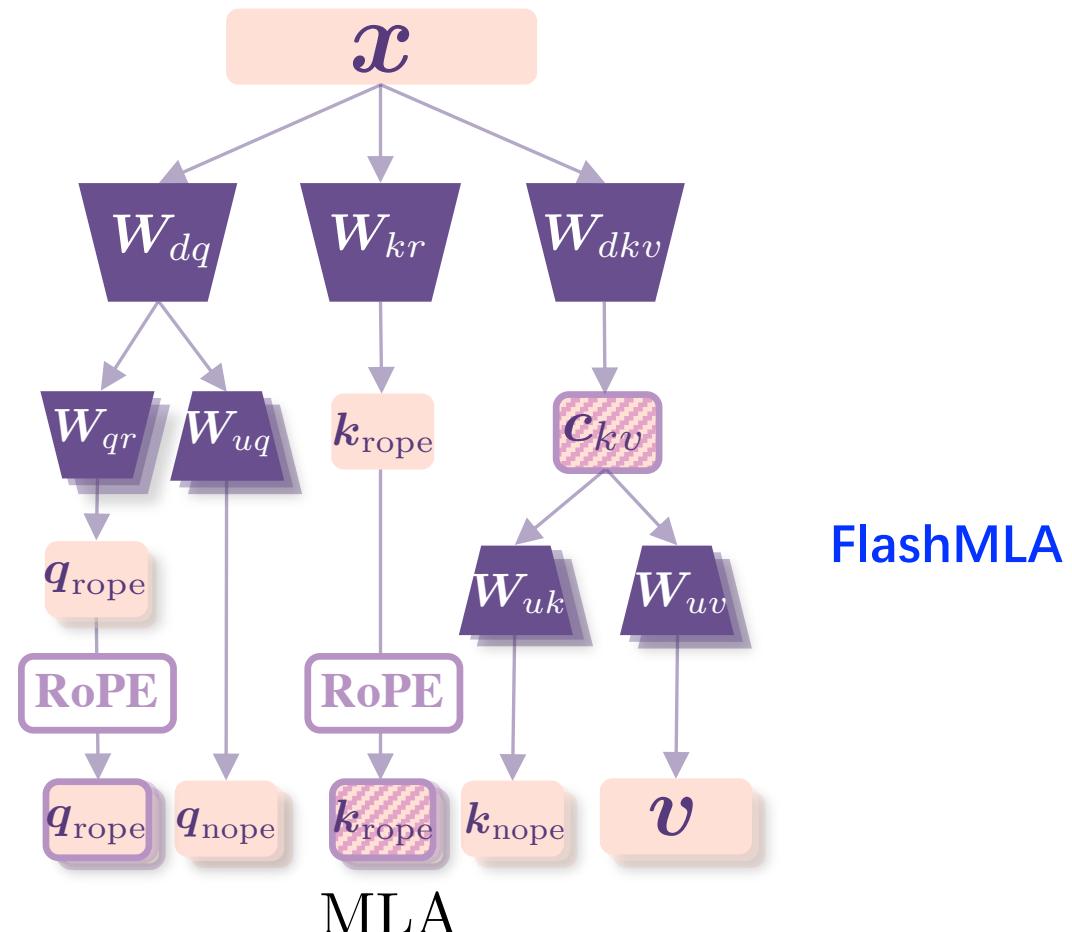
位置无关部分：矩阵吸收

$$\begin{aligned}\mathbf{q}_{i,\text{nope}} \mathbf{k}_{j,\text{nope}}^\top &= (\mathbf{x}_i \mathbf{W}_{dq} \mathbf{W}_{qc}) (\mathbf{c}_{j,kv} \mathbf{W}_{uk})^\top \\ &= \mathbf{x}_i \left(\mathbf{W}_{dq} \mathbf{W}_{qc} \mathbf{W}_{uk}^\top \right) \mathbf{c}_{j,kv}^\top\end{aligned}$$

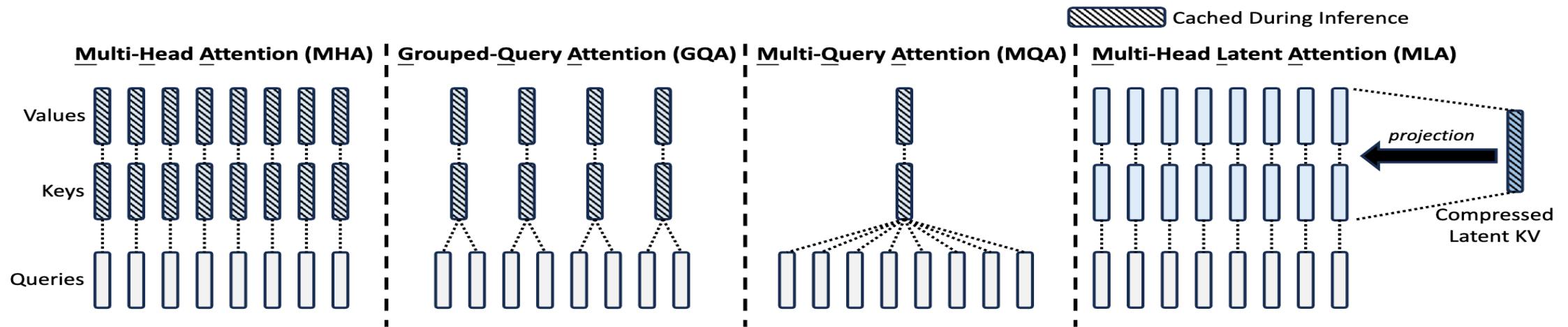
位置相关部分：正常计算

$$\begin{aligned}\mathbf{q}_{i,\text{rope}} \mathbf{k}_{j,\text{rope}}^\top &= (\mathbf{x}_i \mathbf{W}_{dq} \mathbf{W}_{qr} \mathbf{R}_i) (\mathbf{x}_j \mathbf{W}_{kr} \mathbf{R}_j)^\top \\ &= \mathbf{x}_i \left(\mathbf{W}_{dq} \mathbf{W}_{qc} \boxed{\mathbf{R}_{j-i}} \mathbf{W}_{kr}^\top \right) \mathbf{x}_j^\top\end{aligned}$$

无法吸收



MHA vs. MLA



192

$$\mathbf{o}_i^{(h)} = \text{Softmax} \left(\mathbf{q}_{i,\text{rope}}^{(h)} \mathbf{k}_{\leq i, \text{rope}}^{(h)\top} \right) \mathbf{v}_{\leq i}^{(h)},$$

$$\text{MHA}(\mathbf{x}_i) = \left[\mathbf{o}_i^{(1)}, \dots, \mathbf{o}_i^{(n_h)} \right] \mathbf{W}_o,$$

64

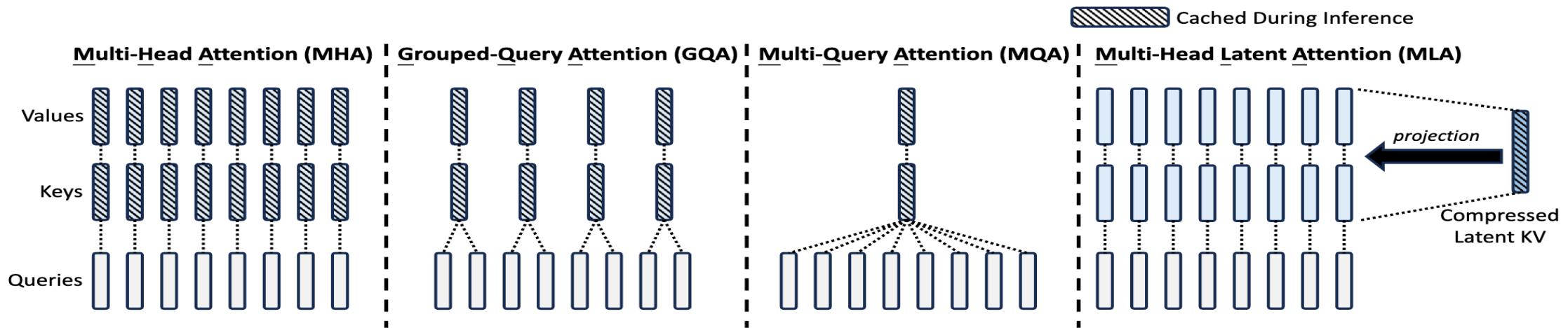
$$\begin{aligned} \mathbf{o}_i^{(h)} &= \text{Softmax} \left(\mathbf{q}_{i,\text{rope}}^{(h)} \mathbf{k}_{\leq i, \text{rope}}^{(h)\top} + \mathbf{q}_{i,\text{nope}} \mathbf{k}_{\leq i, \text{nope}}^{(h)\top} \right) \\ &\cdot \mathbf{v}_{\leq i}^{(h)} \end{aligned}$$

128

$$\text{MLA}(\mathbf{x}_i) = \left[\mathbf{o}_i^{(1)}, \dots, \mathbf{o}_i^{(n_h)} \right] \cdot \mathbf{W}_o. \quad (3)$$

$$\begin{aligned} \mathbf{q}_{i,\text{nope}} \mathbf{k}_{j,\text{nope}}^\top &= (\mathbf{x}_i \mathbf{W}_{dq} \mathbf{W}_{qc}) (\mathbf{c}_{j,kv} \mathbf{W}_{uk})^\top \\ \text{矩阵吸收} &= \mathbf{x}_i \left(\mathbf{W}_{dq} \mathbf{W}_{qc} \mathbf{W}_{uk}^\top \right) \mathbf{c}_{j,kv}^\top \end{aligned}$$

MHA vs. MLA



$$\mathbf{o}_i^{(h)} = \text{Softmax} \left(\mathbf{q}_{i,\text{rope}}^{(h)} \mathbf{k}_{\leq i, \text{rope}}^{(h)\top} \right) \mathbf{v}_{\leq i}^{(h)},$$

$$\text{MHA}(\mathbf{x}_i) = \left[\mathbf{o}_i^{(1)}, \dots, \mathbf{o}_i^{(n_h)} \right] \mathbf{W}_o,$$

$$\mathbf{o}_i^{(h)} = \text{Softmax} \left(\mathbf{q}_{i,\text{rope}}^{(h)} \mathbf{k}_{\leq i, \text{rope}}^{(h)\top} + \mathbf{q}_{i,\text{nope}} \mathbf{k}_{\leq i, \text{nope}}^{(h)\top} \right) \cdot \mathbf{v}_{\leq i}^{(h)}$$

$$\text{MLA}(\mathbf{x}_i) = \left[\mathbf{o}_i^{(1)}, \dots, \mathbf{o}_i^{(n_h)} \right] \cdot \mathbf{W}_o. \quad (3)$$

➤ Q[B, S, 128, 192] K[B, S, 128, 192]

➤ Qr[B, 1, 128, 64] Kr[B, S, 128, 64]

➤ Qc[B, 1, 128, 128] Kc[B, S, 128, 128]
 ➤ Qc[B, 1, 128, 512] Kc[B, S, 1, 512]

Poster

Q&A