

# 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](mailto:taoji@fudan.edu.cn), [jfxiong24@m.fudan.edu.cn](mailto:jfxiong24@m.fudan.edu.cn), [binguo@stu.ecnu.edu.cn](mailto:binguo@stu.ecnu.edu.cn)

# 自我介绍



<https://taoji.eth.limo>



熊静飞 (研一)  
复旦大学



郭彬 (研一)  
华东师范大学

- ▶ Efficient infrastructure via algorithm-hardware co-design
- ▶ Language-centric multimodal in-context learning

时间	单位	经历	导师
2023.09–2025.09	复旦大学	博士后	黄萱菁&邱锡鹏
2017.09–2023.06	华东师范大学 (硕博连读)	博士	吴苑斌&王晓玲
2013.09–2017.06	华东师范大学 (信息学保送)	学士	吴苑斌

## Towards Economical Inference: Enabling DeepSeek's Multi-Head Latent Attention in Any Transformer-based LLMs

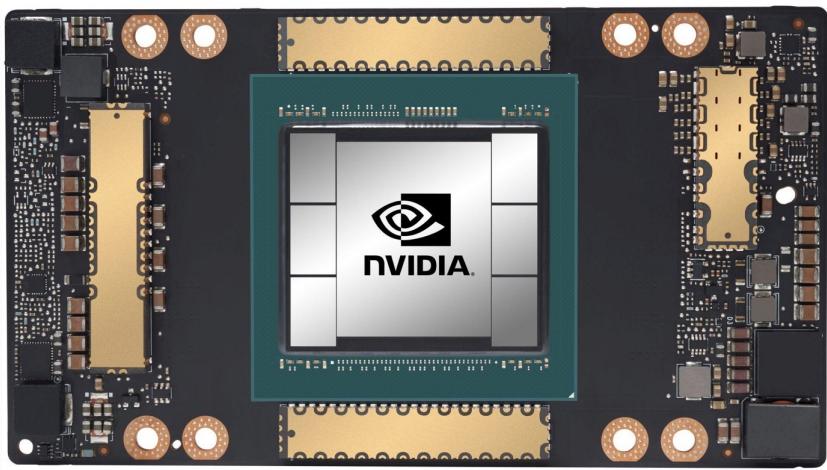
Tao Ji<sup>♦</sup>, Bin Guo<sup>♡</sup>, Yuanbin Wu<sup>♡</sup>,  
Qipeng Guo<sup>◊</sup>, Lixing Shen<sup>♦</sup>, Zhan Chen<sup>♦</sup>, Xipeng Qiu<sup>♦</sup>, Qi Zhang<sup>♦</sup>, Tao Gui<sup>♣✉</sup>  
♦Fudan University ♡East China Normal University ♪Hikvision Inc ◊Shanghai AI Lab  
{taoji, tgui}@fudan.edu.cn {binguo@stu, ybwu@cs}.ecnu.edu.cn

MHA2MLA Public

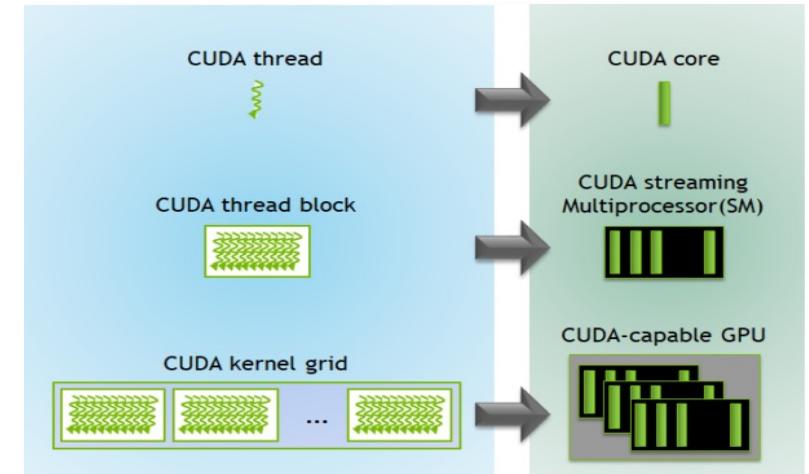
Towards Economical Inference: Enabling DeepSeek's Multi-Head Latent Attention in Any Transformer-based LLMs

● Python ⭐ 178 ⚡ 21

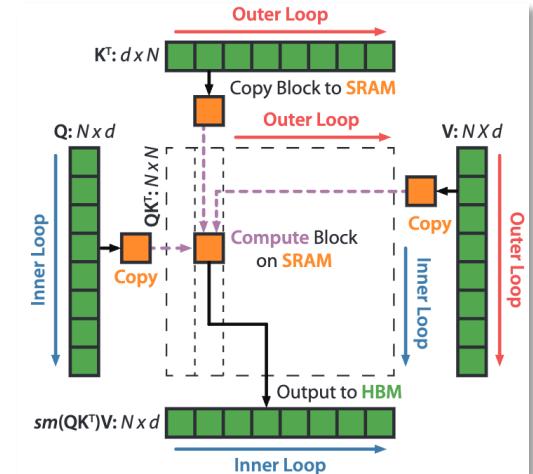
# AI-Infra Seminar



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>



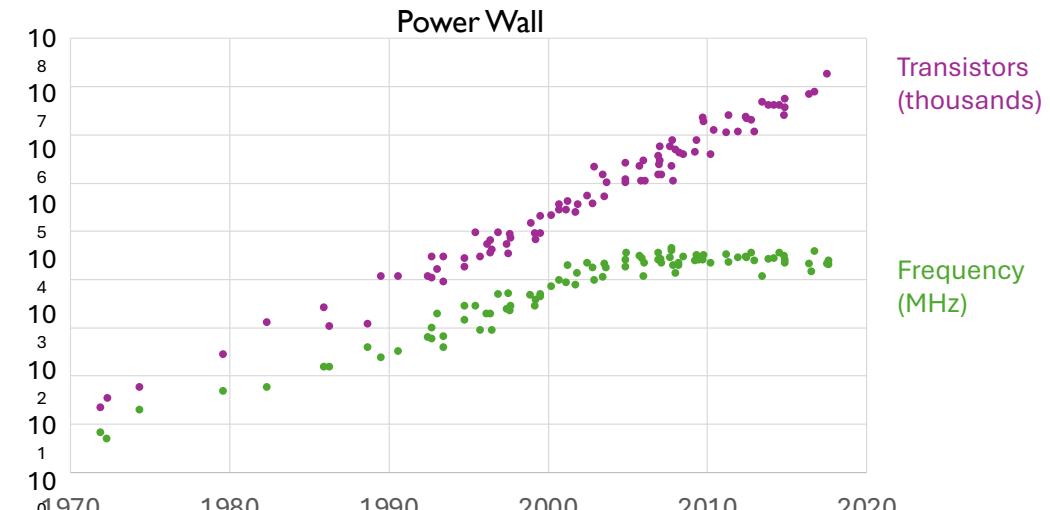
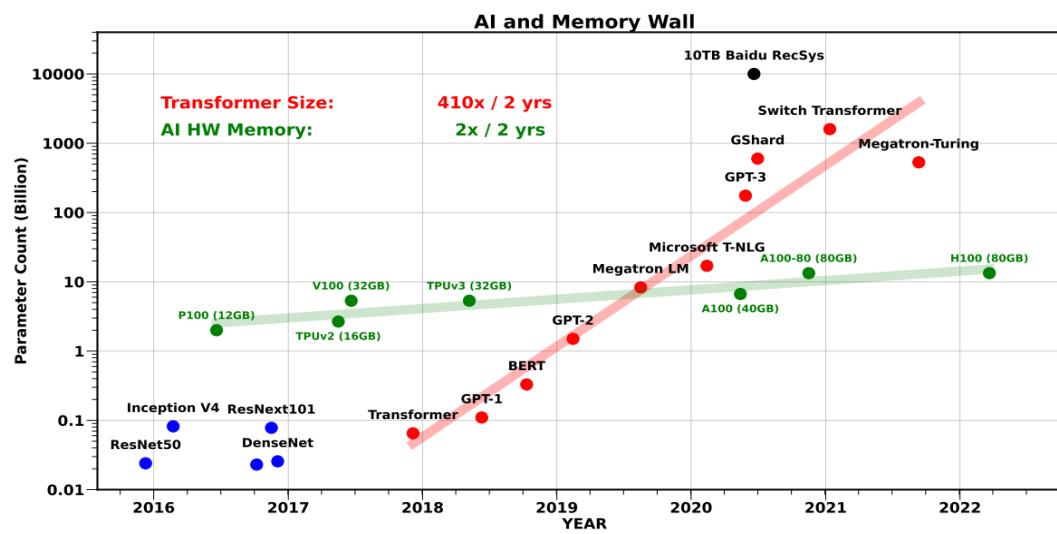
群聊: 并行计算&CUDA讨论班



# AI-Infrastructure

All hardware and software infrastructure used for model training and inference

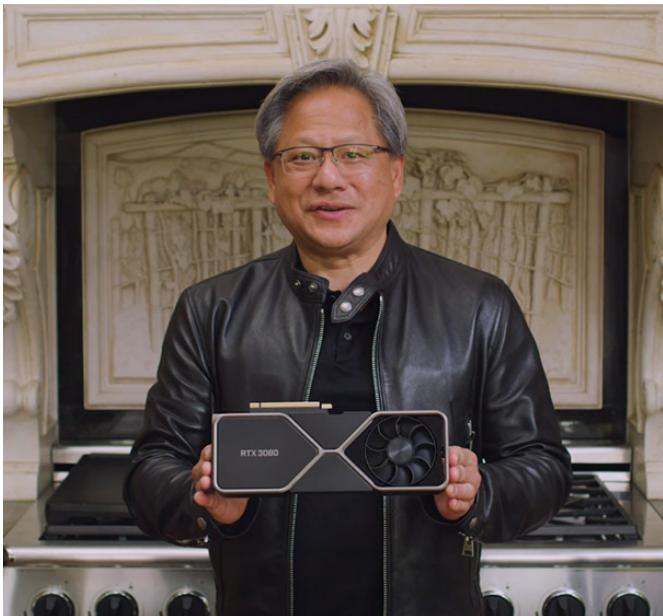
- ▶ model scale ↑↑, GPU scale ↑
- ▶ Transformer hardware-friendly  $\Rightarrow$  computation-bottleneck



Source: M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, C. Batten (1970-2010). K. Rupp (2010-2017).

# AI-Infrastructure

More GPU



More efficient infrastructure



# AI-Infrastructure

## DeepSeek开源周总结

### 第 1 天

FlashMLA：针对 NVIDIA Hopper GPU 的高性能解码内核

### 第 2 天

DeepEP：专为 MoE 模型设计的通信库

### 第 4 天

DualPipe & EPLB：分布式训练优化

### 第 3 天

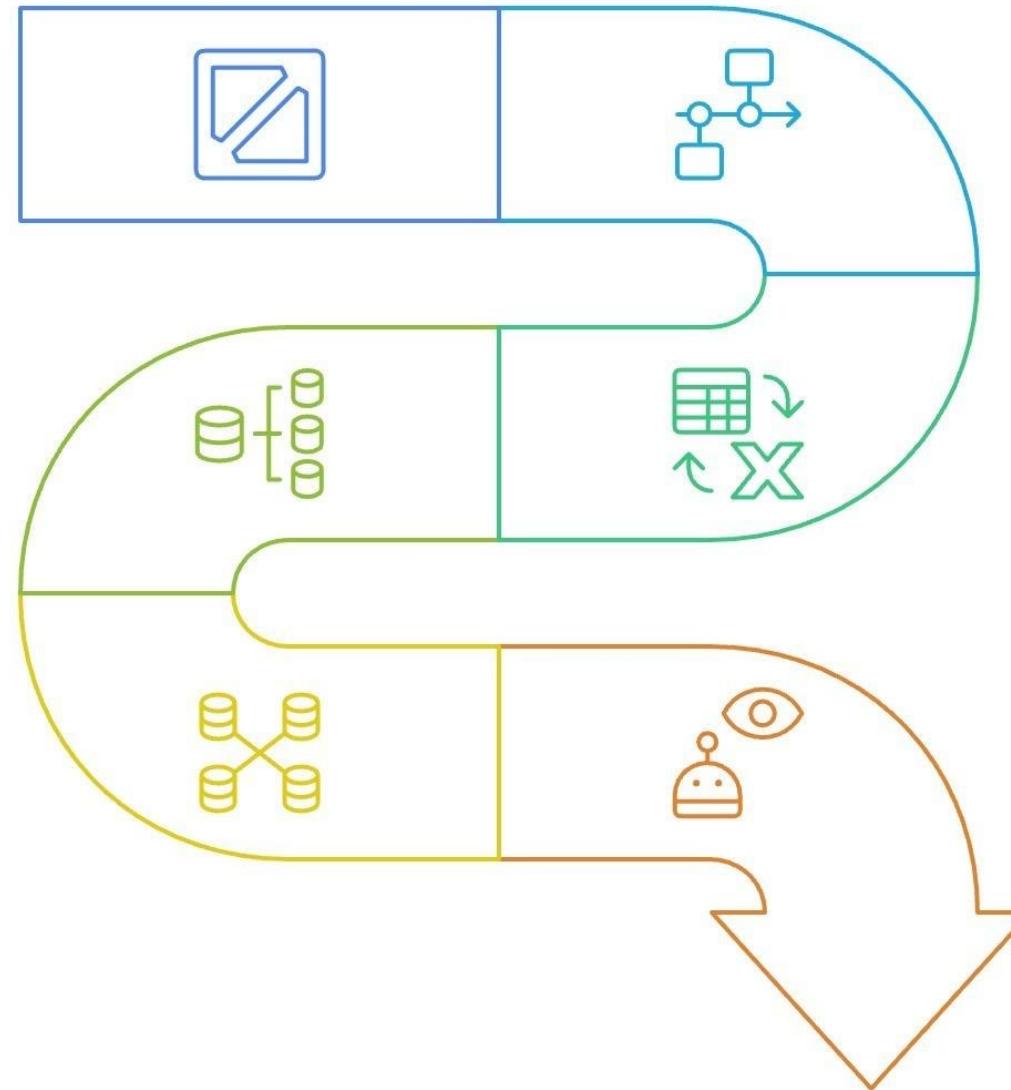
DeepGEMM：FP8 矩阵乘法库

### 第 5 天

3FS：高性能分布式文件系统

### 第 6 天

DeepSeek-V3/R1：AI 推理系统概览



# AI-Infra Conferences



International Conference on Architectural  
Support for Programming Languages and  
Operating Systems  
(ASPLOS)

<https://www.asplos-conference.org>



<https://lmsys.org>

## PPoPP 2025

Sat 1 - Wed 5 March 2025 Las Vegas, Nevada, United States

Attending ▾ Program ▾ Tracks ▾ Organization ▾ PPoPP Mailing List  Search

ACM SIGPLAN Symposium on  
Principles and Practice of Parallel  
Programming 2025

<https://ppopp25.sigplan.org>

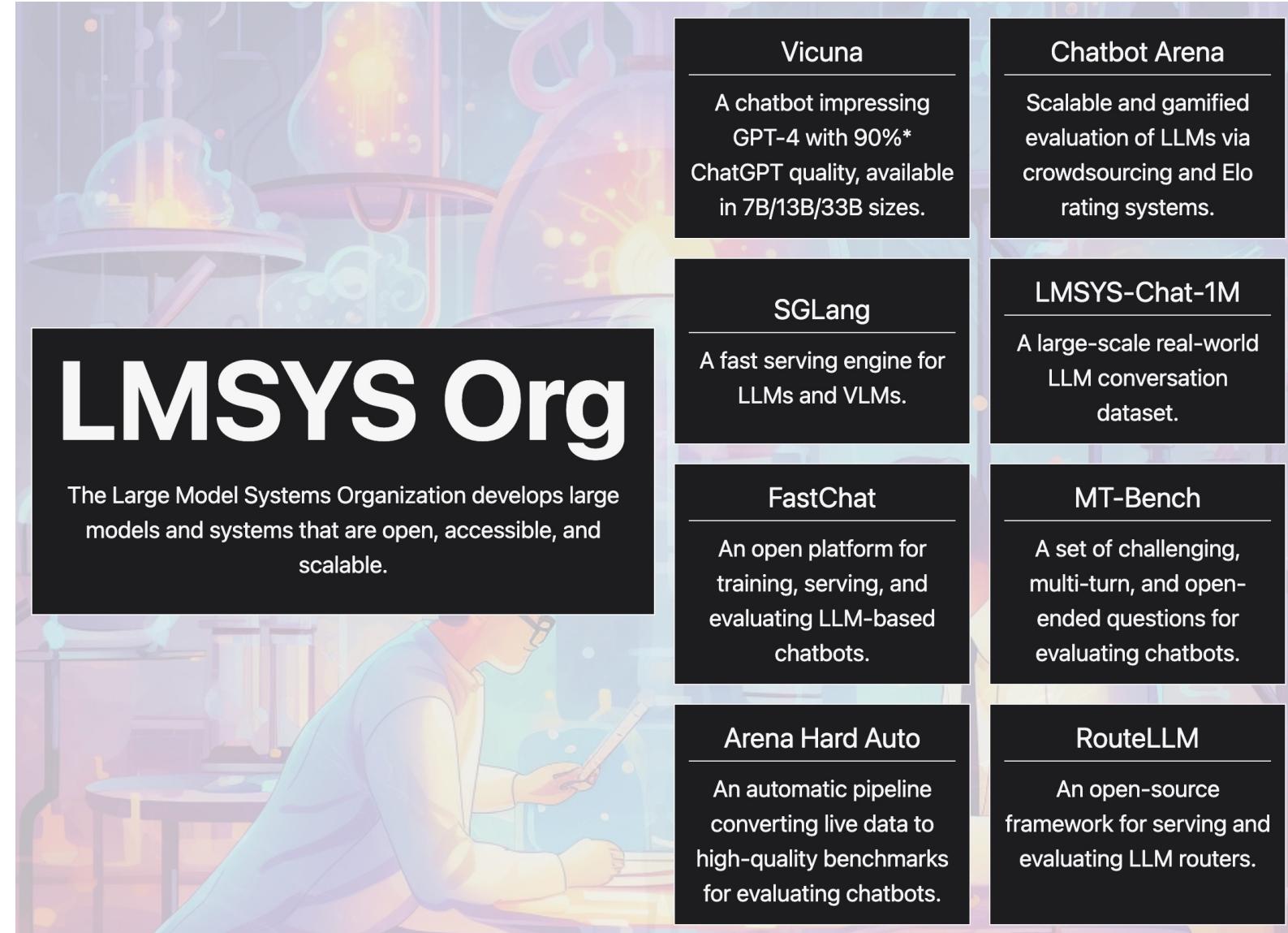


<https://supercomputing.org>

# AI-Infra Research Team

Multi-university in 2023  
UC Berkeley, Stanford,  
UCSD, CMU, MBZUAI

.....



**LMSYS Org**

The Large Model Systems Organization develops large models and systems that are open, accessible, and scalable.

Vicuna	Chatbot Arena
A chatbot impressing GPT-4 with 90%* ChatGPT quality, available in 7B/13B/33B sizes.	Scalable and gamified evaluation of LLMs via crowdsourcing and Elo rating systems.
SGLang	LMSYS-Chat-1M
A fast serving engine for LLMs and VLMs.	A large-scale real-world LLM conversation dataset.
FastChat	MT-Bench
An open platform for training, serving, and evaluating LLM-based chatbots.	A set of challenging, multi-turn, and open-ended questions for evaluating chatbots.
Arena Hard Auto	RouteLLM
An automatic pipeline converting live data to high-quality benchmarks for evaluating chatbots.	An open-source framework for serving and evaluating LLM routers.

# AI-Infra Research Team

[Home](#)[People](#)[Research](#)[Publications](#)

**Catalyst** is an interdisciplinary machine learning and systems research group exploring problems to automate learning systems. Our research spans multiple layers of the machine learning and system stack. Our group is a collaboration between researchers from the [Machine Learning Department](#), [Computer Science Department](#) and [Electrical & Computer Engineering Department](#) at the [Carnegie Mellon University](#).

## Research

### XGrammar

Efficient, Flexible and Portable LLM Structured Generation

[Read more »](#)

### FlexFlow Serve

Low-Latency, High-Performance LLM Serving

[Read more »](#)

### Hyperband and ASHA

Principled early-stopping approaches for hyperparameter optimization

[Read more »](#)

### TidalDecode

A sparse attention framework for large language model decoding

[Read more »](#)

### Cortex

End-to-end compilation of ML applications with dynamic and irregular control flow and data structure accesses

[Read more »](#)

### FlexFlow

Automatically Discovering Fast Parallelization Strategies for DNN Training

[Read more »](#)

### Machine Learning Compilation for Large Language Models

### Apache TVM Stack

TVM: An Automated End-to-End

### TASO

The Tensor Algebra SuperOptimizer for Deep Learning

## Faculty



**Beidi Chen**  
Assistant Professor



**Tianqi Chen**  
Assistant Professor



**Tim Dettmers**  
Assistant Professor



**Greg Ganger**  
Professor



**Phillip B. Gibbons**  
Professor



**Zhihao Jia**  
Assistant Professor



**Todd C. Mowry**  
Professor



**Ameet Talwalkar**  
Assistant Professor



**Eric Xing**  
Professor

# AI-Infra Research Team

UC Berkeley Sky Computing Lab

People Projects Publications News Events Sponsors Contact DARE

## Sky Computing Towards Utility Computing for the Cloud

The Sky Above the Clouds



Read the White Paper

Video: Sky Computing



Watch

DARE: Diversifying Access to Research in Engineering



Learn More

### News

May 21, 2025

Red Hat Launches the IIm-d Community, Powering Distributed Gen AI Inference at Scale

"We are pleased to see Red Hat build upon the established success of vLLM, which originated

### Events

May 16, 2025

Dissertation Talk: Programming Models for Correct and Modular Distributed Systems – Shadaj Laddad

Distributed systems are a fundamental part of modern computing, but they are notoriously

### Publications

April 2025

Smart Casual Verification of the Confidential Consortium Framework.

April 2025

SuperServe: Fine-Grained Inference Serving for Unpredictable Workloads.

### Core Faculty



Natacha Crooks

Core Faculty



Joseph Gonzalez

Core Faculty



Raluca Ada Popa

Core Faculty



Koushik Sen

Core Faculty



Ion Stoica

Core Faculty



Matei Zaharia

Core Faculty

<https://sky.cs.berkeley.edu>

## Tri Dao

Assistant Professor of Computer Science at [Princeton University](#).  
Chief Scientist at [Together AI](#).

[CV](#) (updated 06/2025)

Previously: PhD, [Department of Computer Science, Stanford University](#)

### Research Interests

Machine learning and systems, with a focus on efficient training and inference:

- Hardware-aware algorithms.
- Sequence models with long-range memory.

### Current PhD Students

- [Ted Zadouri](#)
- [Berlin Chen](#)
- [Wentao Guo](#)



tri [at] tridao (dot) me

# AI-Infra Research Team

**IPADS** | 并行与分布式系统研究所



教师们



# AI-Infra Research Team

MADSys [Home](#) News People Projects Publications



## MADSys

The MADSys (Machine Learning, AI, Big Data Systems) group is dedicated to the design, implementation, evaluation, and application of parallel and distributed systems. Our research spans various methods aimed at accelerating data processing. Although the foundational principles like caching, batching, and overlapping are consistent, the strategic and innovative application of these techniques allows system researchers to optimally utilize diverse hardware resources across different scenarios.

## Faculty



[Yongwei Wu](#)



[Kang Chen](#)



[Jinlei Jiang](#)

Jing Zheng



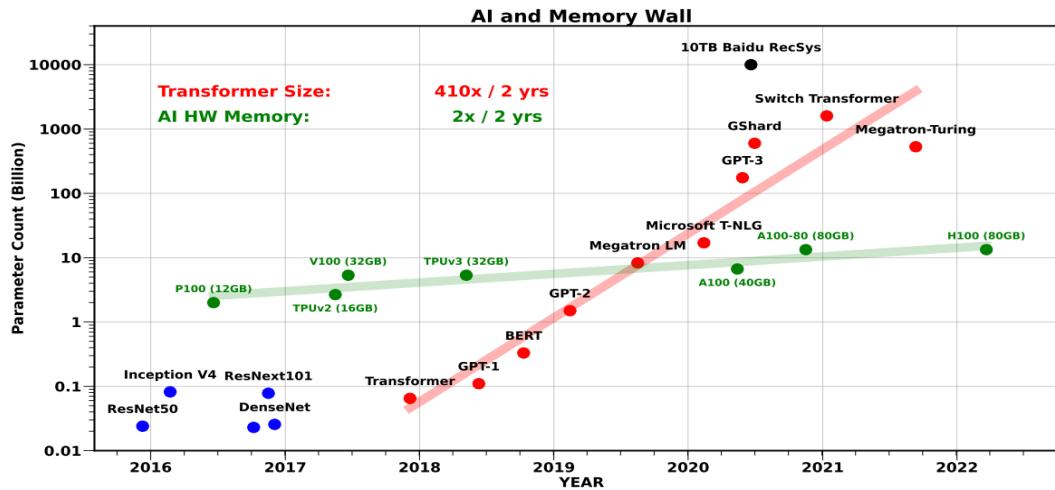
[Mingxing Zhang](#)

Lizhi Deng

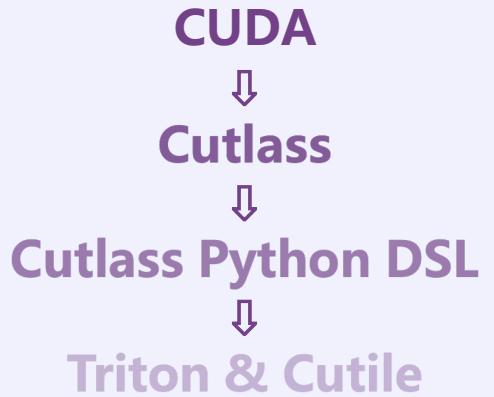


[Yingdi Shan](#)

# AI-Infra Research Team



► more important



► easier

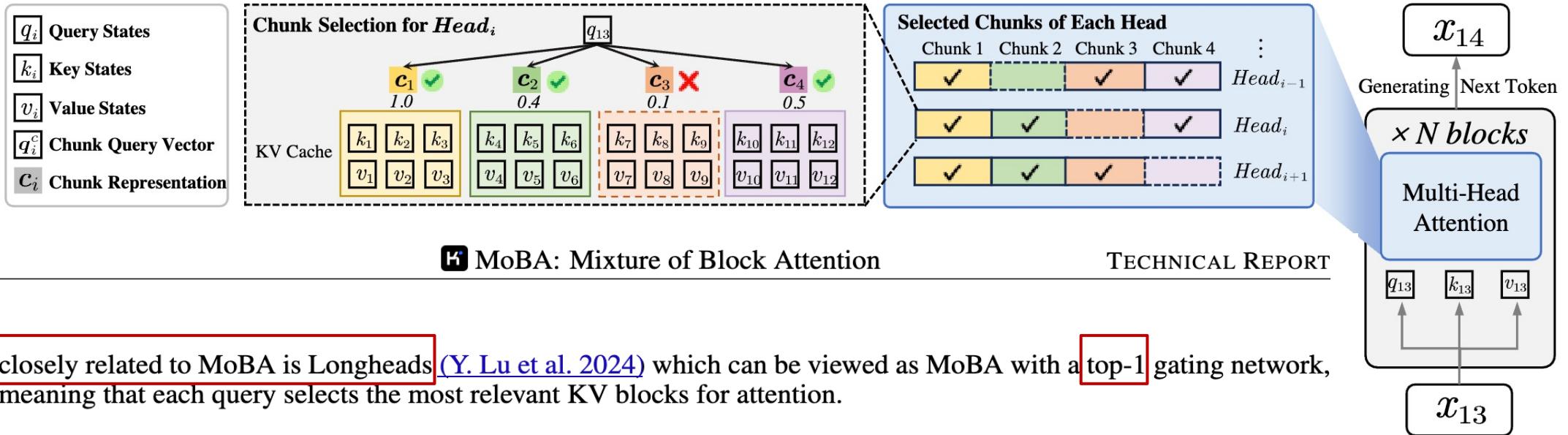


# Why Am I Interested?

## LONGHEADS: Multi-Head Attention is Secretly a Long Context Processor

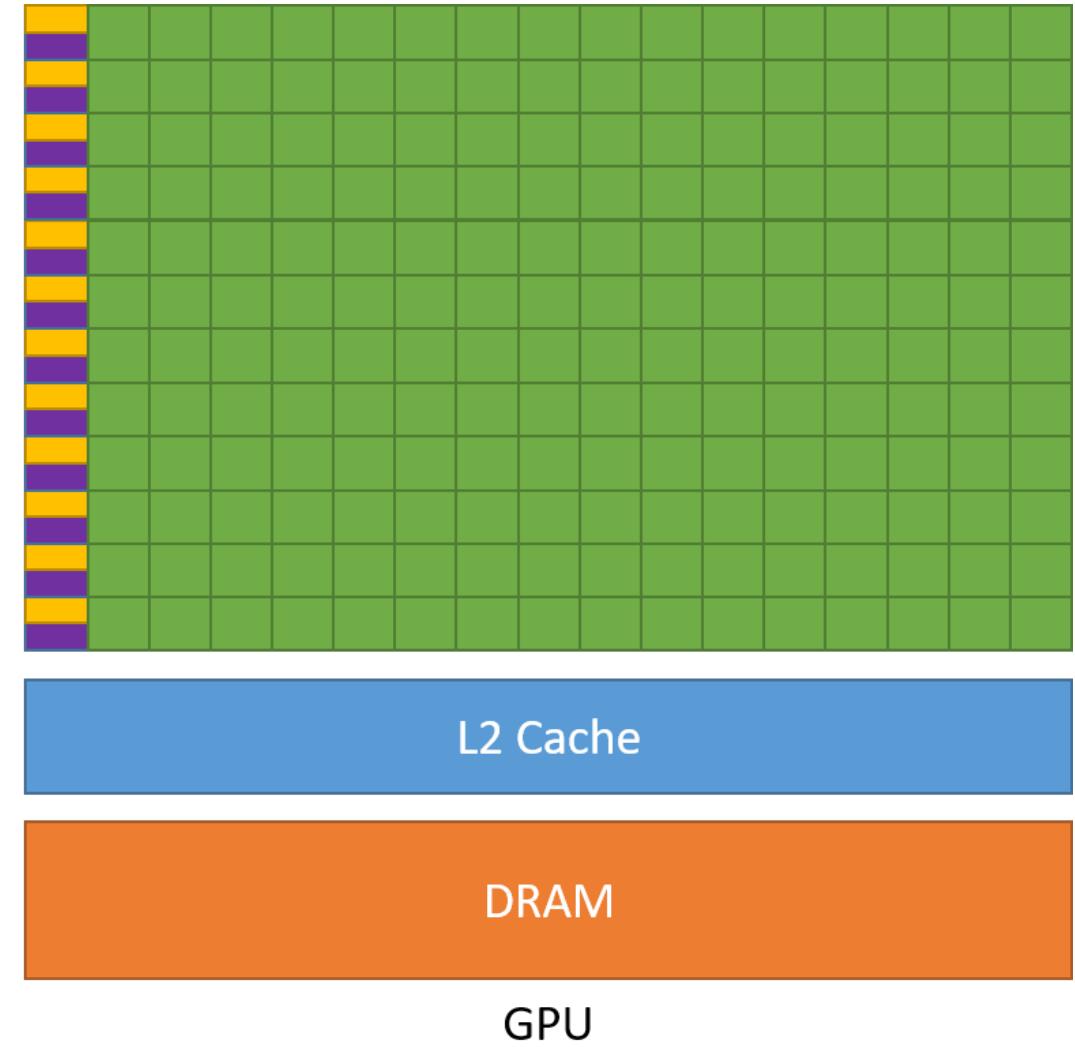
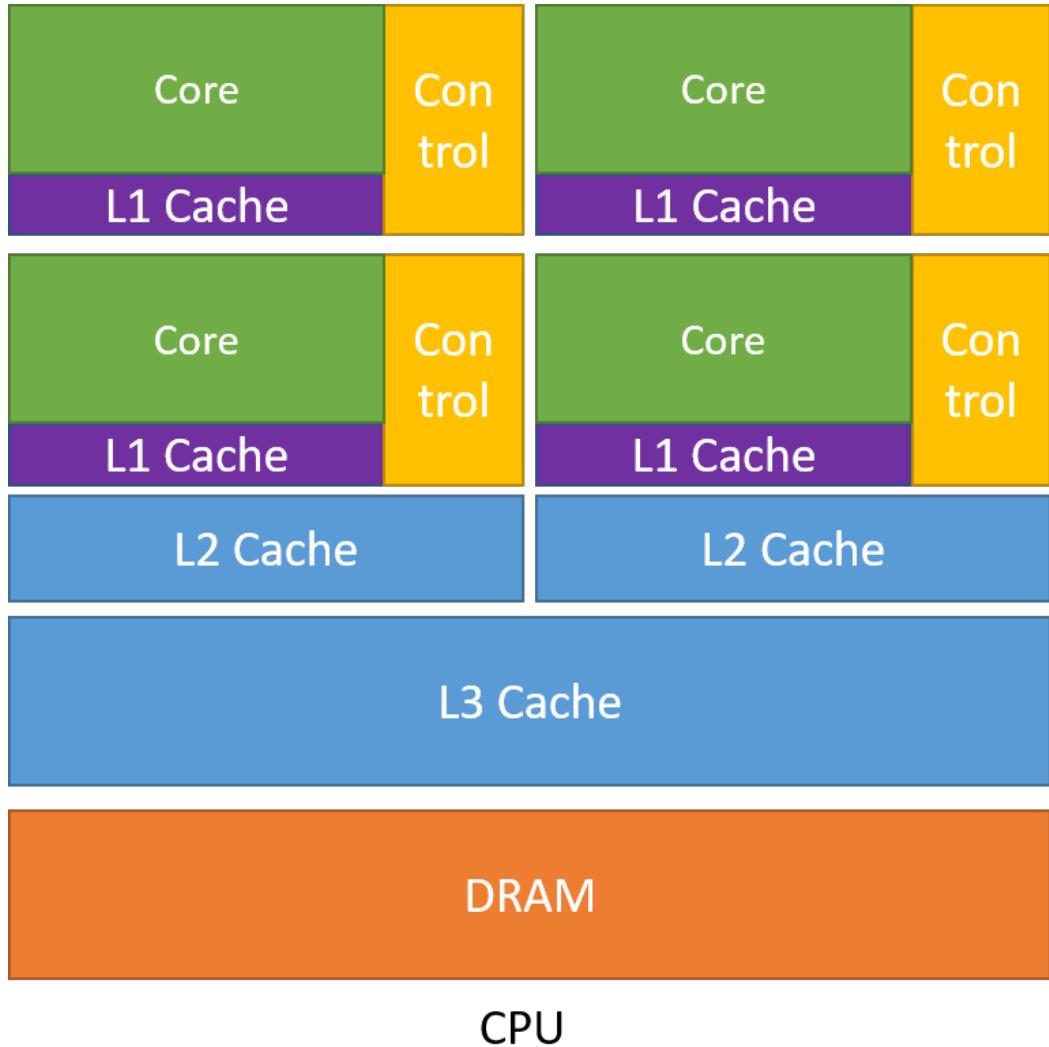
Yi Lu<sup>1\*</sup>, Xin Zhou<sup>1\*</sup>, Wei He<sup>1</sup>, Jun Zhao<sup>1</sup>,  
Tao Ji<sup>1†</sup>, Tao Gui<sup>2†</sup>, Qi Zhang<sup>1†</sup>, Xuanjing Huang<sup>1,3</sup>

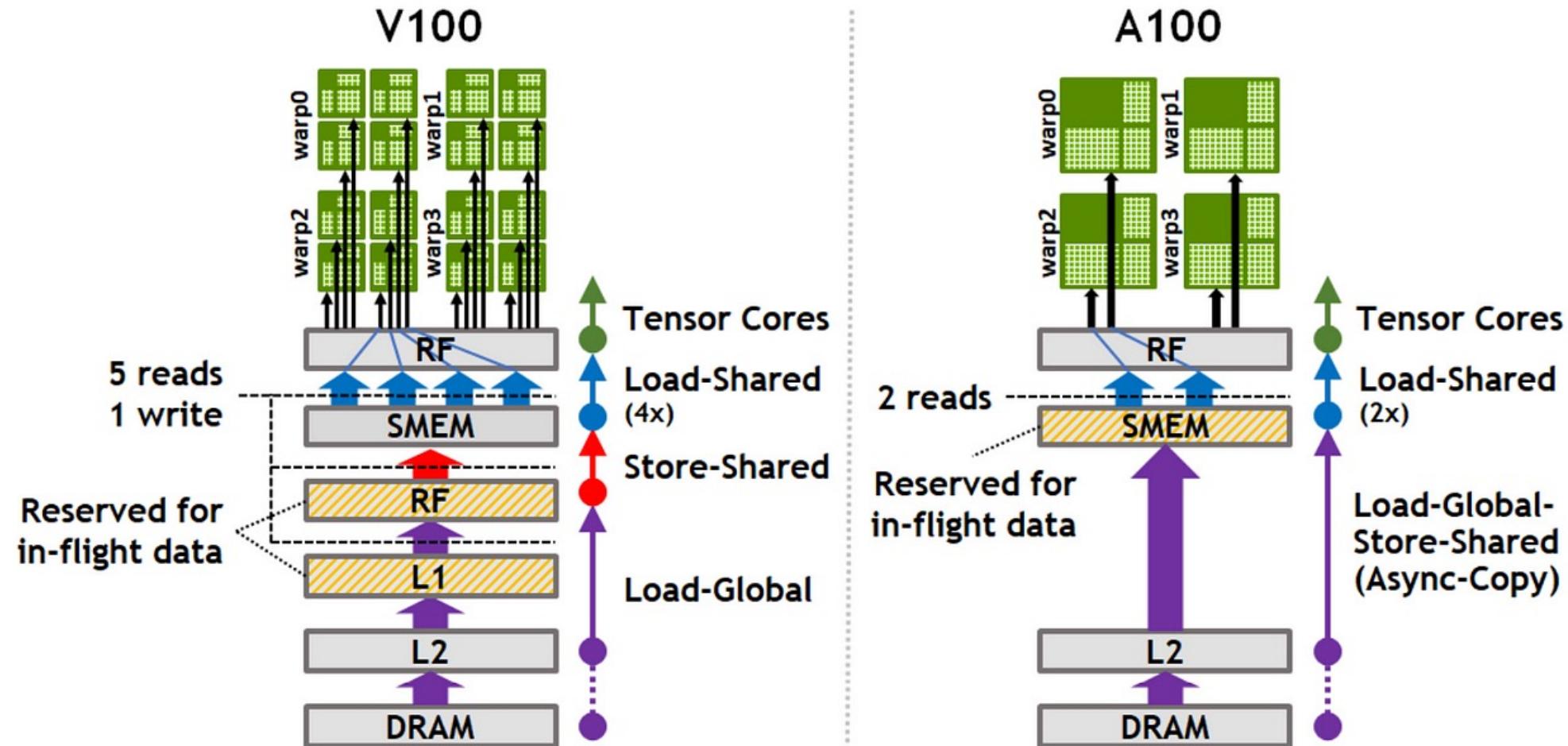
### 上下文分块+稀疏选择



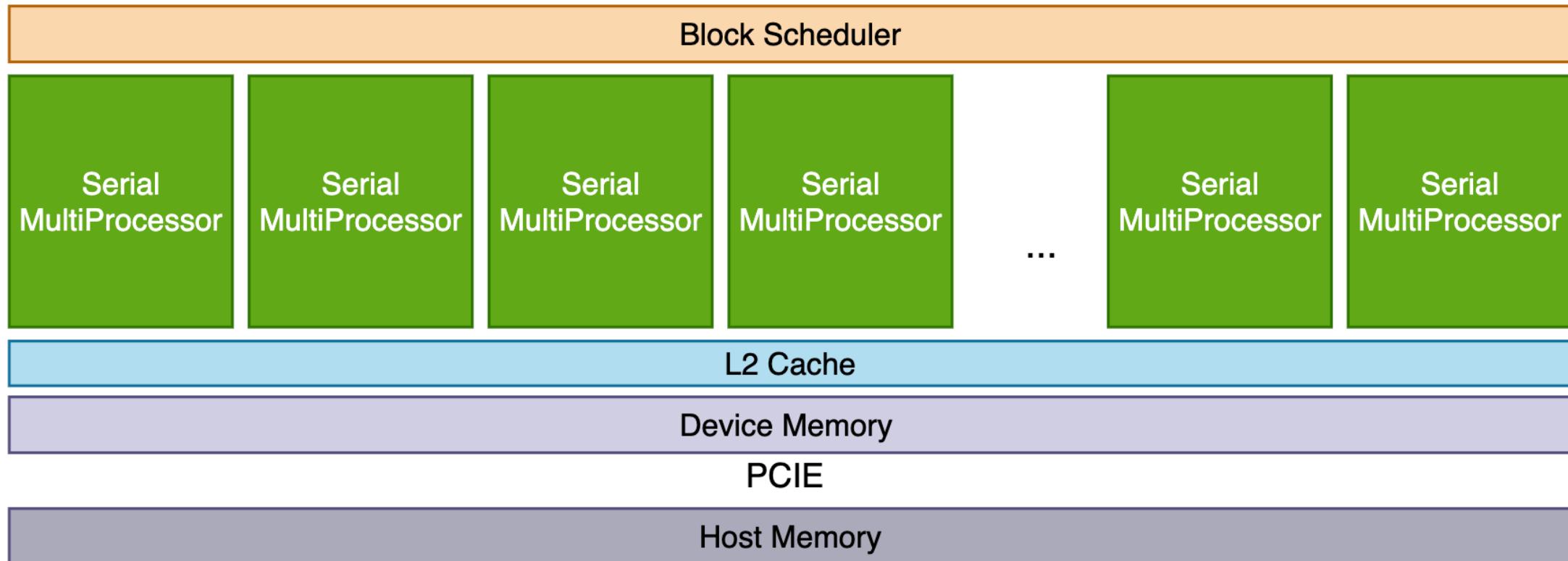
4/3/2.5 4(meta) but findings 🤔

# GPU vs. CPU

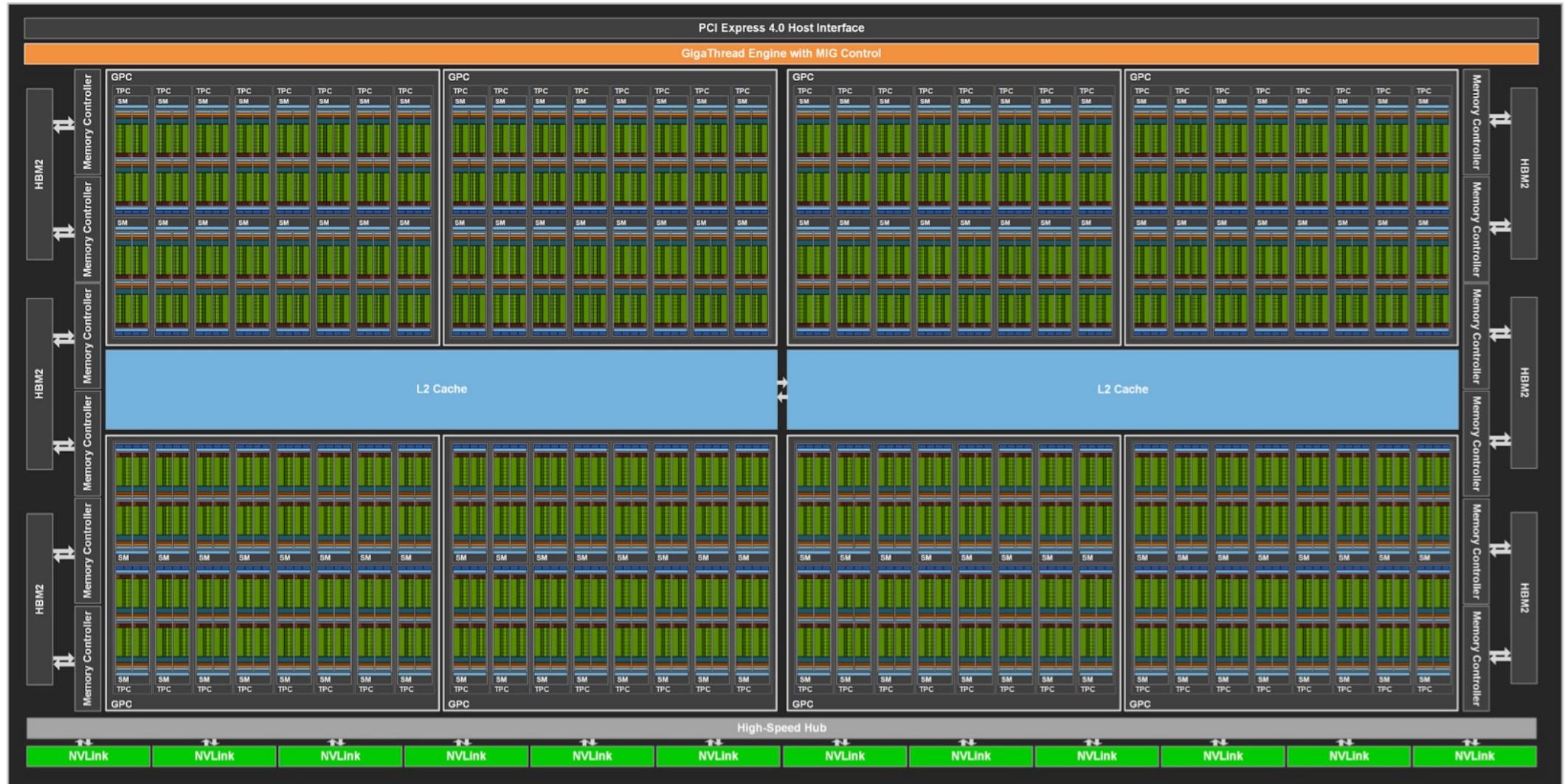


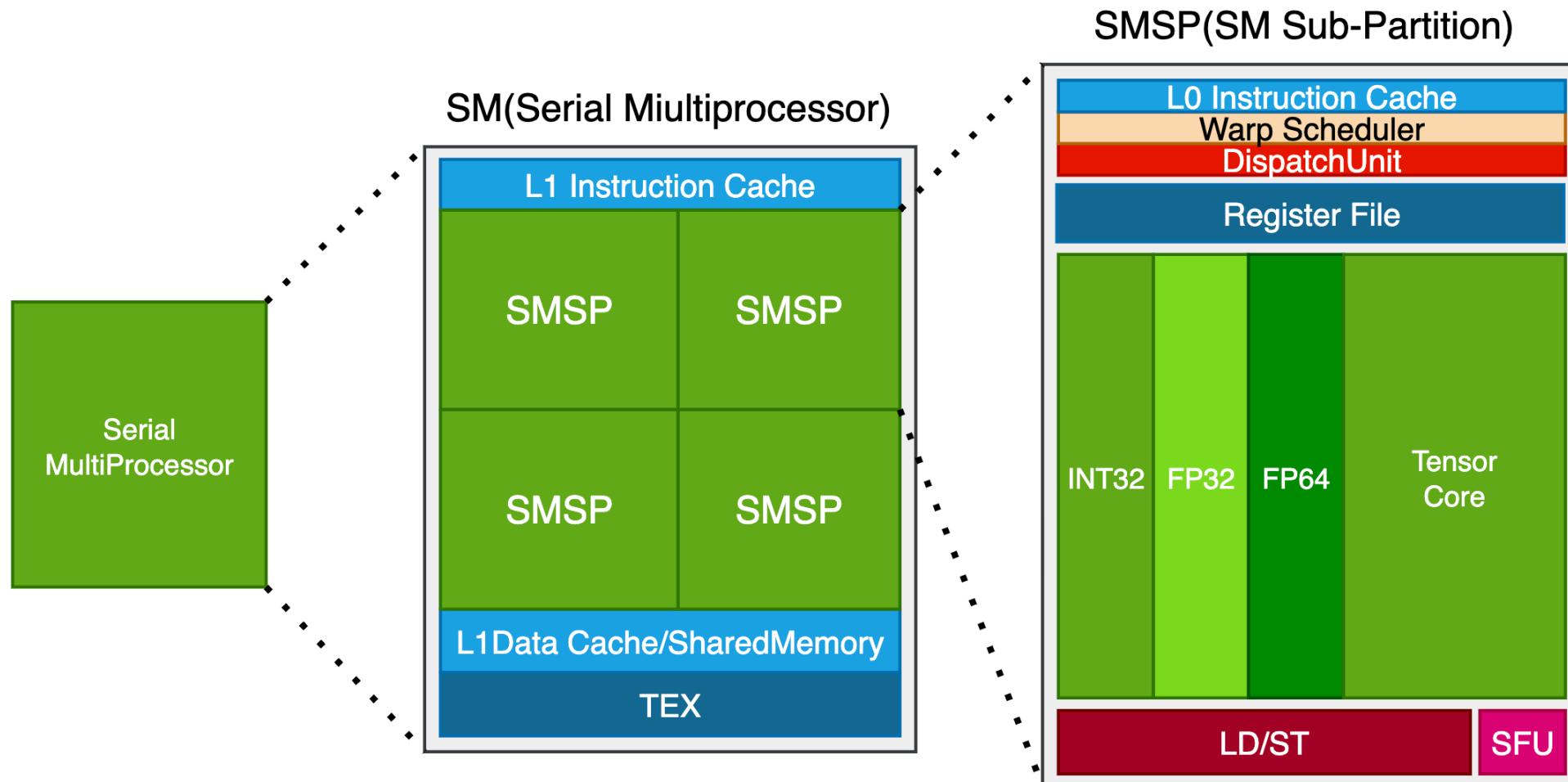


# GPU



# GPU

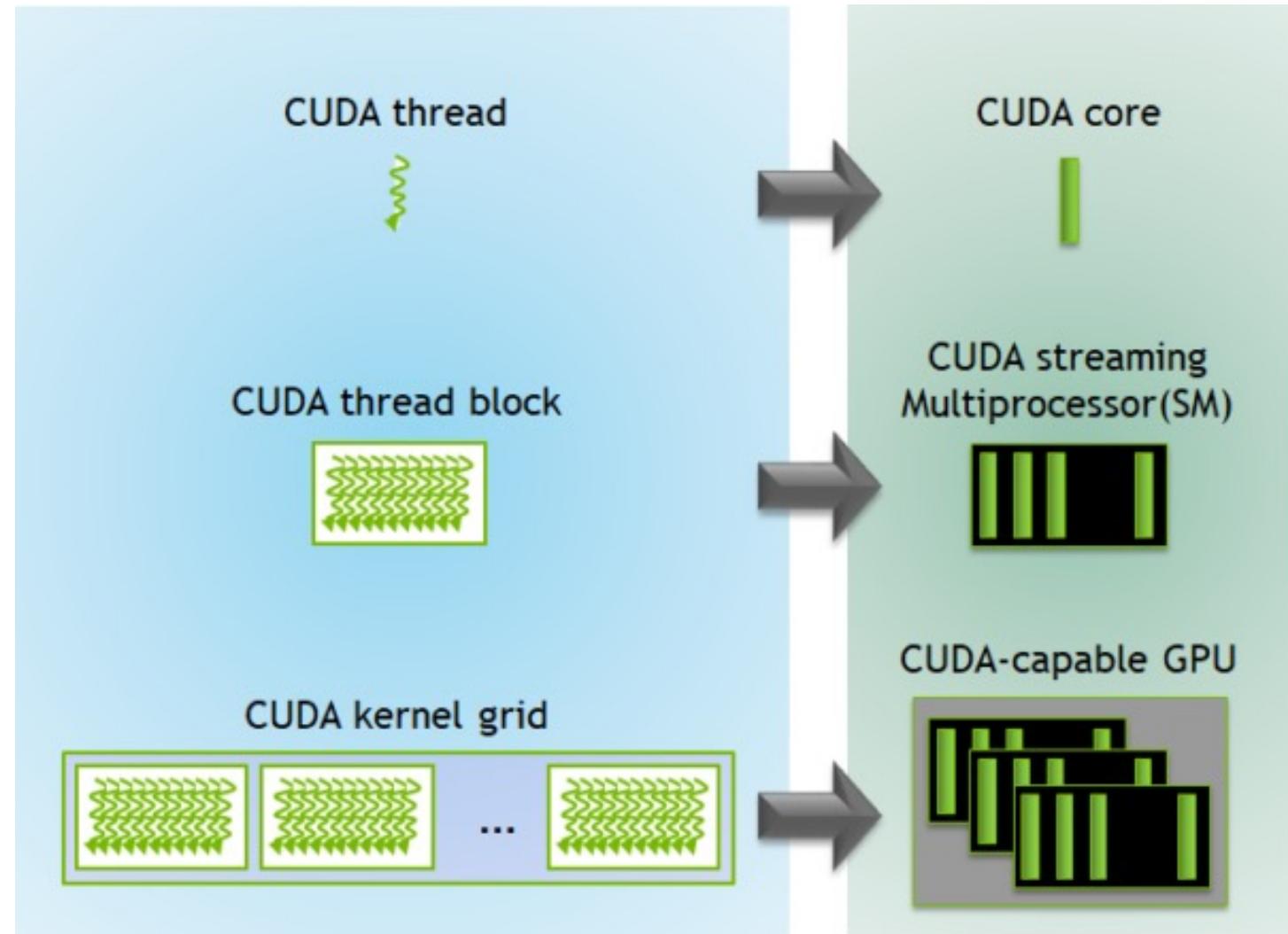




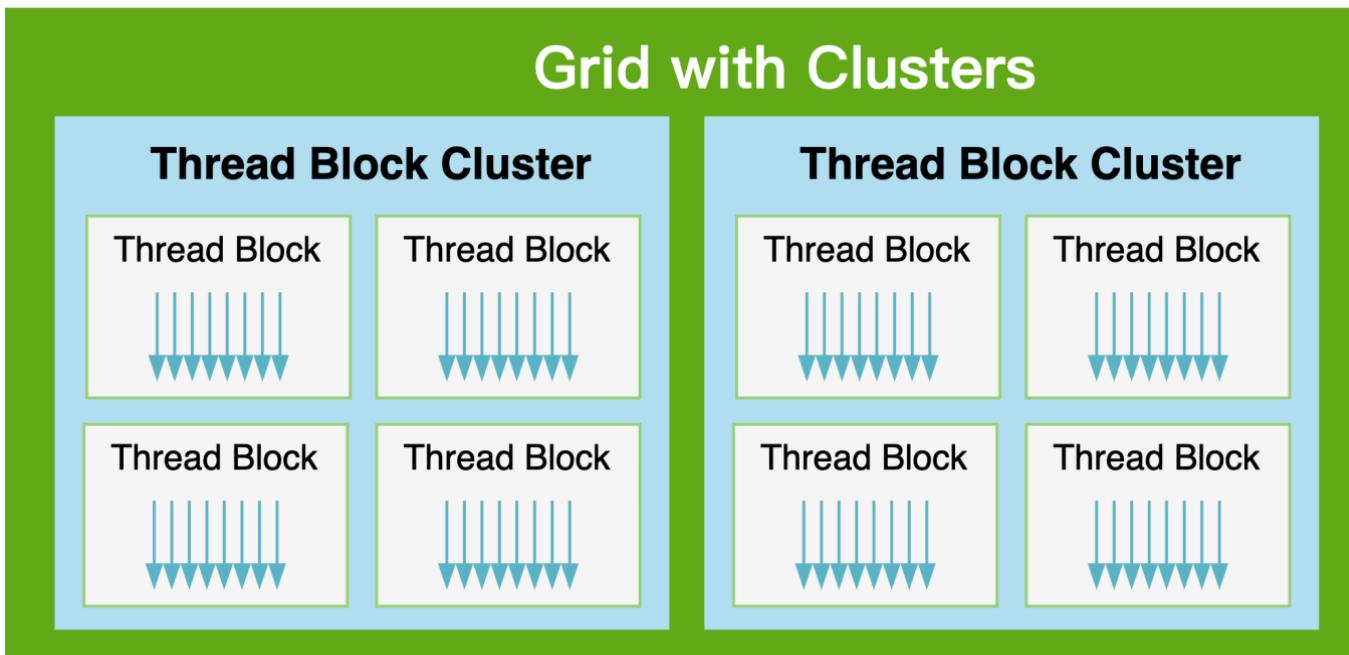
# GPU



# CUDA Programming



# CUDA Programming



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

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

# CUDA Programming: SIMD

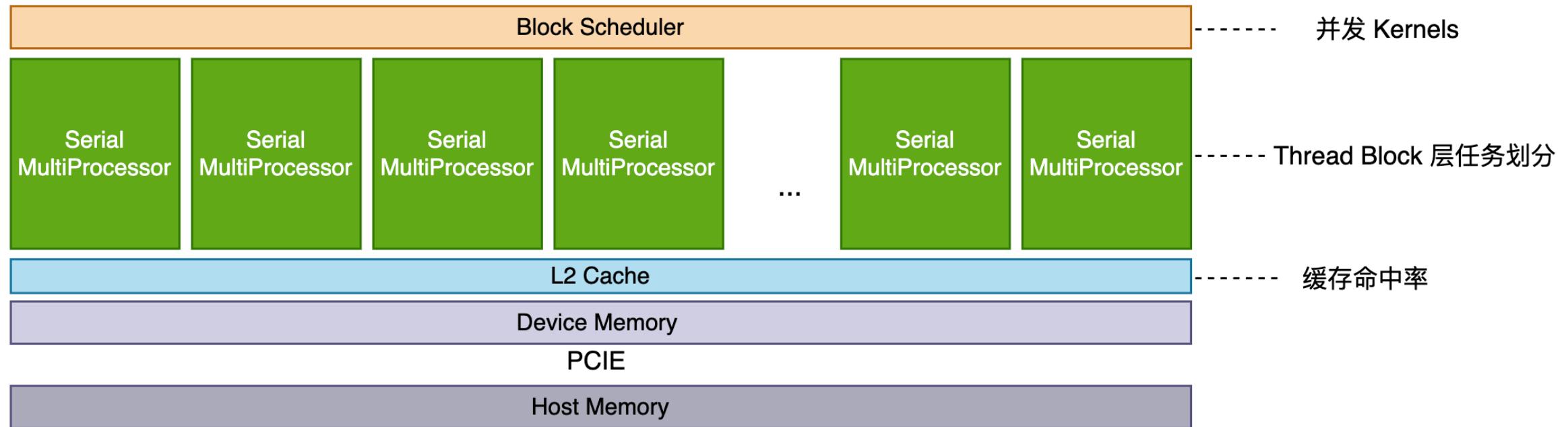


•  
•  
•



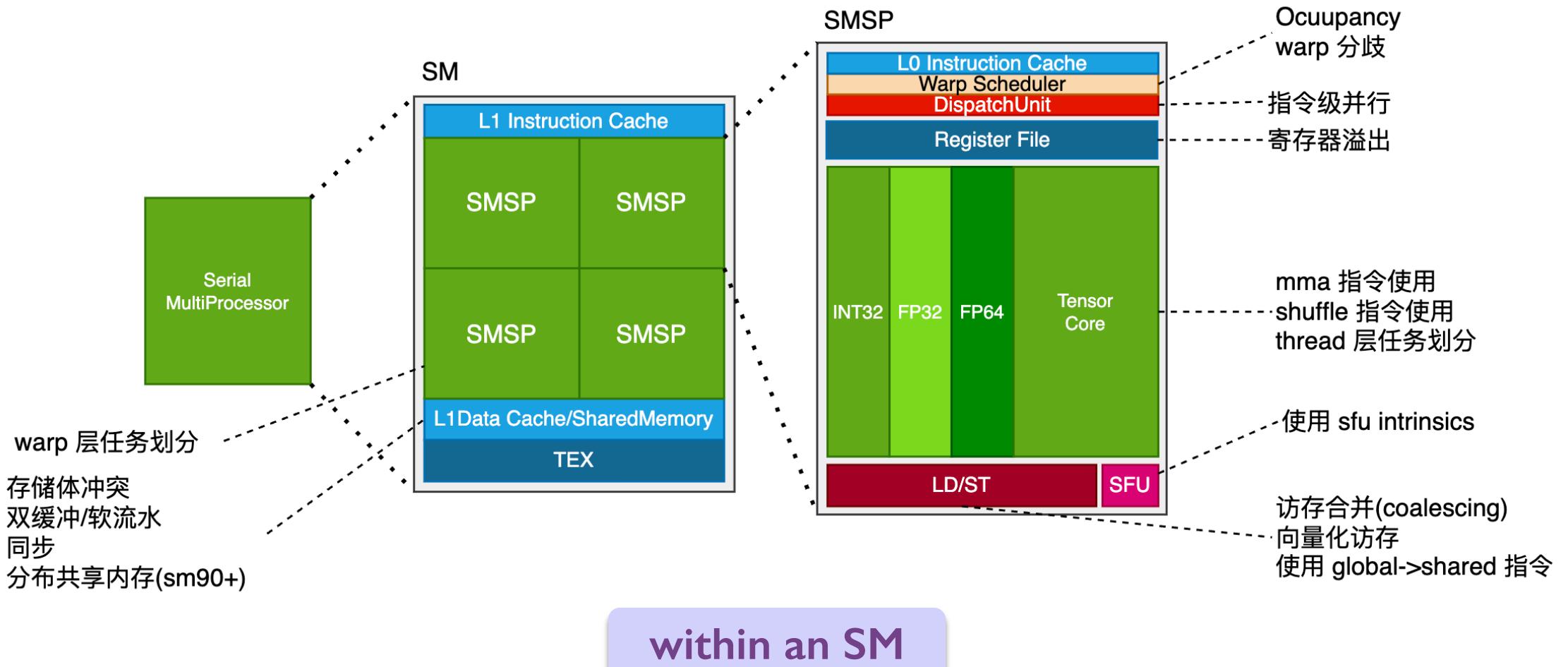
- ▶ ID
- ▶  $R, C = ID // \text{Col}, ID \% \text{Col}$
- ▶

# CUDA Programming



outside an SM

# CUDA Programming



# C++: Hello World!



hello\_world.py

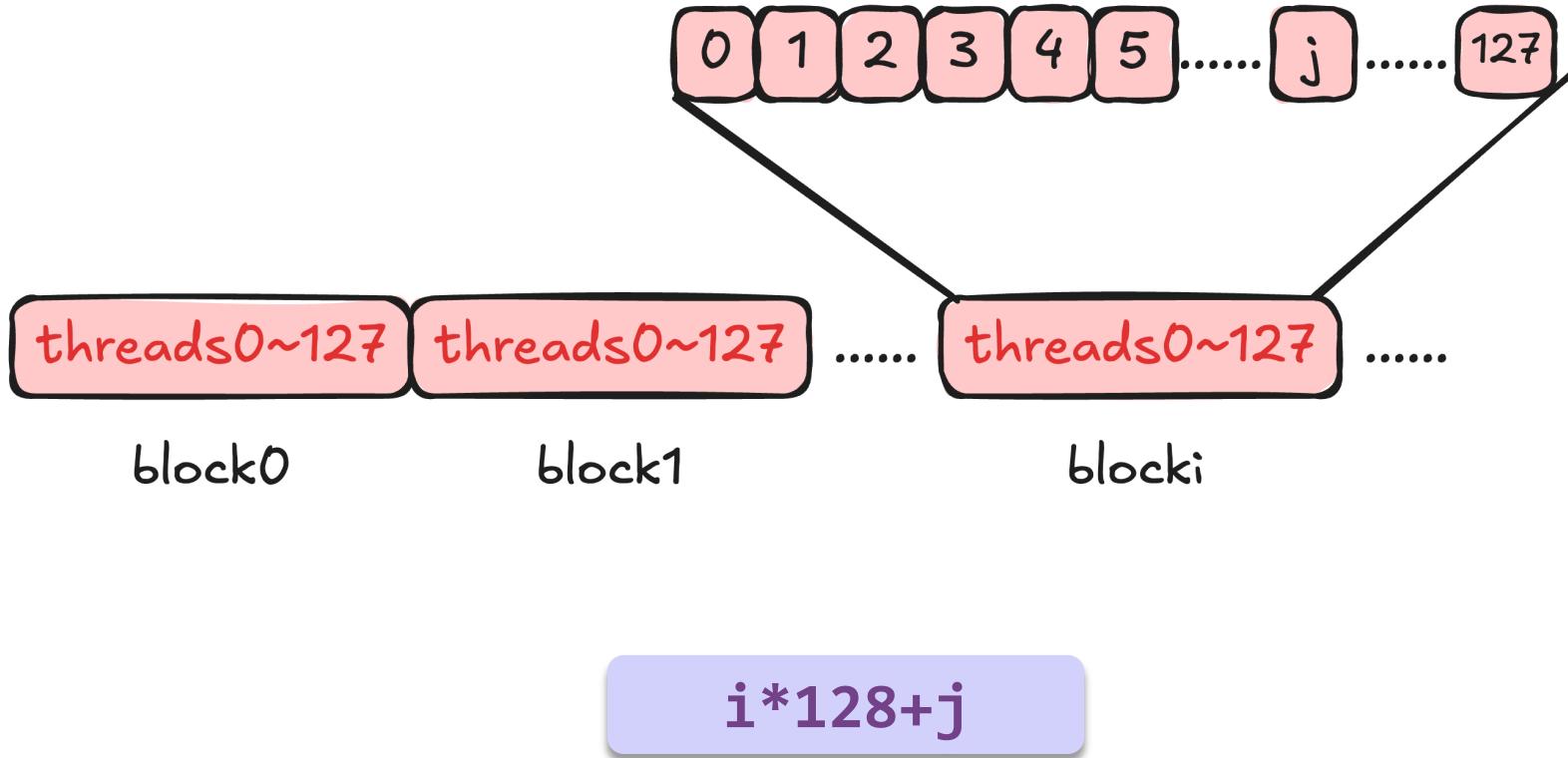
```
1 import torch
2 from torch.utils.cpp_extension import load_inline
3
4 cpp_source = """
5     std::string hello_world() {
6         return "Hello World!";
7     }
8 """
9
10 module = load_inline(
11     name='moss_op',
12     cpp_sources=[cpp_source],
13     functions=['hello_world'],
14     verbose=True,
15 )
16
17 print(module.hello_world())
```

load\_inline

```
1 import torch
2 from torch.utils.cpp_extension import load_inline
3
4 cuda_source = '''
+ 5 __global__ void square_vector_kernel(const float* vector, float* result, int col) {
+ 6     int c = blockIdx.x * blockDim.x + threadIdx.x;
7         ...
8 }
```

修饰符	调用者	执行位置	典型用途
<code>__global__</code>	host	device	GPU kernel, 从 CPU 发起
<code>__device__</code>	device	device	GPU 上运行的辅助函数
<code>__host__</code>	host	host	CPU 上的普通函数 (默认)

```
20
+ 21     return result;
+ 22 }
23 ''
24
25 cpp_source = "torch::Tensor square_vector(torch::Tensor vector);"
26
27 module = load_inline(
28     name='moss_op',
29     cpp_sources=cpp_source,
+ 30     cuda_sources=cuda_source,
31     functions=['square_vector'],
32 )
33
34 a = torch.tensor([1., 2., 3.], device='cuda')
35 print(module.square_vector(a))
```



# Profiler



cuda\_profile.py

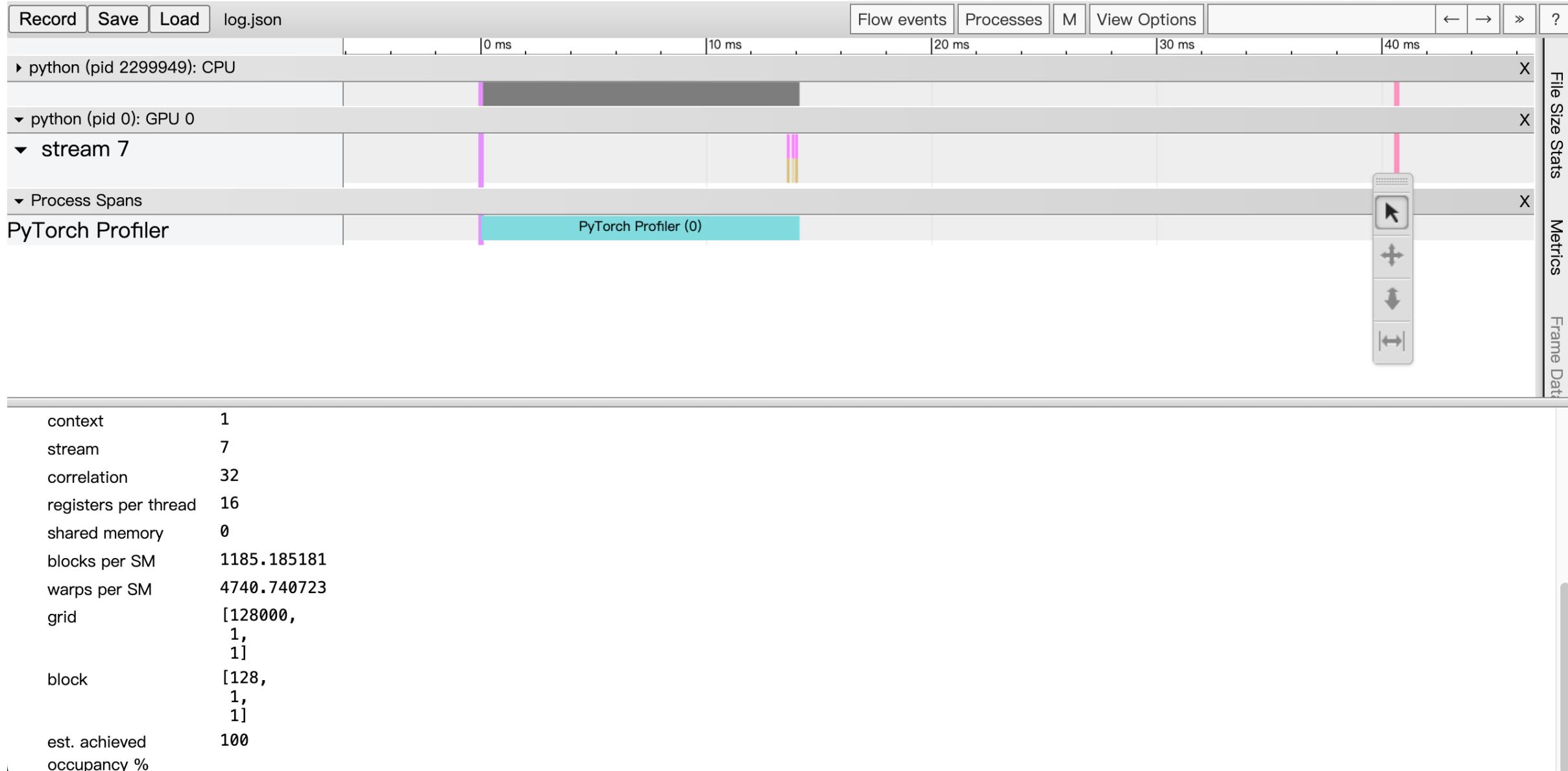
```
1  with profile(
2      activities=[ProfilerActivity.CUDA],
3      schedule=torch.profiler.schedule(wait=2, warmup=2, active=3, repeat=1),
4      on_trace_ready=torch.profiler.tensorboard_trace_handler('./log/moss_op'),
5      record_shapes=True,
6      with_stack=True,
7      with_flops=True
8  ) as prof:
9      for i in range(7):
10          with record_function("square_vector_custom"):
11              module.square_vector(a)
12          prof.step()
13
14 print(prof.key_averages().table(sort_by="cuda_time_total", row_limit=10))
```

```
tensor([9.7066e-01, 6.9254e-02, 5.9329e-03, ..., 1.6538e-03, 3.4866e-05,
       2.4144e-01], device='cuda:0')
```

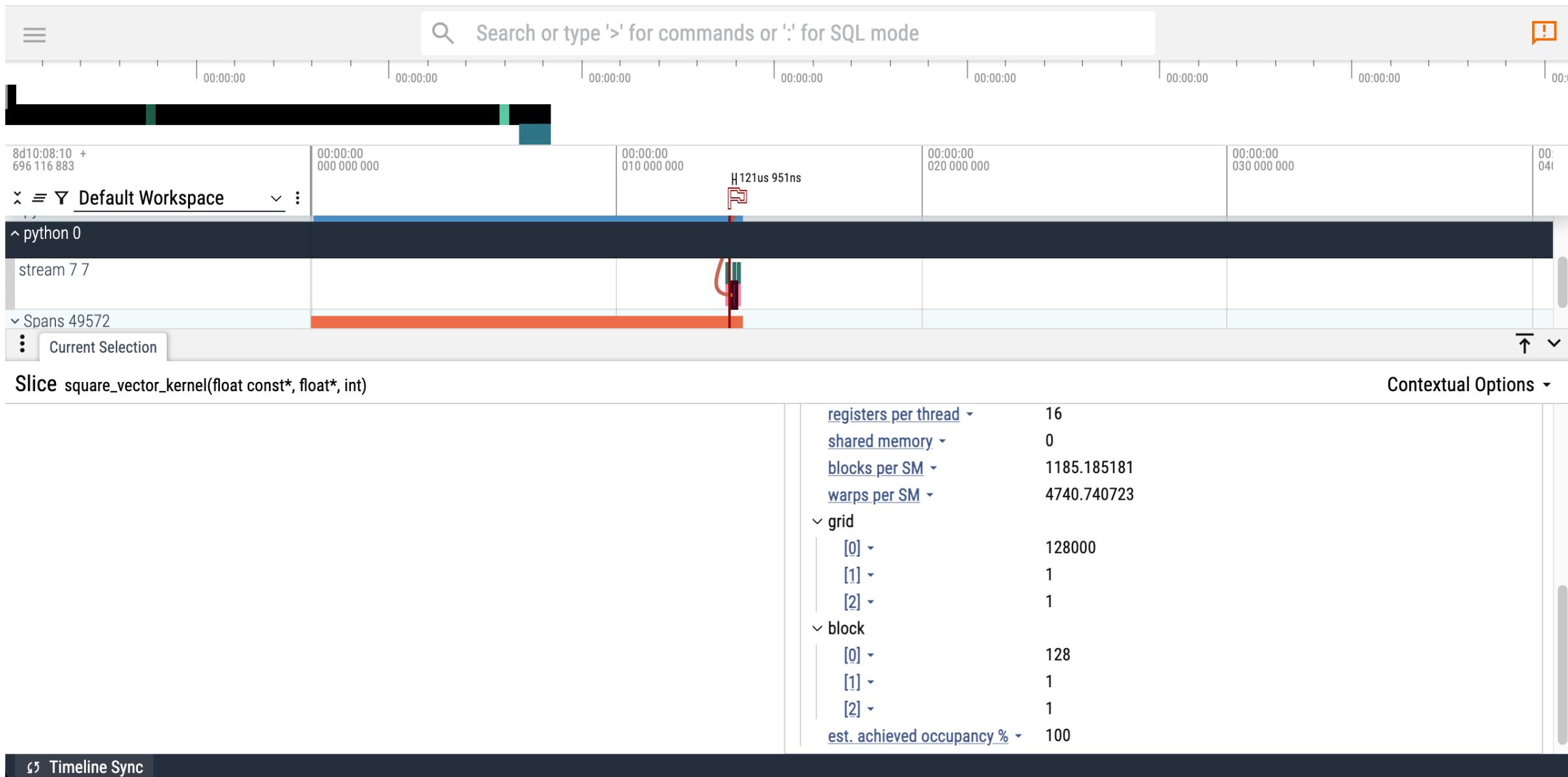
Name	Self CPU %	Self CPU	CPU total %	CPU total	CPU time avg	Self CUDA	Self CUDA %	CUDA total	CUDA time avg	# of Calls
square_vector_kernel(float const*, float*, int)	0.00%	0.000us	0.00%	0.000us	0.000us	364.991us	100.00%	364.991us	121.664us	3
cudaLaunchKernel	24.99%	35.837us	24.99%	35.837us	11.946us	0.000us	0.00%	0.000us	0.000us	3
cudaDeviceSynchronize	75.01%	107.547us	75.01%	107.547us	107.547us	0.000us	0.00%	0.000us	0.000us	1

```
Self CPU time total: 143.384us
Self CUDA time total: 364.991us
```

# Profiler



# Profiler



## 作业一：Profiler LLMs Module

- Attention
- LayerNorm
- Feedforward
- LMhead

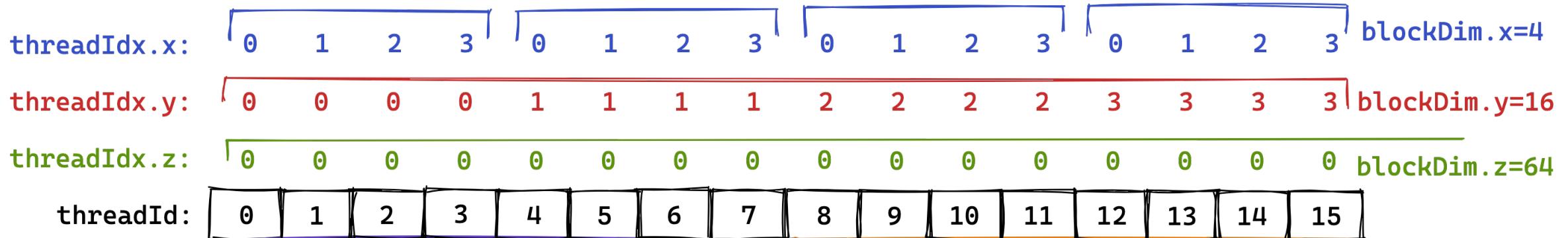
.....

# 2D/3D Block



blocks.cu

```
1 dim3 threads_per_block(4, 16, 64);  
2 dim3 blocks_per_grid(cdiv(N, threads_per_block.x), cdiv(M, threads_per_block.y), cdiv(K, threads_per_block.z));
```



threadId = threadIdx.x + blockDim.x\*threadIdx.y + blockDim.x\*blockDim.y\*threadIdx.z

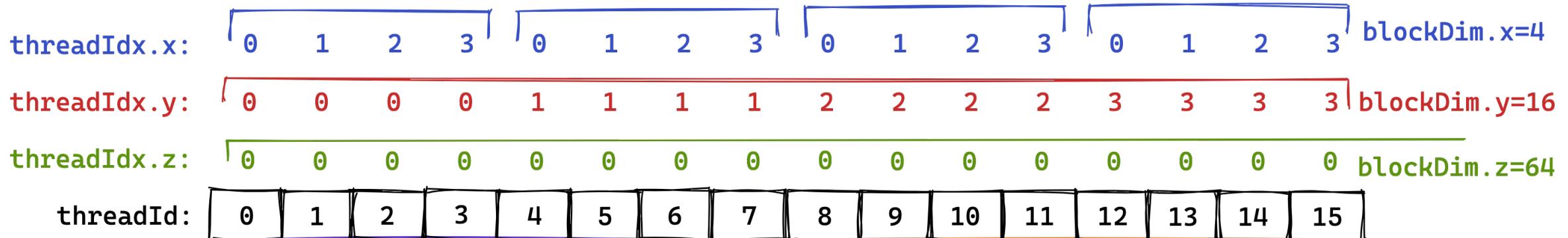
$$4 * 16 * 64 = 4096$$

# 2D/3D Block



blocks.cu

```
1 dim3 threads_per_block(4, 16, 64);  
2 dim3 blocks_per_grid(cdiv(N, threads_per_block.x), cdiv(M, threads_per_block.y), cdiv(K, threads_per_block.z));
```



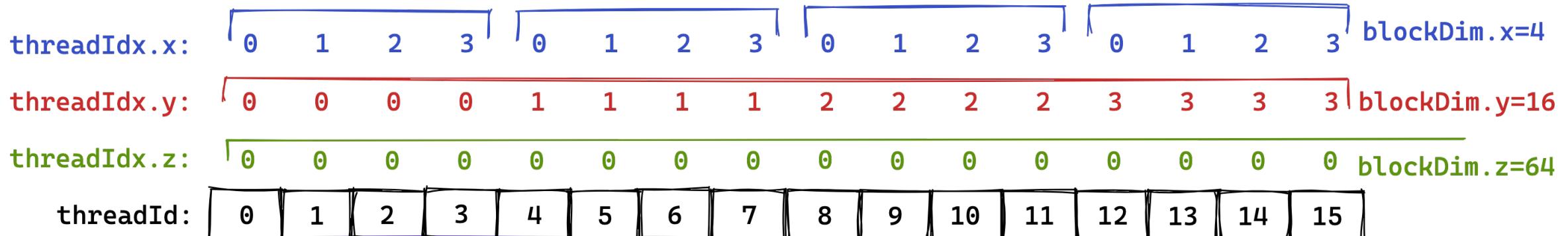
限制类型	数值	说明
每个 block 最多线程数	1024	即: <code>blockDim.x * blockDim.y * blockDim.z &lt;= 1024</code>
每个维度的最大线程数	<code>blockDim.x &lt;= 1024</code> <code>blockDim.y &lt;= 1024</code> <code>blockDim.z &lt;= 64</code>	各维度上限依 GPU 架构而定 (较常见值)

# 2D/3D Block



blocks.cu

```
1 dim3 threads_per_block(4, 16, 64);
2 dim3 blocks_per_grid(cdiv(N, threads_per_block.x), cdiv(M, threads_per_block.y), cdiv(K, threads_per_block.z));
```



check\_device.py

```
1 import torch
2
3 prop = torch.cuda.get_device_properties(0)
4
5 print("Device name:", prop.name)
6 print("Max threads per block:", prop.max_threads_per_block)
7 print("Max threads per SM:", prop.max_threads_per_multi_processor)
8 print("Warp size:", prop.warp_size)
9 print("Number of SMs:", prop.multi_processor_count)
```

# SiLU activation



warp\_coalesced.py

```
1 rows, cols = 8192, 8192
2 x = torch.randn((rows, cols), device="cuda")
3
4 with profile(
5     activities=[ProfilerActivity.CUDA],
6     schedule=torch.profiler.schedule(wait=1, warmup=1, active=2),
7     on_trace_ready=torch.profiler.tensorboard_trace_handler("./log/warp_coalesced"),
8     record_shapes=False,
9     with_stack=False,
10 ) as prof:
11     for _ in range(4):
12         with record_function("silu_torch"), torch.no_grad():
13             out0 = torch.nn.functional.silu(x)
14         with record_function("silu_coalesced"):
15             out1 = module.silu_coalesced(x)
16         with record_function("silu_uncoalesced"):
17             out2 = module.silu_uncoalesced(x)
18         prof.step()
19
20 assert torch.allclose(out1, out0) and torch.allclose(out2, out0)
21 print(prof.key_averages().table(sort_by="cuda_time_total", row_limit=10))
```

# SiLU activation

Loading extension module moss\_op...

Name	Self CPU %	Self CPU	CPU total %	CPU total	CPU time avg	Self CUDA	Self CUDA %	CUDA total	CUDA time avg
silu_uncoalesced_kernel(float const*, float*, int, i...	0.00%	0.000us	0.00%	0.000us	0.000us	7.175ms	81.10%	7.175ms	2.392ms
silu_coalesced_kernel(float const*, float*, int, int...	0.00%	0.000us	0.00%	0.000us	0.000us	1.067ms	12.06%	1.067ms	533.326us
void at::native::vectorized_elementwise_kernel<4, at...	0.00%	0.000us	0.00%	0.000us	0.000us	605.246us	6.84%	605.246us	302.623us
cudaLaunchKernel	1.15%	69.085us	1.15%	69.085us	11.514us	0.000us	0.00%	0.000us	0.000us
cudaDeviceSynchronize	98.85%	5.954ms	98.85%	5.954ms	5.954ms	0.000us	0.00%	0.000us	0.000us

Self CPU time total: 6.024ms

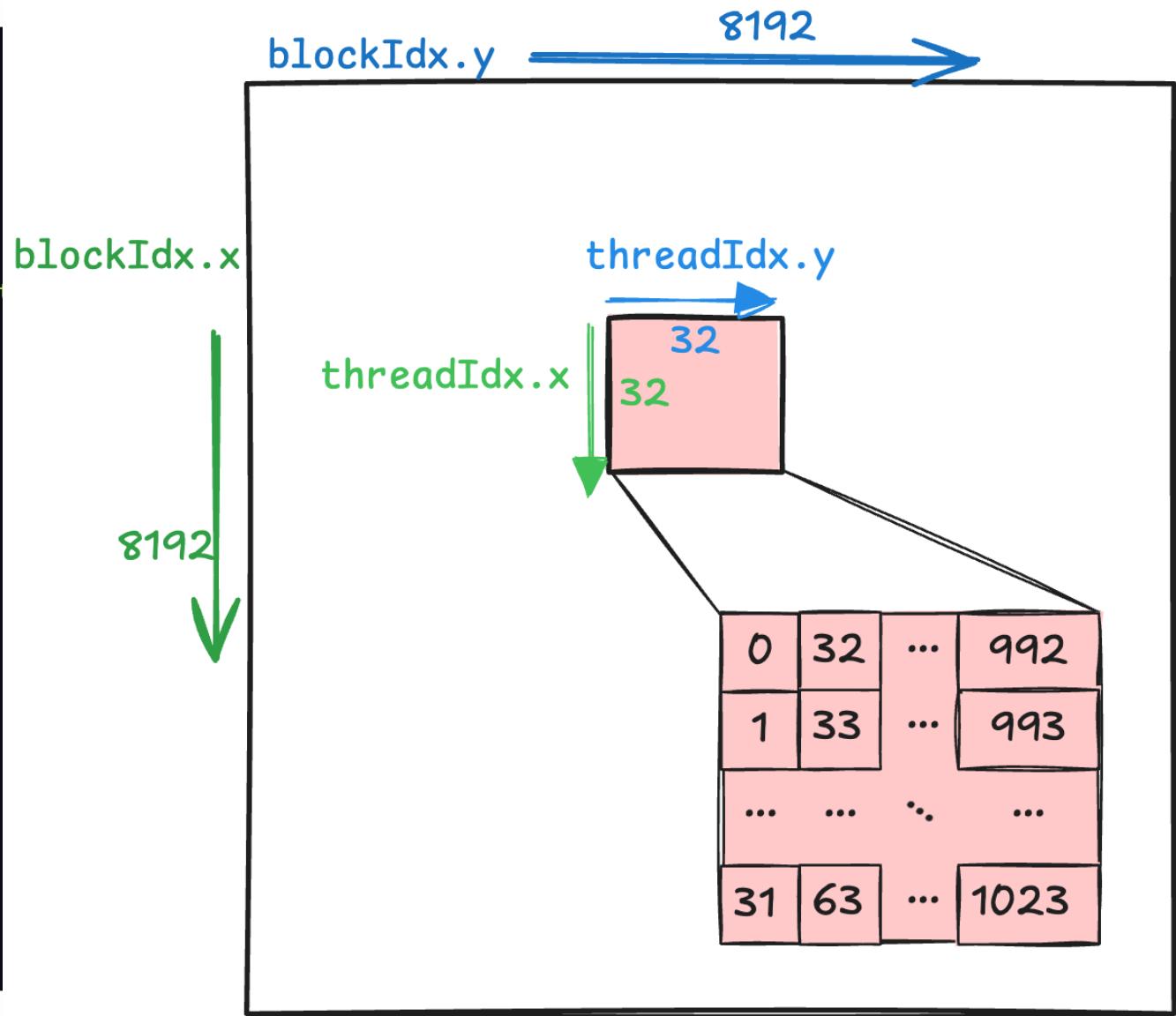
Self CUDA time total: 8.847ms

# SiLU activation

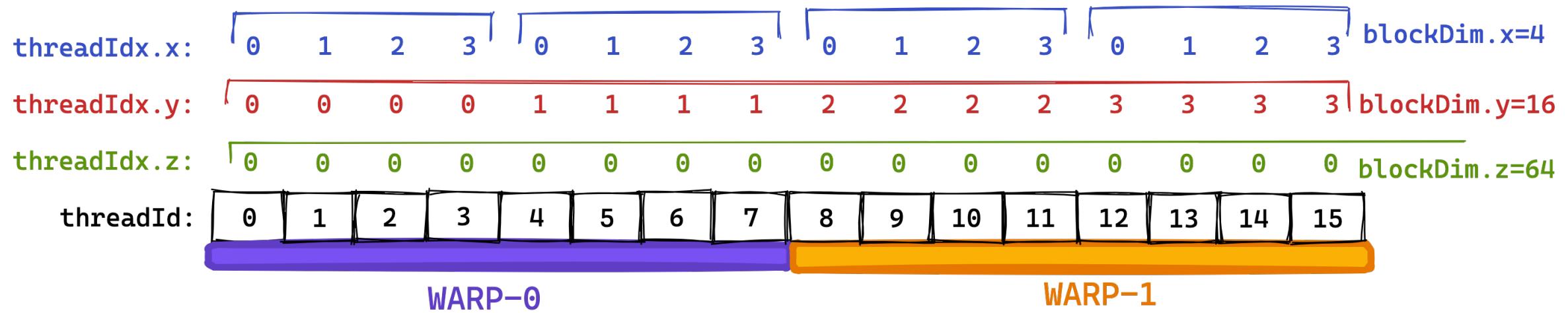


warp\_coalesced.py

```
1  cuda_source = """
2  __device__ float silu(float x) {
3      return x / (1.0f + expf(-x));
4  }
5
6  __global__ void silu_uncoalesced_kernel(const float* input, float*
7      int col = blockIdx.y * blockDim.y + threadIdx.y;
8      int row = blockIdx.x * blockDim.x + threadIdx.x;
9      if (row < rows && col < cols) {
10          int idx = row * cols + col;
11          output[idx] = silu(input[idx]);
12      }
13  }
14
15 torch::Tensor silu_uncoalesced(torch::Tensor input) {
16     const int rows = input.size(0);
17     const int cols = input.size(1);
18     auto output = torch::empty_like(input);
19
20     dim3 threads(32, 32);
21     dim3 blocks((rows + threads.x - 1) / threads.x,
22                 (cols + threads.y - 1) / threads.y);
23     silu_uncoalesced_kernel<<<blocks, threads>>>(
24         input.data_ptr<float>(), output.data_ptr<float>(), rows,
25     return output;
26 }
```

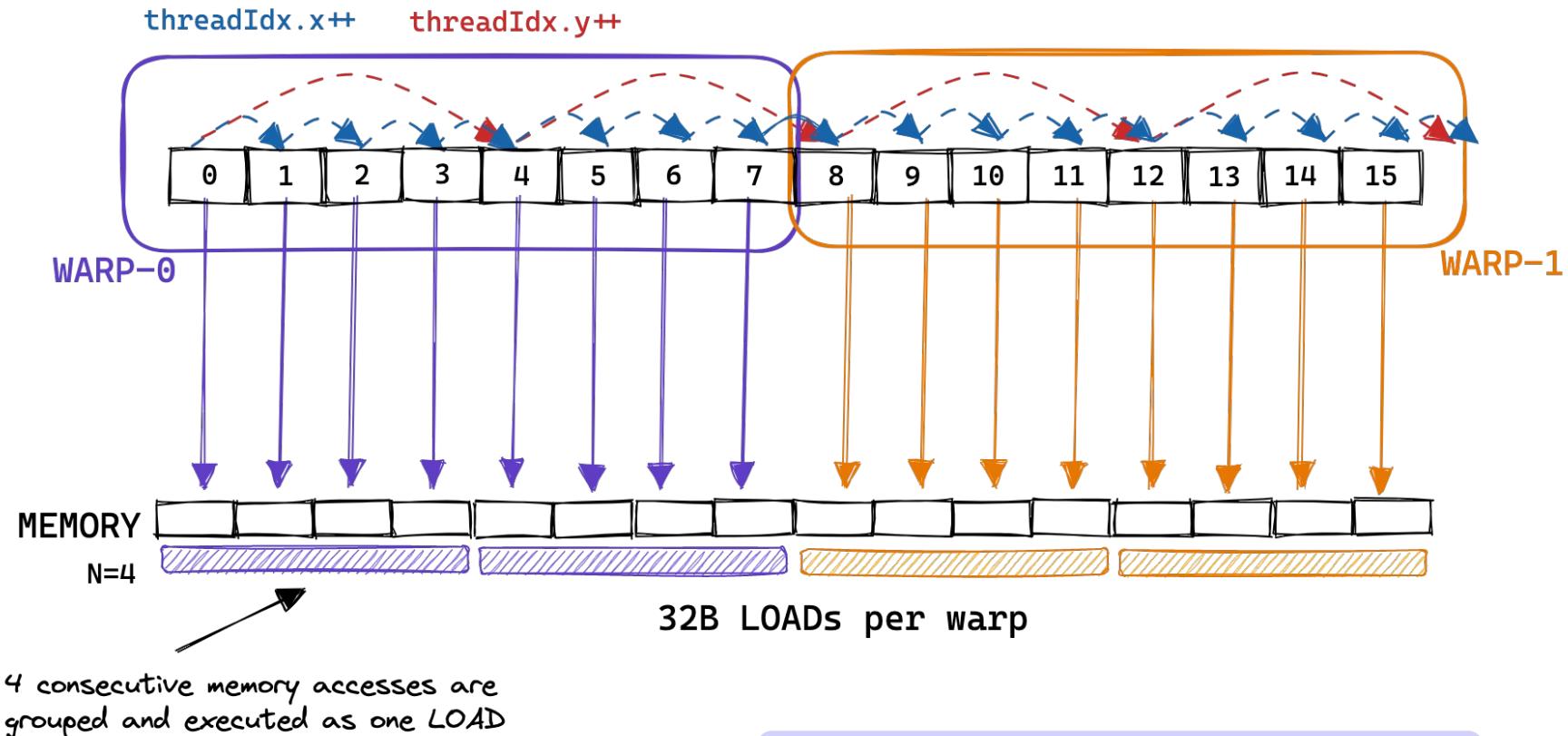


# Warp: 线程束/线程簇

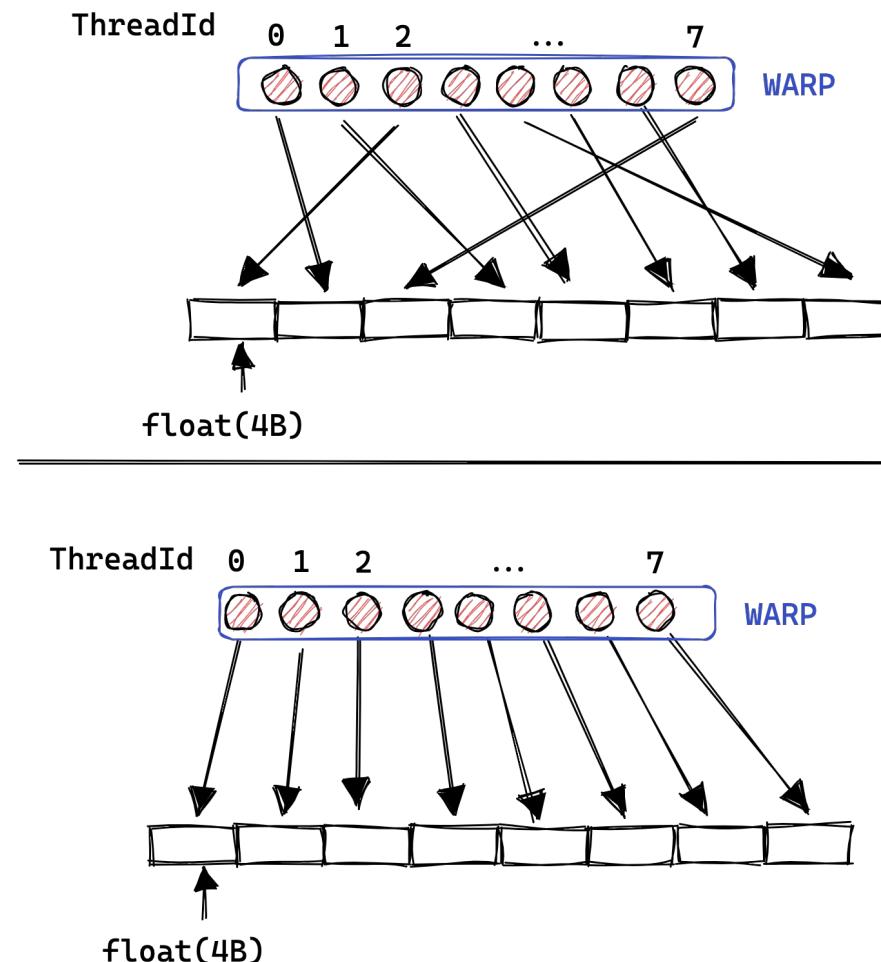


Warp=32

# Warp: 线程束/线程簇



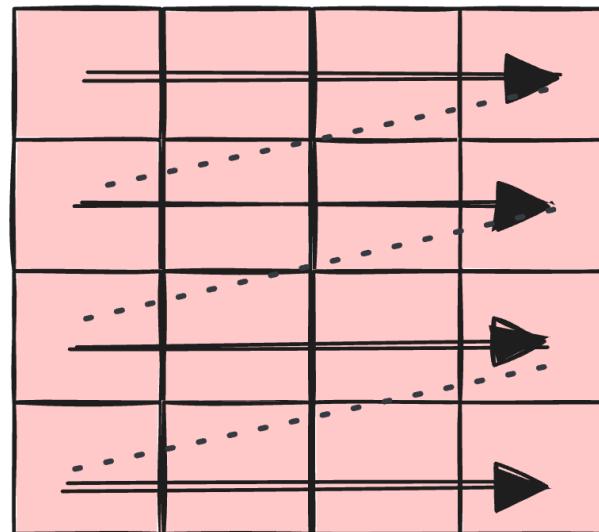
# Warp: 线程束/线程簇



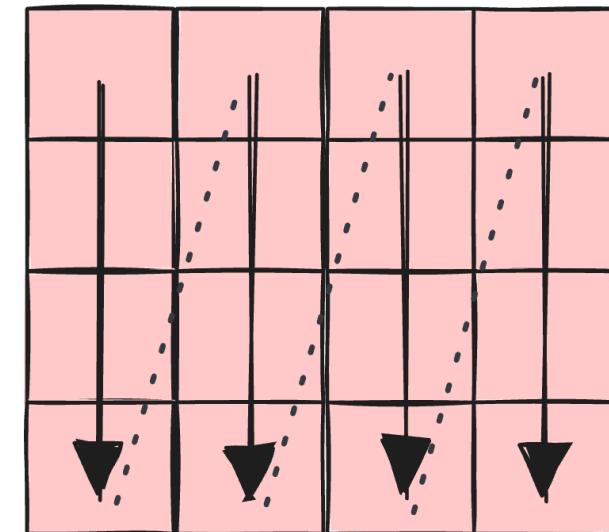
Both of these access patterns can be coalesced!

# Contiguous

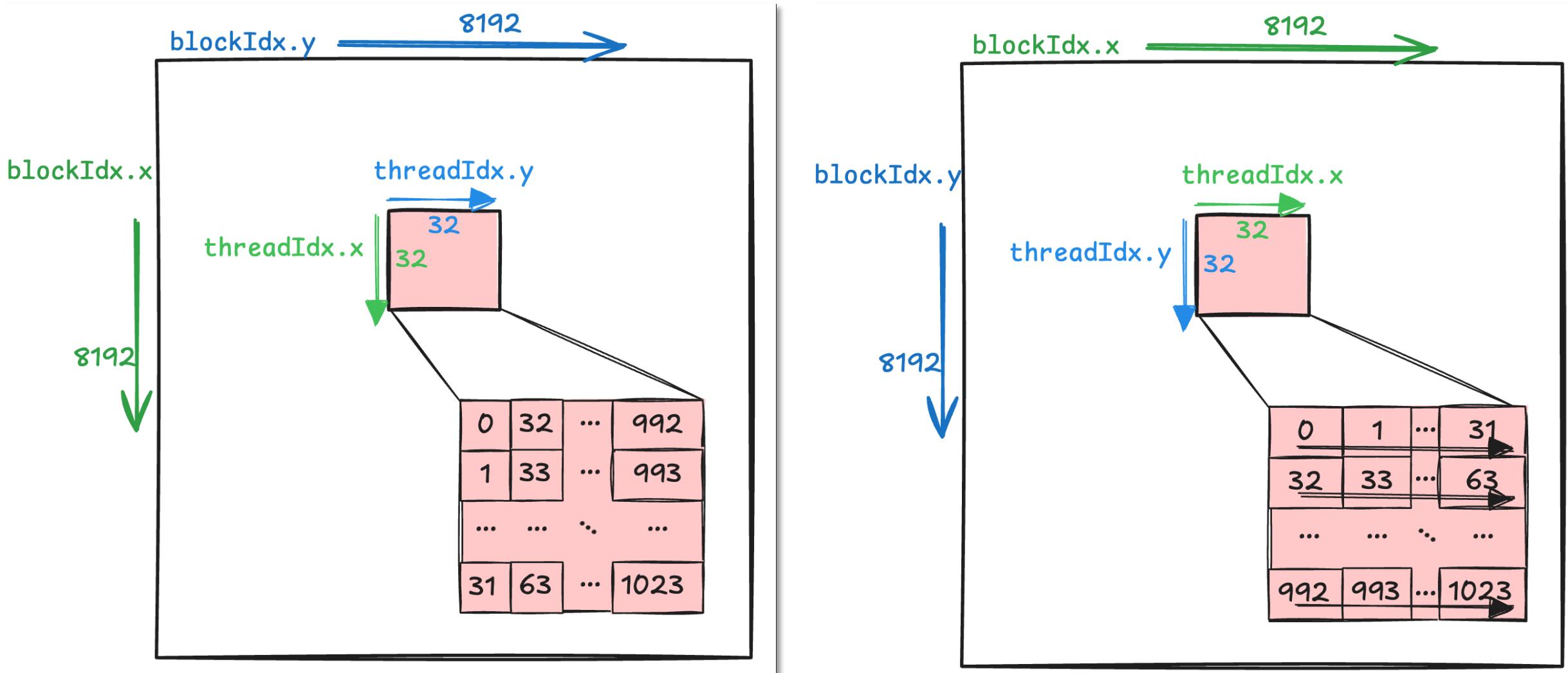
row-major (pytorch)



column-major (MATLAB)



# 访存合并: Coalesced



# 访存合并：Coalesced



warp\_coalesced.py

```
 1  __global__ void silu_coalesced_kernel(const float* input, float* output, int rows, int cols) {
+ 2      int col = blockIdx.x * blockDim.x + threadIdx.x;
+ 3      int row = blockIdx.y * blockDim.y + threadIdx.y;
 4      if (row < rows && col < cols) {
 5          int idx = row * cols + col;
 6          output[idx] = silu(input[idx]);
 7      }
 8  }
 9
10 torch::Tensor silu_coalesced(torch::Tensor input) {
11     const int rows = input.size(0);
12     const int cols = input.size(1);
13     auto output = torch::empty_like(input);
14
15     dim3 threads(32, 32);
+ 16     dim3 blocks((cols + threads.x - 1) / threads.x,
+ 17                  (rows + threads.y - 1) / threads.y);
18     silu_coalesced_kernel<<<blocks, threads>>>(
19         input.data_ptr<float>(), output.data_ptr<float>(), rows, cols);
20     return output;
21 }
```

# 访存合并：Coalesced

Loading extension module moss\_op...

Name	Self CPU %	Self CPU	CPU total %	CPU total	CPU time avg	Self CUDA	Self CUDA %	CUDA total	CUDA time avg
silu_uncoalesced_kernel(float const*, float*, int, i...	0.00%	0.000us	0.00%	0.000us	0.000us	7.175ms	81.10%	7.175ms	2.392ms
silu_coalesced_kernel(float const*, float*, int, int...	0.00%	0.000us	0.00%	0.000us	0.000us	1.067ms	12.06%	1.067ms	533.326us
void at::native::vectorized_elementwise_kernel<4, at...	0.00%	0.000us	0.00%	0.000us	0.000us	605.246us	6.84%	605.246us	302.623us
cudaLaunchKernel	1.15%	69.085us	1.15%	69.085us	11.514us	0.000us	0.00%	0.000us	0.000us
cudaDeviceSynchronize	98.85%	5.954ms	98.85%	5.954ms	5.954ms	0.000us	0.00%	0.000us	0.000us

Self CPU time total: 6.024ms

Self CUDA time total: 8.847ms

# 向量化访存: Vectorized

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



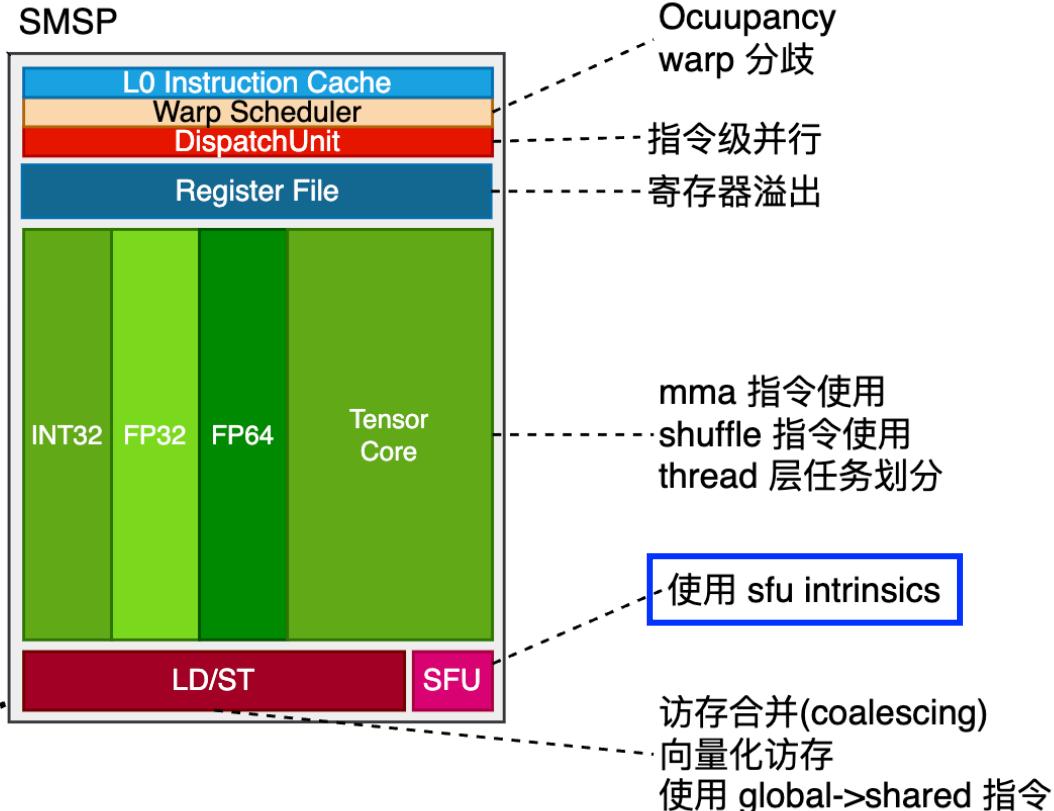
vectorized.py

```
1 __global__ void silu_coalesced_float4_kernel(const float* input, float* output, int rows, int cols) {
2     int row = blockIdx.y * blockDim.y + threadIdx.y;
3     int col4 = blockIdx.x * blockDim.x + threadIdx.x;
4     int col = col4 << 2;
5
6     if (row < rows && col + 3 < cols) {
7         const float4* input_v4 = reinterpret_cast<const float4*>(input);
8         float4 val = input_v4[row * (cols >> 2) + col4];
9
10        val.x = silu(val.x);
11        val.y = silu(val.y);
12        val.z = silu(val.z);
13        val.w = silu(val.w);
14
15        float4* output_v4 = reinterpret_cast<float4*>(output);
16        output_v4[row * (cols >> 2) + col4] = val;
17    }
18 }
19
20 torch::Tensor silu_coalesced_float4(torch::Tensor input) {
21     TORCH_CHECK(input.size(1) % 4 == 0, "Column size must be divisible by 4 for float4");
22
23     int rows = input.size(0);
24     int cols = input.size(1);
25
26     auto output = torch::empty_like(input);
27
28     dim3 threads(32, 32);
29     dim3 blocks((cols / 4 + threads.x - 1) / threads.x, (rows + threads.y - 1) / threads.y);
30     silu_coalesced_float4_kernel<<<blocks, threads>>>
31         input.data_ptr<float>(), output.data_ptr<float>(), rows, cols);
32     return output;
33 }
```

# 向量化访存: Vectorized

	Name	Self CPU %	Self CPU	CPU total %	CPU total	CPU time avg	Self CUDA	Self CUDA %	CUDA total	CUDA time avg
	silu_coalesced_float4_kernel(float const*, float*, i...	0.00%	0.000us	0.00%	0.000us	0.000us	638.077us	51.25%	638.077us	319.039us
	void at::native::vectorized_elementwise_kernel<4, at...	0.00%	0.000us	0.00%	0.000us	0.000us	607.037us	48.75%	607.037us	303.519us
	cudaLaunchKernel	4.70%	45.861us	4.70%	45.861us	11.465us	0.000us	0.00%	0.000us	0.000us
	cudaDeviceSynchronize	95.30%	929.920us	95.30%	929.920us	929.920us	0.000us	0.00%	0.000us	0.000us

# SFU Functions



函数	描述
<code>_sinf(x)</code>	快速计算 $\sin(x)$
<code>_cosf(x)</code>	快速计算 $\cos(x)$
<code>_expf(x)</code>	快速计算 $\exp(x)$
<code>_logf(x)</code>	快速计算 $\log(x)$
<code>_powf(x, y)</code>	快速计算 $x^y$
<code>_sqrtf(x)</code>	快速计算 $\sqrt{x}$
<code>_fdividef(a, b)</code>	快速除法

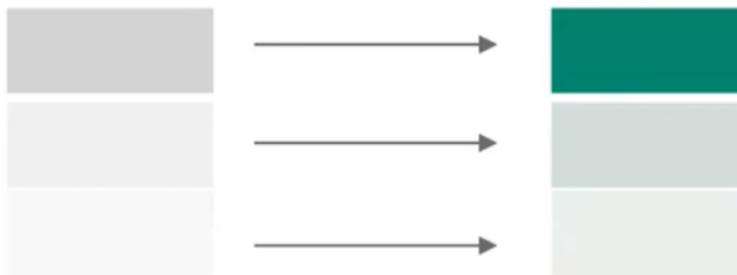
# SFU Functions



sfu.py

```
- 1  __device__ float silu(float x) {  
- 2      return x / (1.0f + expf(-x));  
- 3  }  
4  
+ 5  __device__ float silu_sfu(float x) {  
+ 6      return x / (1.0f + __expf(-x));  
+ 7  }
```

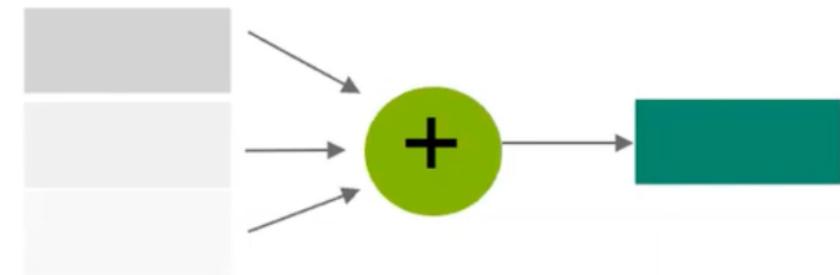
# Transformation vs Reduction



Transformation:

e.g.  $c[i] = a[i] + 10;$

Thread strategy: one thread per output point

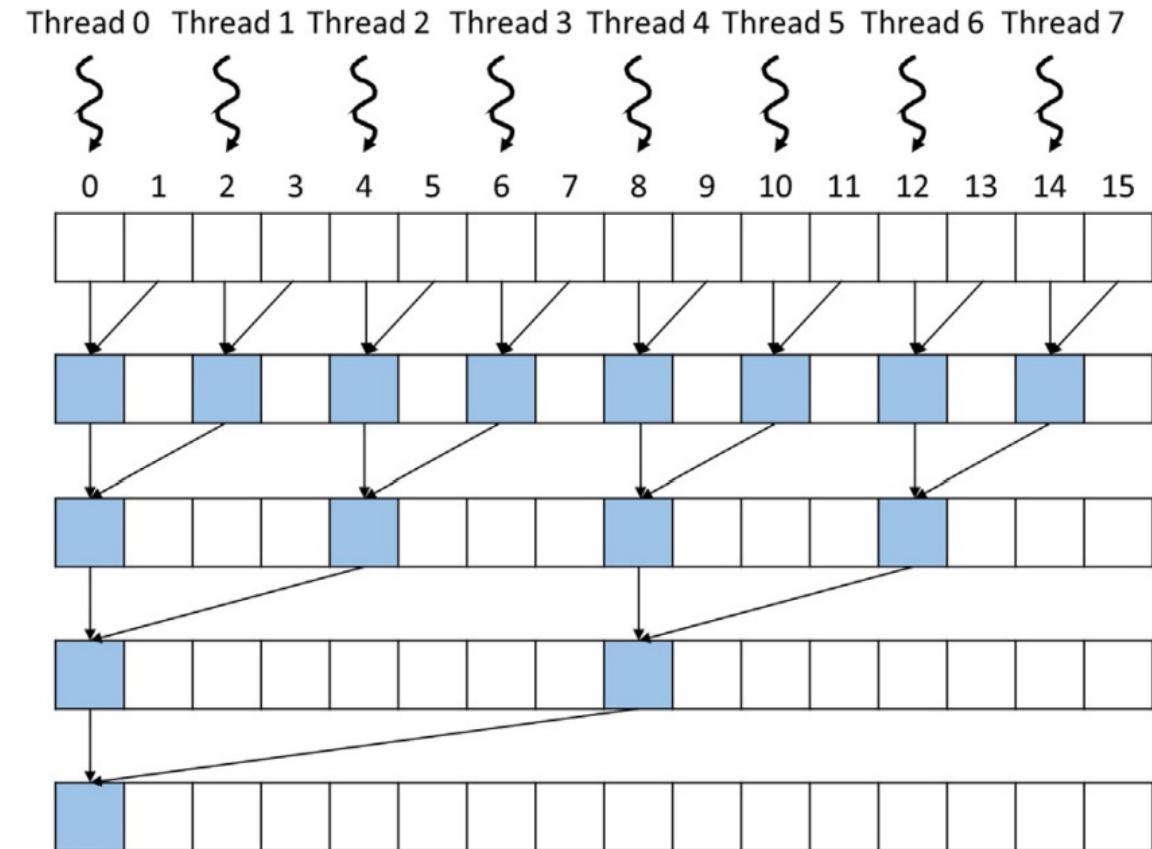
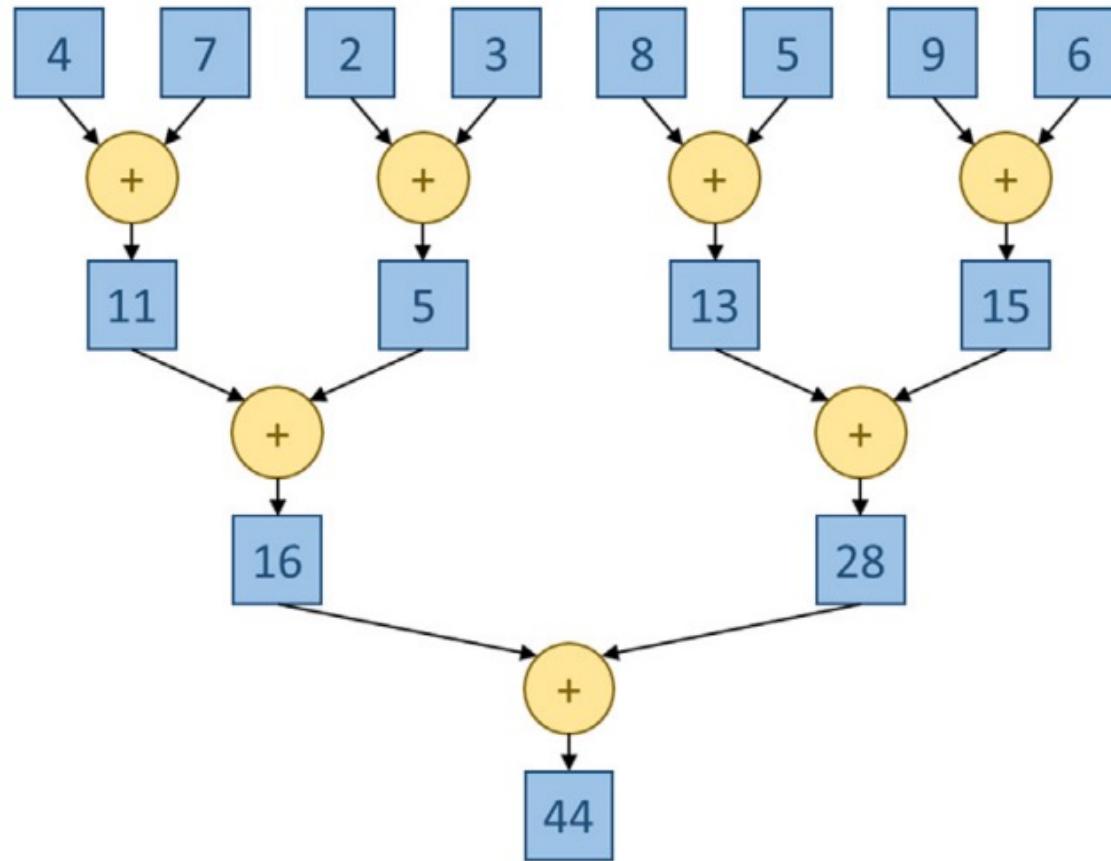


Reduction:

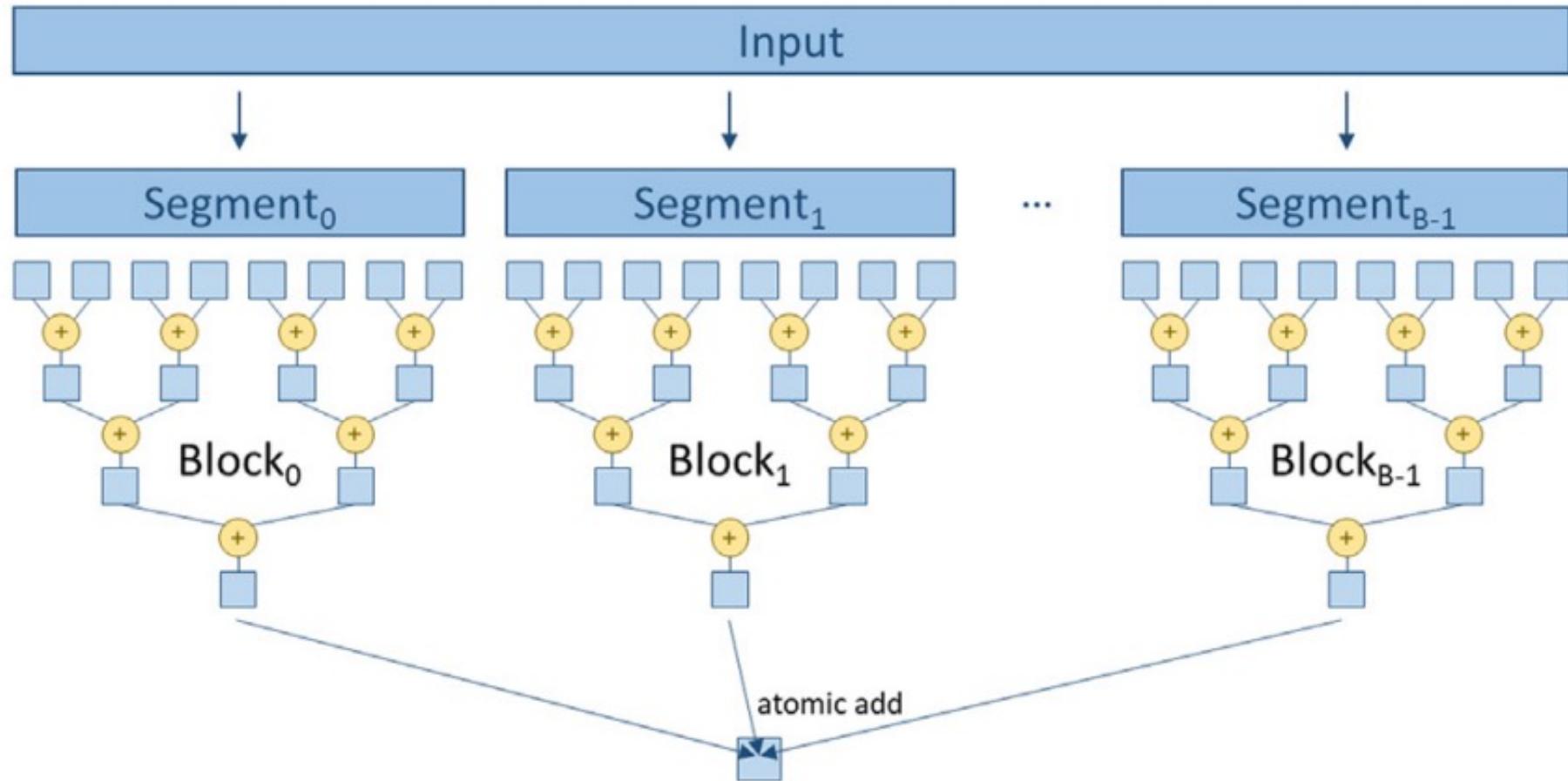
e.g.  $*c = \sum a[i]$

Thread strategy: ??

# Reduction Tree



# Reduction Tree



## naive\_reduce.py

```

1  cuda_source = """
2  __global__ void reduce_neighbor_atomic_kernel(const float* input, float* output, int N) {
3      __shared__ float sdata[128];
4
5      uint tid = threadIdx.x;
6      uint i = blockIdx.x * blockDim.x * 2 + tid;
7
8      float val = 0.0f;
9      if (i < N) val += input[i];
10     if (i + blockDim.x < N) val += input[i + blockDim.x];
11     sdata[tid] = val;
12
13     __syncthreads();
14
15     for (int stride = 1; stride < blockDim.x; stride <= 1) {
16         if (tid % (stride<<1) == 0 && tid + stride < blockDim.x)
17             sdata[tid] += sdata[tid + stride];
18         __syncthreads();
19     }
20
21     if (tid == 0) {
22         atomicAdd(output, sdata[0]);
23     }
24 }
25
26 torch::Tensor reduce_neighbor_atomic(torch::Tensor input) {
27     int N = input.size(0);
28     int threads = 128;
29     int blocks = (N + threads * 2 - 1) / (threads * 2);
30
31     auto output = torch::zeros({}, input.options());
32
33     reduce_neighbor_atomic_kernel<<<blocks, threads>>>(
34         input.data_ptr<float>(), output.data_ptr<float>(), N
35     );
36
37     return output;
38 }
39 """

```

## naive\_reduce.py

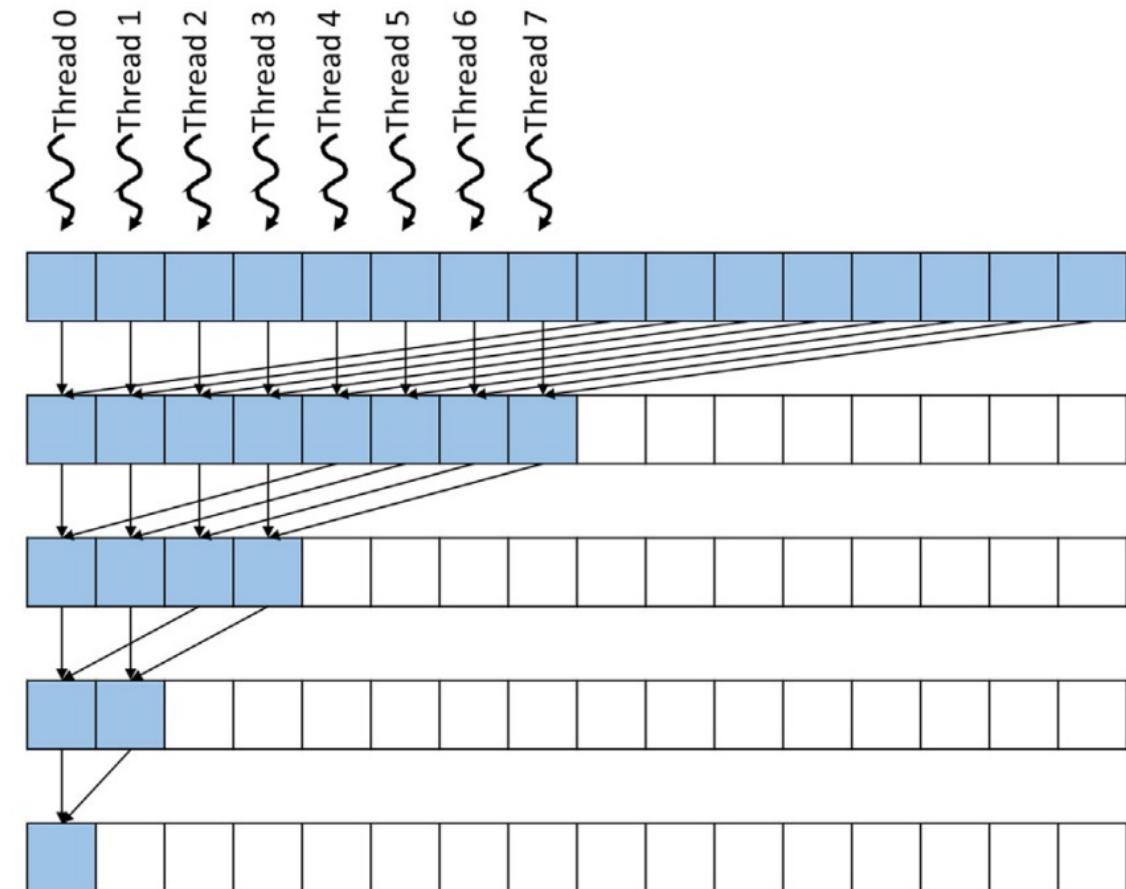
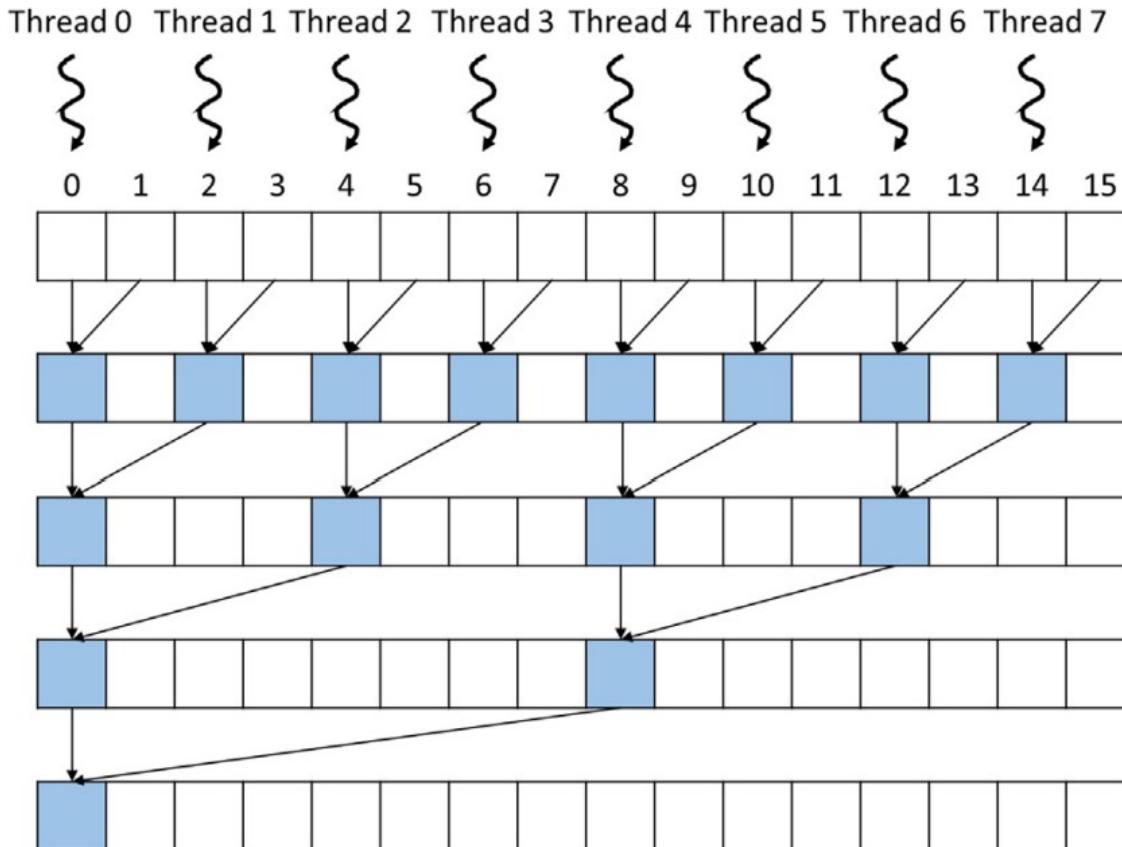
```

1  cpp_source = "torch::Tensor reduce_neighbor_atomic(torch::Tensor input);"
2
3  module = load_inline(
4      name="moss_op",
5      cpp_sources=cpp_source,
6      cuda_sources=cuda_source,
7      functions=["reduce_neighbor_atomic"],
8      verbose=False,
9  )
10
11 N = 1 << 20
12 x = torch.randn(N, device="cuda")
13
14 with profile(
15     activities=[ProfilerActivity.CUDA],
16     schedule=torch.profiler.schedule(wait=1, warmup=1, active=2),
17     on_trace_ready=torch.profiler.tensorboard_trace_handler("./log/reduce_atomic_prof"),
18     record_shapes=False,
19     with_stack=False,
20 ) as prof:
21     for i in range(4):
22         with record_function("cuda_reduce_atomic"), torch.no_grad():
23             out1 = module.reduce_neighbor_atomic(x)
24         with record_function("torch_sum"), torch.no_grad():
25             out0 = torch.sum(x)
26         prof.step()
27
28     assert torch.allclose(out0, out1)
29     print(prof.key_averages().table(sort_by="cuda_time_total", row_limit=10))

```

# Reduction Tree

- ▶ A lot of threads will be inactive :(



# Reduction



reorder\_reduce.py

```
1  cuda_source = """
2  __global__ void reduce_reorder_atomic_kernel(const float* input, float* output, int N) {
3      __shared__ float sdata[1024];
4
5      uint tid = threadIdx.x;
6      uint i = blockIdx.x * blockDim.x * 2 + tid;
7
8      float val = 0.0f;
9      if (i < N) val += input[i];
10     if (i + blockDim.x < N) val += input[i + blockDim.x];
11     sdata[tid] = val;
12
13     __syncthreads();
14
+ 15     for (int s = blockDim.x / 2; s > 0; s >>= 1) {
+ 16         if (tid < s)
+ 17             sdata[tid] += sdata[tid + s];
18         __syncthreads();
19     }
20
21     if (tid == 0) {
22         atomicAdd(output, sdata[0]);
23     }
24 }
25
26 torch::Tensor reduce_reorder_atomic(torch::Tensor input) {
27     int N = input.size(0);
28     int threads = 1024;
29     int blocks = (N + threads * 2 - 1) / (threads * 2);
30
31     auto output = torch::zeros({}, input.options());
32
33     reduce_reorder_atomic_kernel<<<blocks, threads>>>(
34         input.data_ptr<float>(), output.data_ptr<float>(), N
35     );
36
37     return output;
38 }
39 """
```

# Reduction

Name	Self CPU %	Self CPU	CPU total %	CPU total	CPU time avg	Self CUDA	Self CUDA %	CUDA total	CUDA time avg
reduce_reorder_atomic_kernel(float const*, float*, int*)	0.00%	0.000us	0.00%	0.000us	0.000us	21.536us	47.09%	21.536us	10.768us
void at::native::reduce_kernel<512, 1, at::native::ReduceType::Sum>(Tensor const*, Tensor*, int)	0.00%	0.000us	0.00%	0.000us	0.000us	18.368us	40.17%	18.368us	9.184us
void at::native::vectorized_elementwise_kernel<4, at::native::ElementWiseOpType::Add>(Tensor const*, Tensor const*, Tensor*, int)	0.00%	0.000us	0.00%	0.000us	0.000us	3.265us	7.14%	3.265us	1.632us
Memset (Device)	0.00%	0.000us	0.00%	0.000us	0.000us	2.560us	5.60%	2.560us	1.280us
cudaLaunchKernel	71.60%	83.565us	71.60%	83.565us	13.928us	0.000us	0.00%	0.000us	0.000us
cudaMemsetAsync	18.79%	21.930us	18.79%	21.930us	10.965us	0.000us	0.00%	0.000us	0.000us
cudaDeviceSynchronize	9.61%	11.216us	9.61%	11.216us	11.216us	0.000us	0.00%	0.000us	0.000us

# Reduction

## 作业二： + Vectorized



reorder\_reduce.py

```
1  cuda_source = """
2  __global__ void reduce_reorder_atomic_kernel(const float* input, float* output, int N) {
3      __shared__ float sdata[1024];
4
5      uint tid = threadIdx.x;
6      uint i = blockIdx.x * blockDim.x * 2 + tid;
7
8      float val = 0.0f;
9      if (i < N) val += input[i];
10     if (i + blockDim.x < N) val += input[i + blockDim.x];
11     sdata[tid] = val;
12
13     __syncthreads();
14
+ 15     for (int s = blockDim.x / 2; s > 0; s >>= 1) {
+ 16         if (tid < s)
+ 17             sdata[tid] += sdata[tid + s];
18         __syncthreads();
19     }
20
21     if (tid == 0) {
22         atomicAdd(output, sdata[0]);
23     }
24 }
25
26 torch::Tensor reduce_reorder_atomic(torch::Tensor input) {
27     int N = input.size(0);
28     int threads = 1024;
29     int blocks = (N + threads * 2 - 1) / (threads * 2);
30
31     auto output = torch::zeros({}, input.options());
32
33     reduce_reorder_atomic_kernel<<<blocks, threads>>>(
34         input.data_ptr<float>(), output.data_ptr<float>(), N
35     );
36
37     return output;
38 }
39 """
```

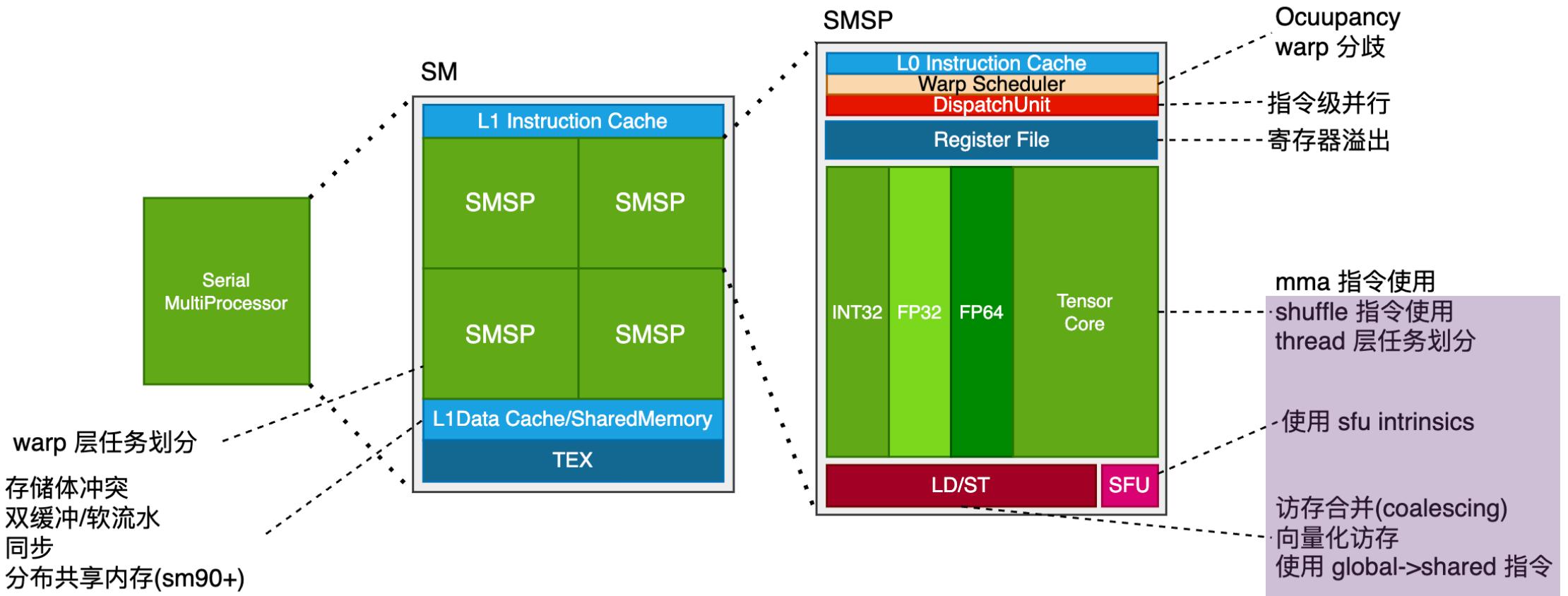
# Shuffle指令



4\_shuffle\_reduce.py

```
1  cuda_source = """
2  __inline__ __device__ float warp_reduce_sum(float val) {
3      for (int8_t offset = 16; offset > 0; offset >>= 1)
4          val += __shfl_down_sync(0xffffffff, val, offset);
5      return val;
6  }
7
8  __global__ void reduce_shuffle_atomic_kernel(const float* input, float* output, int N) {
9      __shared__ float shared[32];
10
11     float val = 0.0f;
12     int idx = blockIdx.x * blockDim.x * 2 + threadIdx.x;
13
14     if (idx < N) val += input[idx];
15     if (idx + blockDim.x < N) val += input[idx + blockDim.x];
16
17     val = warp_reduce_sum(val);
18
19     if ((threadIdx.x & 31) == 0)
20         shared[threadIdx.x >> 5] = val;
21
22     __syncthreads();
23
24     float block_sum = 0.0f;
25     if (threadIdx.x < 32) {
26         block_sum = shared[threadIdx.x];
27         block_sum = warp_reduce_sum(block_sum);
28     }
29
30     if (threadIdx.x == 0)
31         atomicAdd(output, block_sum);
32 }
```

# Conclusion



# Research Proposal

- ▶ Efficient infrastructure via algorithm-hardware co-design
- ▶ Language-centric multimodal in-context learning

**FlexMLA**

Optimize + Test

**SlimFFN**

Optimize + Test

**FoX**

Reproduction

[taoji@fudan.edu.cn](mailto:taoji@fudan.edu.cn)  
A5029, 江湾叉二

# Q&A