

Triton, Gluon, and the Future of Tile-Based Programming Models

Keren Zhou
kzhou6@gmu.edu

Evolution of AI Applications



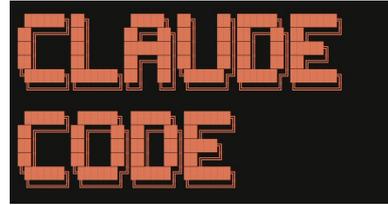
10 Years Ago



3 Years Ago



2 Years Ago



Last Year



Last Month



Past

Now

AI Software Stack



Outline

- Parallel Systems
- Triton
- Gluon and Layouts
- Rethink Tile-based Programming Models
- Conclusions

Parallel Systems

Overview

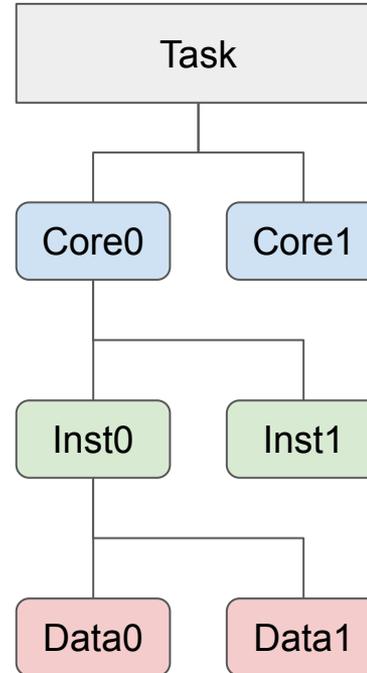
- Multi-core CPUs
- GPUs
- Accelerators

CPU Architectures

- x86 (Intel, AMD)
- ARM (Advanced RISC Machine – common in mobile and embedded)
- RISC-V (Open-source RISC architecture, rising in adoption)
- PowerPC (IBM, used in older Macs, embedded systems)
- MIPS (Used in routers, embedded systems)

Parallelisms on CPUs

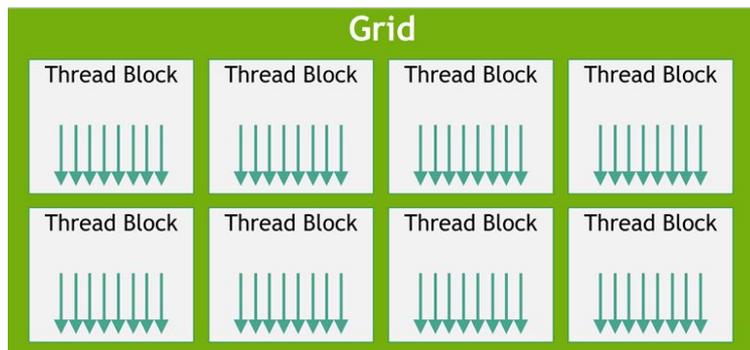
- Thread parallelism
- Instruction parallelism
- SIMD parallelism



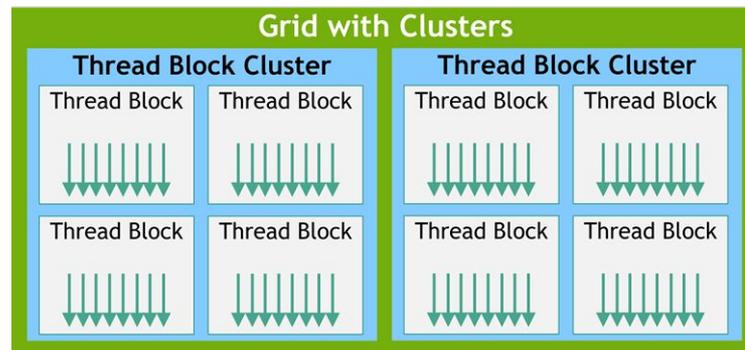
GPUs

- Grid (kernel) parallelism
- Thread block parallelism
- Thread block cluster parallelism
- Warp parallelism
- Thread parallelism
- Instruction parallelism
- SIMD parallelism

NVIDIA GPUs



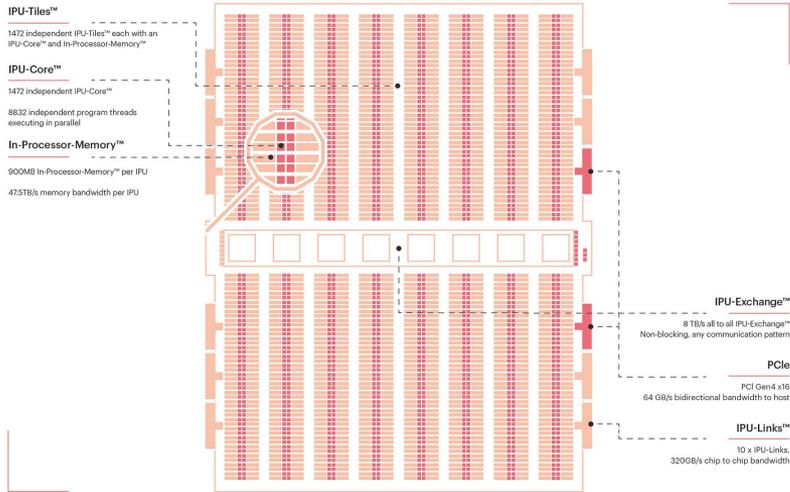
Before Hopper



Since Hopper

Accelerators

Different accelerators have different abstractions



Graphcore

Dataflow Architecture for Terabyte Sized Models



DataScale SN10-8R
1/4 Rack System

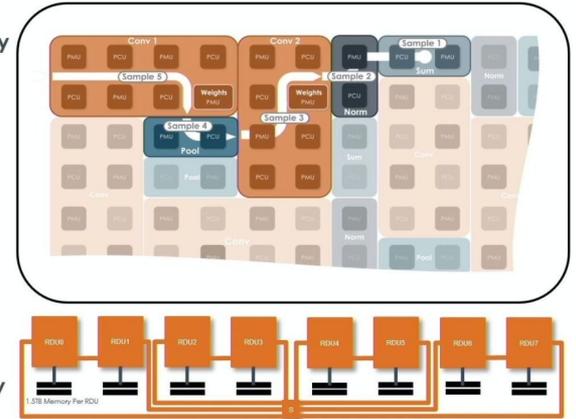
Dataflow Efficiency

+

Compute
Capability

+

Large
Memory Capacity



Sambanova

Thread-Based Programming

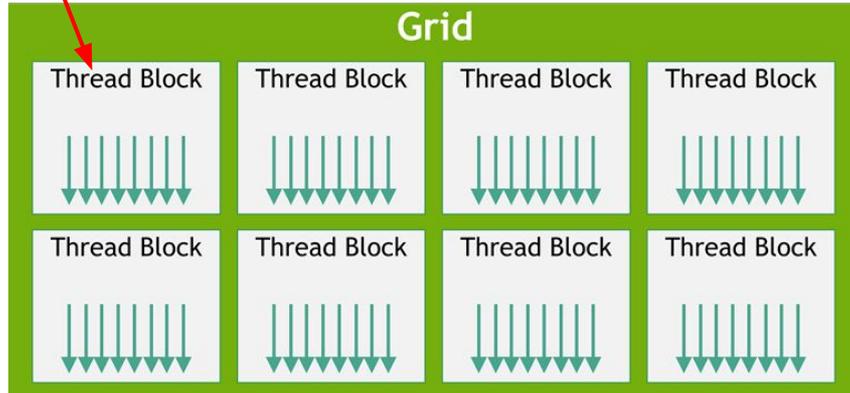
- Users specify the behavior of each thread
- The number of blocks and threads within each block is controlled on the host side

Behavior of Each Thread

```
vecAdd

__global__ void vectorAdd(const float *A, const float *B, float *C, int numElements) {
    int i = blockDim.x * blockIdx.x + threadIdx.x;
    if (i < numElements) {
        C[i] = A[i] + B[i];
    }
}
```

snappify.com

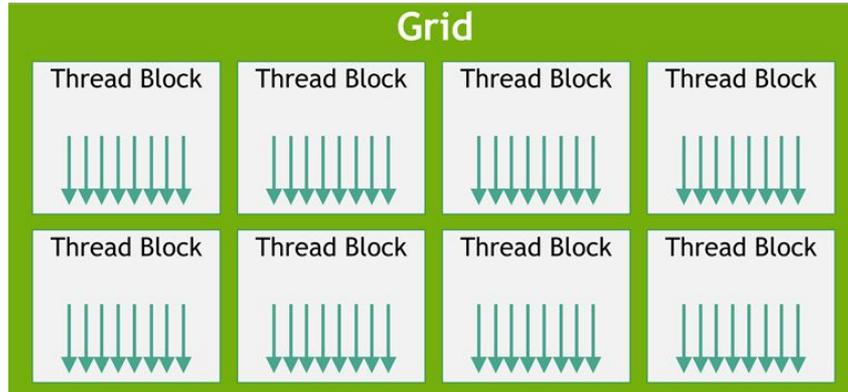


Number of Blocks and Threads

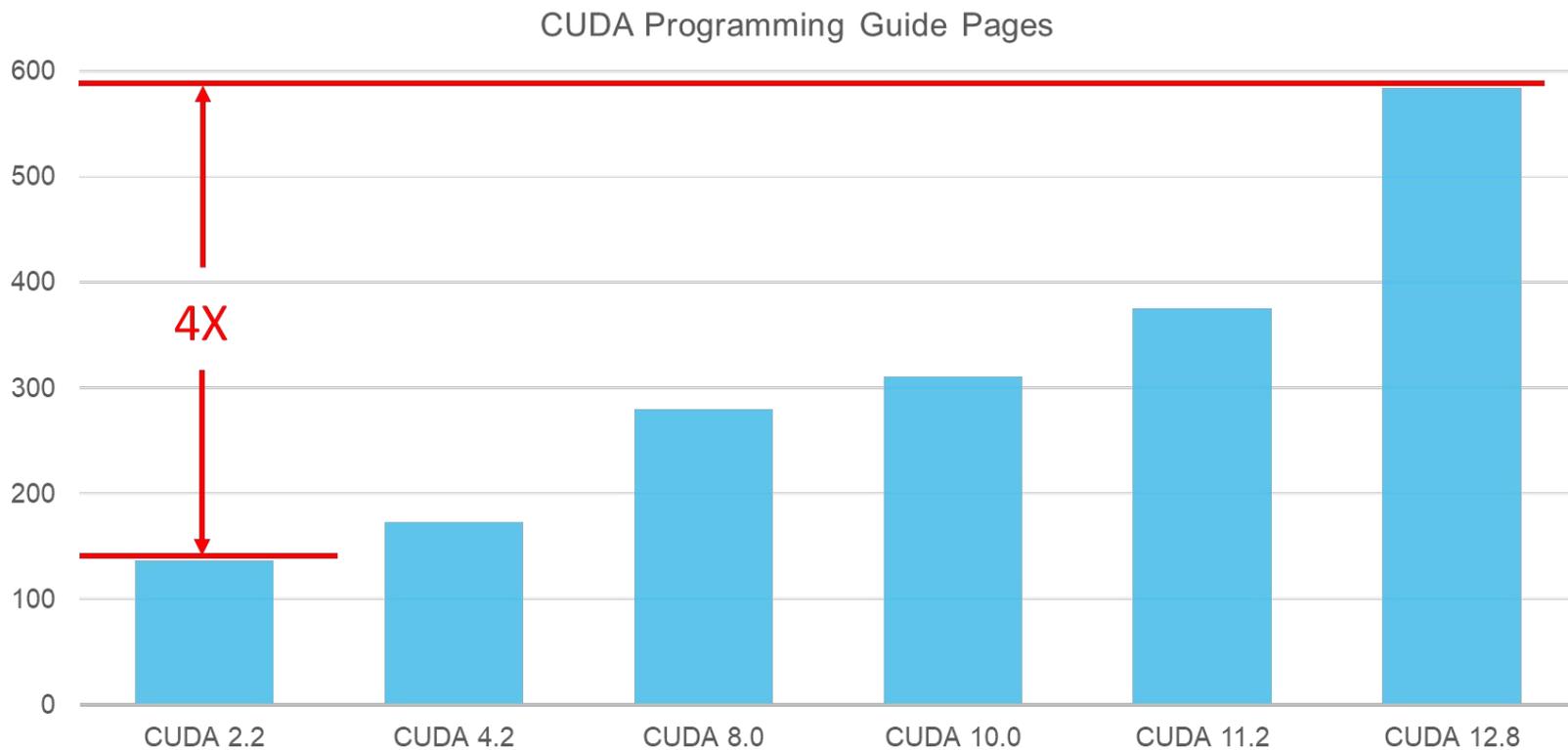
```
vecAdd

int threadsPerBlock = 256;
int blocksPerGrid = (numElements + threadsPerBlock - 1) / threadsPerBlock;
vectorAdd<<<blocksPerGrid, threadsPerBlock>>>(d_A, d_B, d_C, numElements);
```

snappify.com



CUDA Programming Guide



The “Problems” with Full CUDA Specification

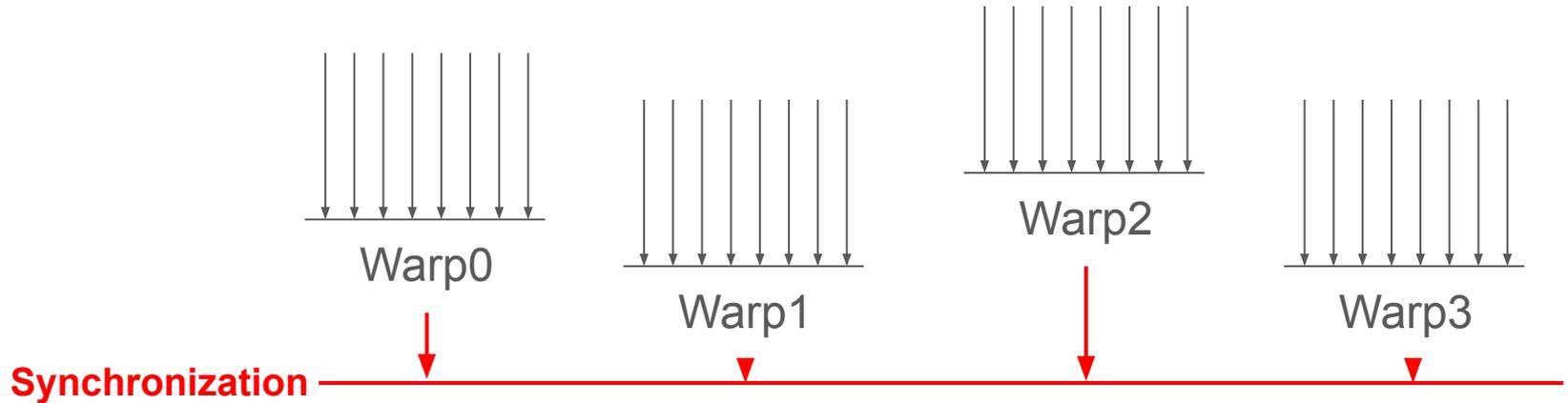
- Unnecessary to get high performance in common scenarios
- Increasingly complex compute units/memory hierarchy
- Sophisticated user interfaces for tensor core instructions
- Hidden/unspecified behaviors

How many of you are familiar with “Programmatic
Dependent Launch”?

Triton

Motivation

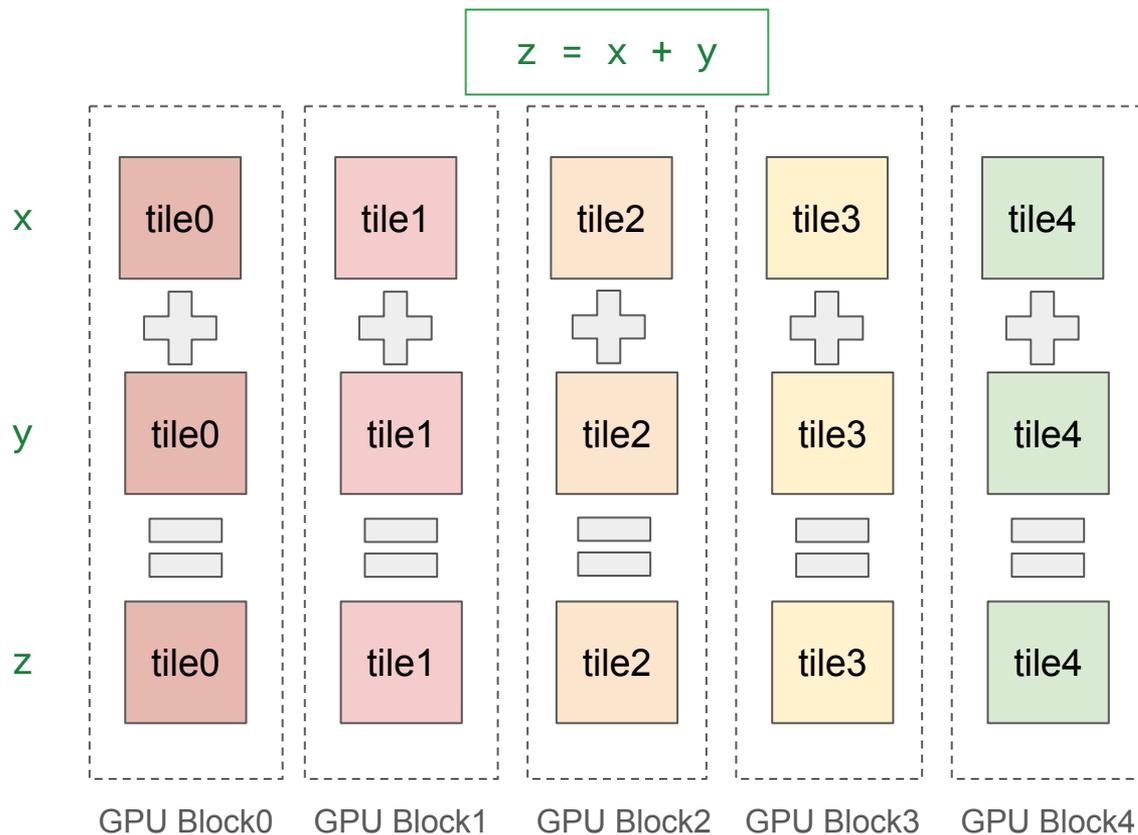
- Threads within a warp are executed in lock steps
 - This is the behavior of many simulators
 - Though not accurate with “independent thread scheduling”
- Warps within a thread block access the same programmable shared memory
- We can coarsen the execution unit as using warps or thread blocks



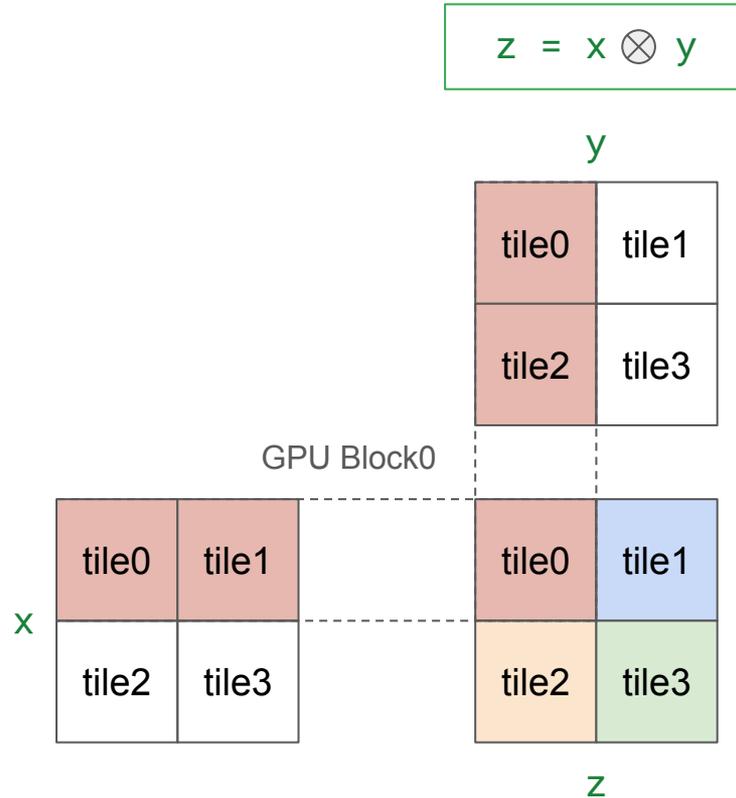
Core Concepts - Tiling/Blocking

- A tile is a block of data elements such as a submatrix or subarray processed collectively by a couple of threads
- In the context of GPU, tiles are utilized to load chunks of data into shared memory or registers
 - Coalesced global memory accesses
 - Reduced bank conflicts
 - Matching tensor core layouts
 - Promoting cache utilizing

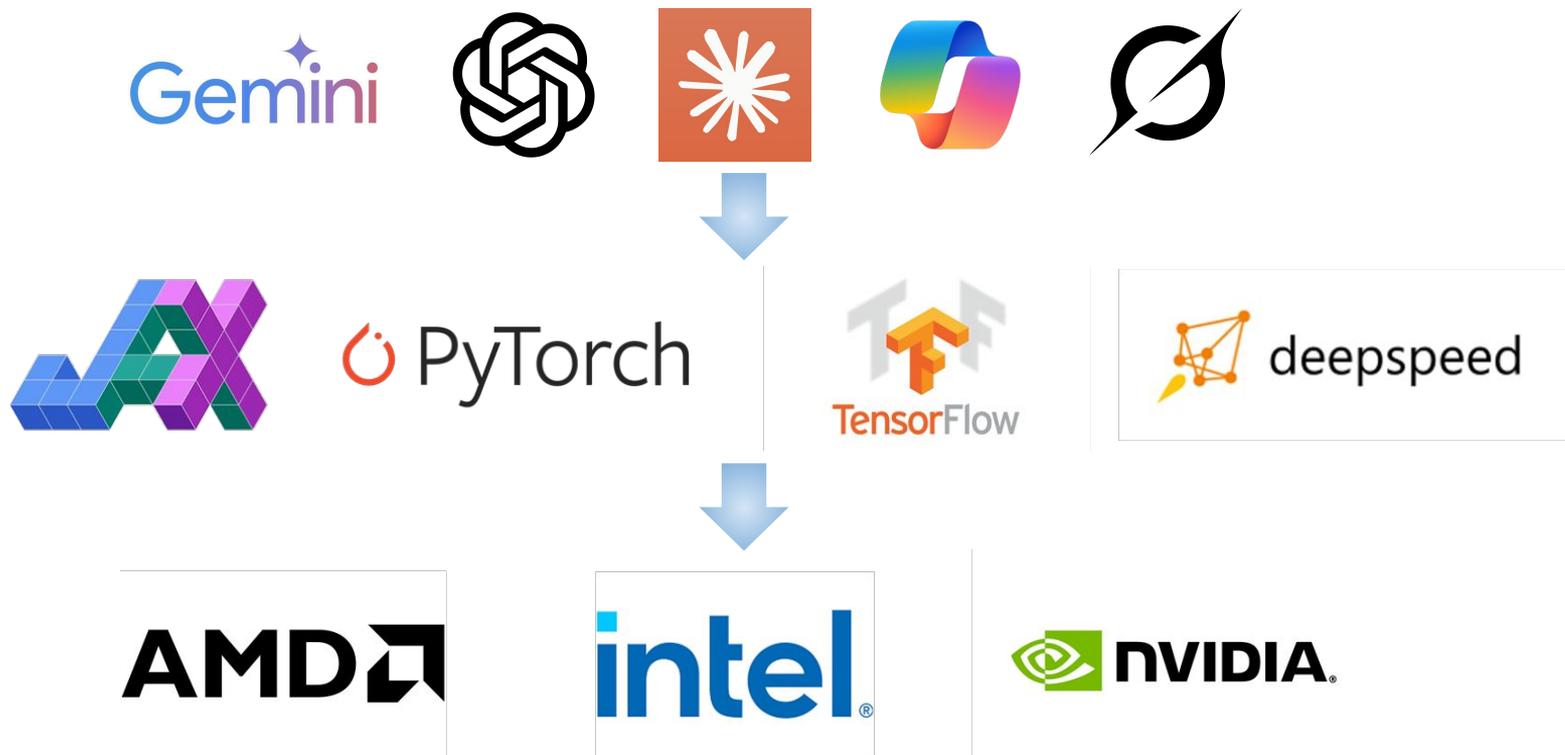
Example - Vector Addition



Example - Matrix Multiplication



AI System Software Stack



Why Triton?

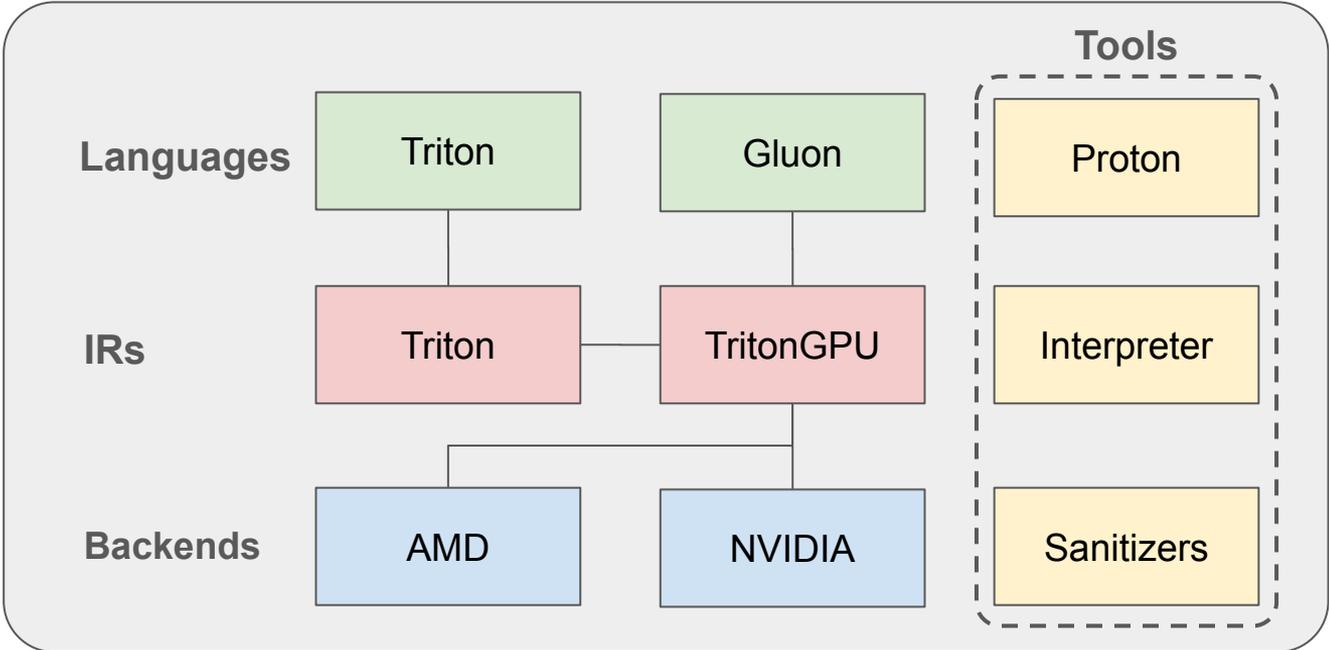


deepspeed

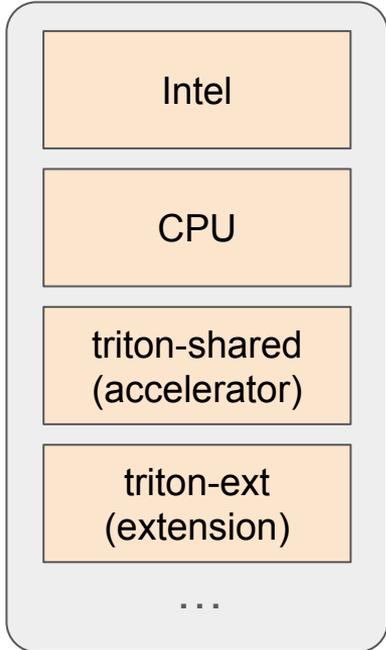


Ecosystem

In-tree Modules



Out-of-tree Modules



Triton Language

- Python-like language designed for high flexibility and performance in deep learning applications
 - Support tensor interface similar to PyTorch
 - Uses Python-like syntax
- Compared to CUDA/ROCm, Triton simplifies GPU programming
 - Only requiring knowledge that a kernel is divided into multiple blocks (Triton programs)
 - Most underlying details are handled by the compiler

Triton vs CUDA

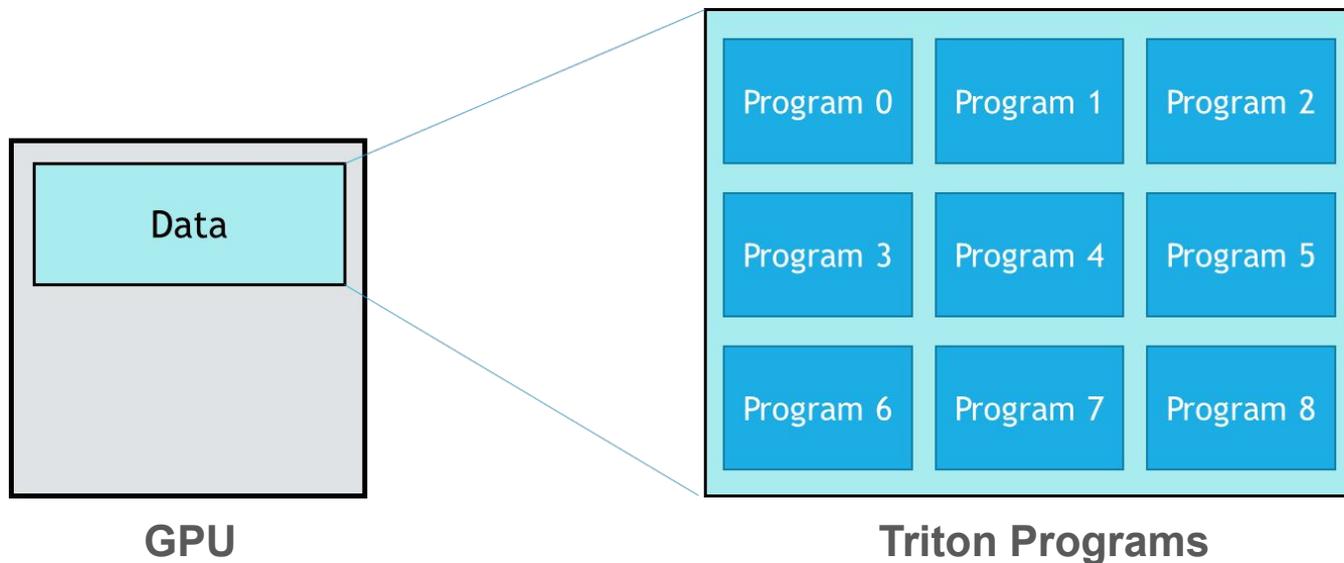
	CUDA	Triton
Memory	Global/Shared/Local	Automatic
Parallelism	Threads/Blocks/Warps	Mostly Blocks*
Tensor Core	Manual	Automatic
Vectorization	.8/.16/.32/.64/.128	Automatic
Async SIMT	Support	Limited
Device Function	Support	Support

Using Triton, you only need to know that a program is divided into multiple blocks

*In Triton 3.3, we added warp specialization so technically not all warps issue the same code

SPMD Model

- Each program executes the “same” code
 - Only program ids are different



A Simple Triton Program - Kernel Code

```
z: dim0 x dim1 = x: dim0 x dim1 + y: dim0 x dim1
```

```
Kernel decorator — 1 @triton.jit
2 def add_kernel(x_ptr, y_ptr, z_ptr, dim0, dim1,
3               BLOCK_DIM0: tl.constexpr, BLOCK_DIM1: tl.constexpr):
Programming model — 4     pid_x = tl.program_id(axis=0)
5     pid_y = tl.program_id(axis=1)
6     block_start = pid_x * BLOCK_DIM0 * dim1 + pid_y * BLOCK_DIM1
Creation ops — 7     offsets_dim0 = tl.arange(0, BLOCK_DIM0)[:,None]
8     offsets_dim1 = tl.arange(0, BLOCK_DIM1)[None, :]
9     offsets = block_start + offsets_dim0 * dim1 + offsets_dim1
10    masks = (offsets_dim0 < dim0) & (offsets_dim1 < dim1)
Memory ops — 11    x = tl.load(x_ptr + offsets, mask=masks)
12    y = tl.load(y_ptr + offsets, mask=masks)
13    output = x + y
14    tl.store(z_ptr + offsets, output, mask=masks)
```

A Simple Triton Program - Kernel Launch

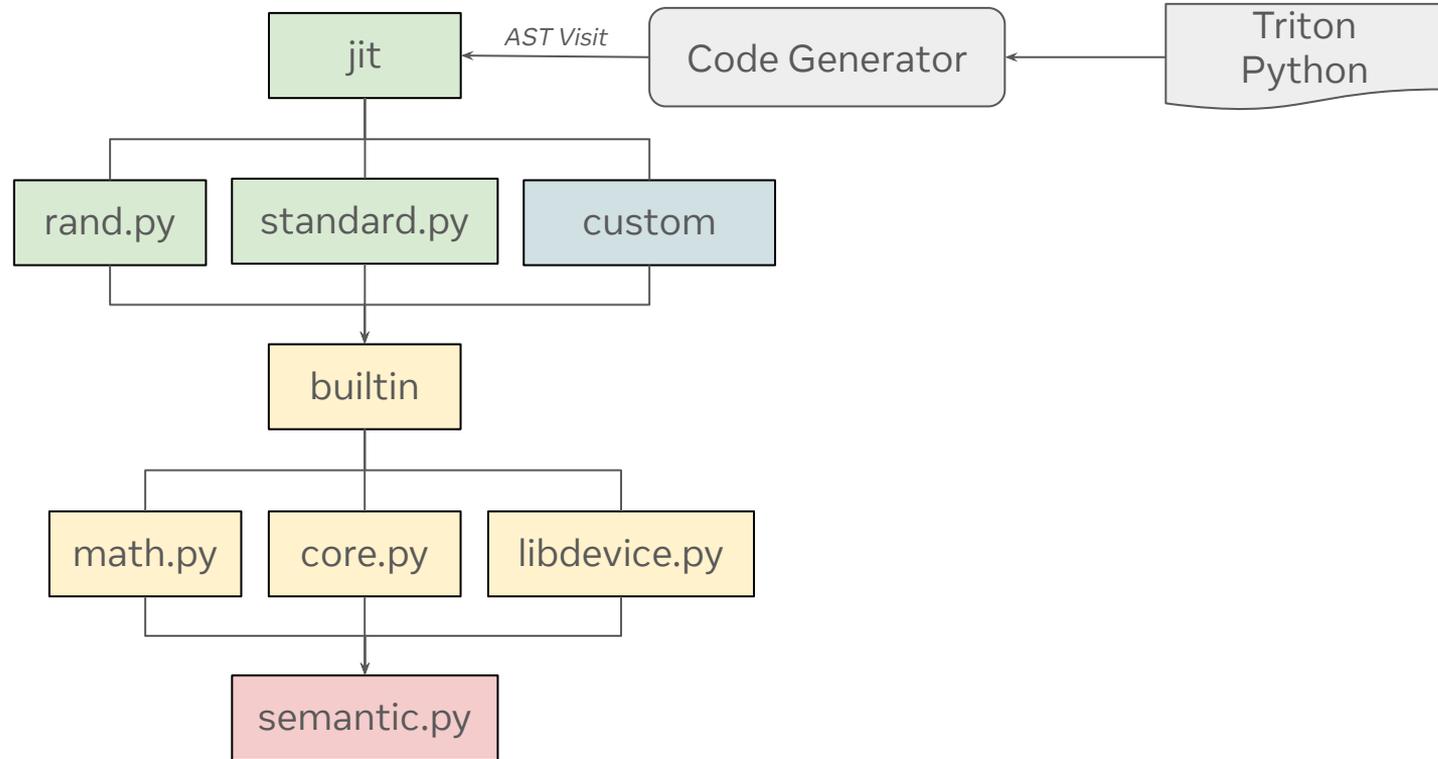
```
z: dim0 x dim1 = x: dim0 x dim1 + y: dim0 x dim1
```



Host Code

```
1 x = torch.randn((4096,1), device="cuda")
2 y = torch.randn((4096,1), device="cuda")
3 z = torch.empty((4096,1), device="cuda")
4 programs = 4096 // 128
5 add_kernel[(programs,)](x, y, z, 4096, 1, 128, 1)
```

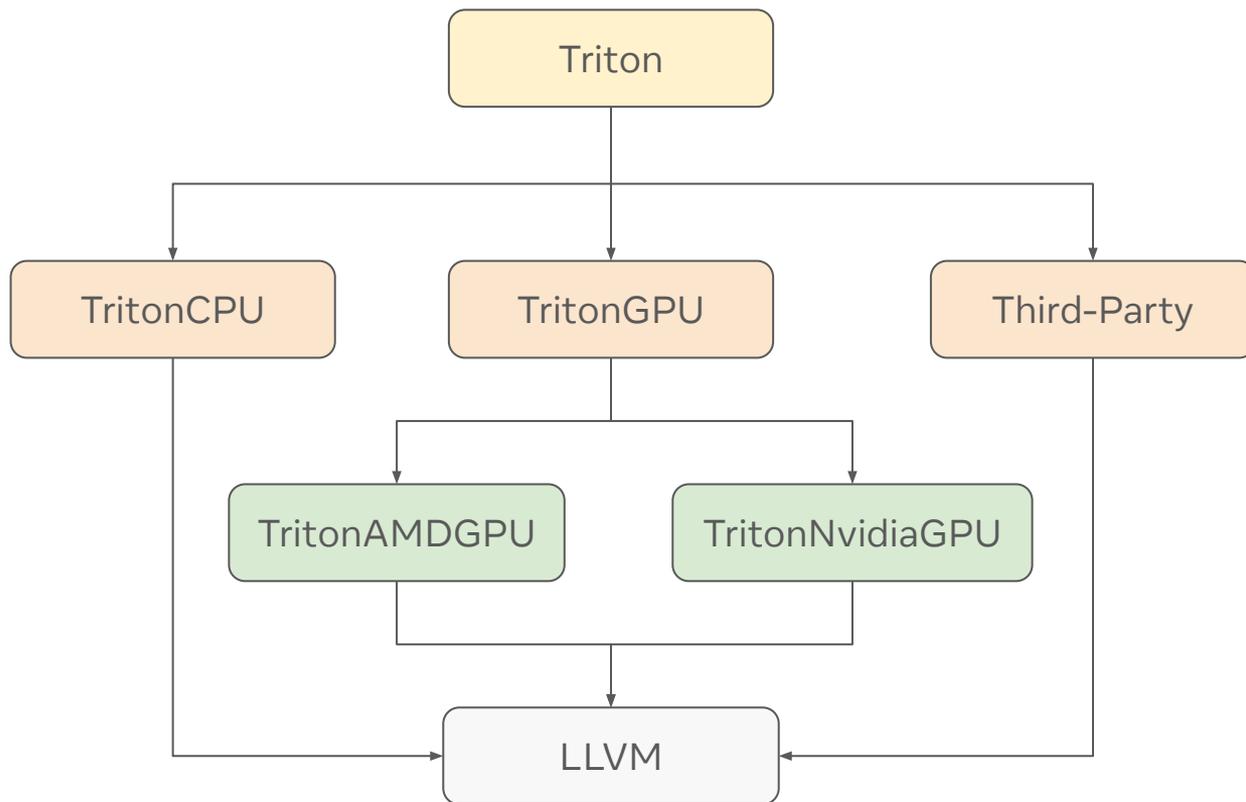
Triton Frontend



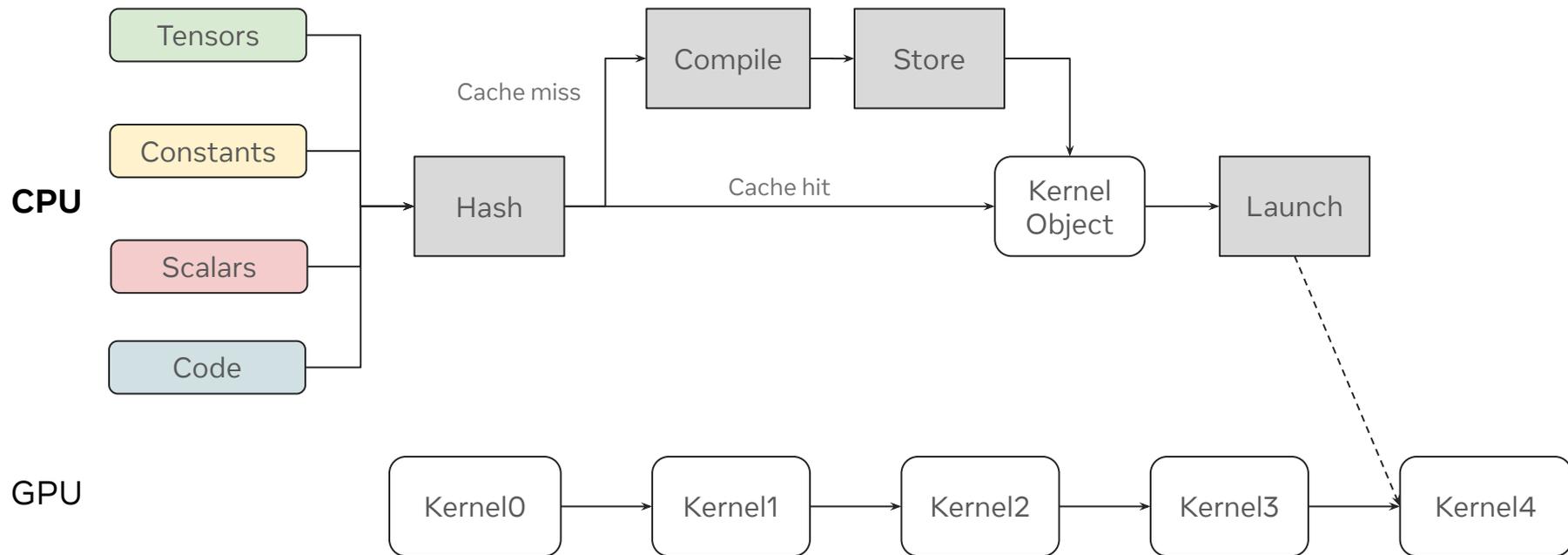
MLIR: Multi-Level IR Compiler Framework

- Building reusable and extensible compiler infrastructure
 - Address software fragmentation
 - Improve compilation for heterogeneous hardware
 - Reduce the cost of building domain specific compilers
- Key concept: Dialects
 - Each dialect is given a unique namespace that is prefixed to each defined attribute/operation/type
 - For example, the *Affine* dialect defines the namespace: `affine`
 - MLIR allows for multiple dialects exist in the same IR
 - A dialect can be converted to another under conversion rules

Triton Backends



Step-by-Step Compilation: JIT Compilation



Step-by-Step Compilation: TritonIR (TTIR)

```
tt.func                                Mangled function name
@matmul_kernel__Pfp32_Pfp32_Pfp32_i32_i32_i32_i32_i32_i32_i32_i32_i32__12c64_13c64_14c64_15c8
(%arg0: !tt.ptr<f32> {tt.divisibility = 16 : i32}, ...) {
  %cst = arith.constant dense<true> : tensor<64x64xi1> ← Tensor
  %c64 = arith.constant 64 : i32 ← Scalar
  %c0 = arith.constant 0 : i32
  %0 = tt.get_program_id * : i32 ← Triton operation
  %1 = arith.addi %arg3, %c63_i32 : i32 ← Arith operation
  %2 = arith.divsi %1, %c64_i32 : i32
  %3 = arith.addi %arg4, %c63_i32 : i32
  %4 = arith.divsi %3, %c64_i32 : i32
  %5 = arith.muli %4, %c8_i32 : i32
  %6 = arith.divsi %0, %5 : i32
  %7 = arith.muli %6, %c8_i32 : i32
```

Step-by-Step Compilation: TritonGPU IR (TTGIR)

Specialized “layouts”

```
#blocked = #ttg.blocked<{sizePerThread = [8, 1], threadsPerWarp = [8, 4], warpsPerCTA = [1, 4], order = [0, 1]}>
#blocked1 = #ttg.blocked<{sizePerThread = [1, 8], threadsPerWarp = [4, 8], warpsPerCTA = [4, 1], order = [1, 0]}>
#mma = #ttg.nvidia_mma<{versionMajor = 2, versionMinor = 0, warpsPerCTA = [4, 1], instrShape = [16, 8]}>
module attributes {"ttg.num-ctas" = 1 : i32, "ttg.num-warps" = 4 : i32} {
// CHECK-LABEL: tt.func @load_two_users
  tt.func @load_two_users(%arg0: !tt.ptr<f16> {tt.divisibility = 16 : i32}, %arg1: !tt.ptr<f16> {tt.divisibility = 16 : i32})
-> (tensor<128x16xf32, #mma>, tensor<128x64xf32, #mma>) {
  %cst = arith.constant dense<0> : tensor<1x16xi32, #blocked>
  %cst_0 = arith.constant dense<0> : tensor<128x1xi32, #blocked1>
  %c0_i64 = arith.constant 0 : i64
  %c0_i32 = arith.constant 0 : i32
  %cst_1 = arith.constant dense<0.000000e+00> : tensor<128x16xf32, #mma>
  %cst_2 = arith.constant dense<0.000000e+00> : tensor<128x64xf32, #mma>
```

Example Layouts - MMA

R\C	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	T0:{a0,a1}	T1:{a0,a1}	T2:{a0,a1}	T3:{a0,a1}	T0:{a4,a5}	T1:{a4,a5}	T2:{a4,a5}	T3:{a4,a5}								
1	T4:{a0,a1}	T5:{a0,a1}	T6:{a0,a1}	T7:{a0,a1}	T4:{a4,a5}	T5:{a4,a5}	T6:{a4,a5}	T7:{a4,a5}								
2	→				→											
..	←				←											
7	T28:{a0,a1}	T29:{a0,a1}	T30:{a0,a1}	T31:{a0,a1}	T28:{a4,a5}	T29:{a4,a5}	T30:{a4,a5}	T31:{a4,a5}								
8	T0:{a2,a3}	T1:{a2,a3}	T2:{a2,a3}	T3:{a2,a3}	T0:{a6,a7}	T1:{a6,a7}	T2:{a6,a7}	T3:{a6,a7}								
9	T4:{a2,a3}	T5:{a2,a3}	T6:{a2,a3}	T7:{a2,a3}	T4:{a6,a7}	T5:{a6,a7}	T6:{a6,a7}	T7:{a6,a7}								
10	→				→											
..	←				←											
15	T28:{a2,a3}	T29:{a2,a3}	T30:{a2,a3}	T31:{a2,a3}	T28:{a6,a7}	T29:{a6,a7}	T30:{a6,a7}	T31:{a6,a7}								

%laneid:{fragments}

Key Transformation Passes

- Remove layout conversion
 - Rewrite the `ConvertLayoutOps` to reduce the number of operations and to prefer favorable layouts like `BlockedEncodingAttr` layout for "expensive" loads and stores (good for coalescing) and `NvidiaMmaEncodingAttr` otherwise (good for tensor ops)
- Accelerate matmul
 - Optimize the input/output layout of `dot` instructions to make them compatible hardware accelerators
- Automatic warp specialization
 - Analyze the loops in the kernel and attempt to create a partition schedule so that different warps handles different code regions
- Pipeline
 - Apply software pipelining to loops in the module based on number of stages. This may convert some load into asynchronous loads, and multi-buffer the data.

More info can be found in:
[triton/Dialect/TritonGPU/Transforms/Passes.td](#)

Gluon and Layouts

Problems with Triton

- Performance
 - Triton uses a synchronous, block-level abstraction
 - Modern GPU architectures are warp-specialized and asynchronous
- Maintenance
 - Divergence across GPU architectures from different vendors
 - Makes compiler backends complex and difficult to maintain

Gluon Semantics

- Tailored for GPU architectures
 - Triton was designed for any accelerators/CPUs, though not all backends have implemented
- Each GPU architecture provides its own set of high-level operations
- Allows access to architecture-specific hardware intrinsics
- Uses very lightweight GPU conversion passes
- Users must manually assign a layout to each tile

Gluon Language

- Data types
 - Basic types
 - inherited from Triton
 - Additional types
 - `shared_memory_descriptor_type`
 - `tensor_memory_descriptor_type`
 - `distributed_type`
- Operations
 - Basic operations
 - Inherited from Triton
 - Shared memory operations
 - `shared_load`, `shared_store`, `shared_gather`, `shared_scatter`
 - Execution
 - `warp_specialize`
- ISA
 - `hopper.tma`, `hopper.cluster`, `hopper.mbarrier`
 - `blackwell.clc`, `blackwell.tma`, `blackwell.float2`
 - `ampere.mbarrier`, `ampere.async_copy`

Matmul Example

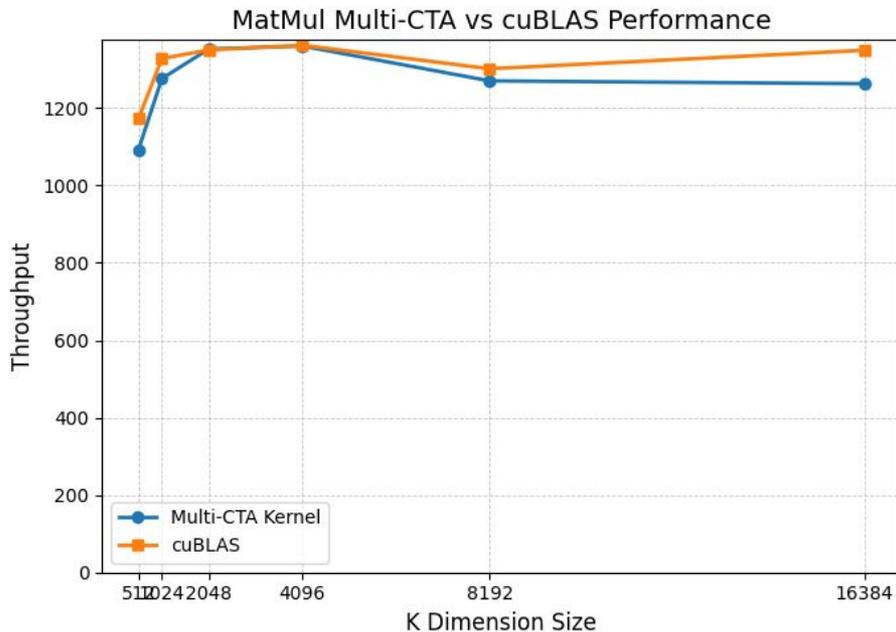
- Partition workloads into 7 warps
 - epilogue partition (4)
 - Load partition (1)
 - mma partition (1)
 - clc partition (1)
- Enable the 2-cta mode to reduce shared memory usage
- Enable cluster launch control to accelerate work-stealing of persistent processing

```
1  gl.warp_specialize([
2      (matmul_epilogue_partition, (p, )),
3      (matmul_load_partition, (p, )),
4      (matmul_mma_partition, (p, )),
5      (matmul_clc_partition, (p, )),
6  ], [1, 1, 1], [24, 24, 24])
```

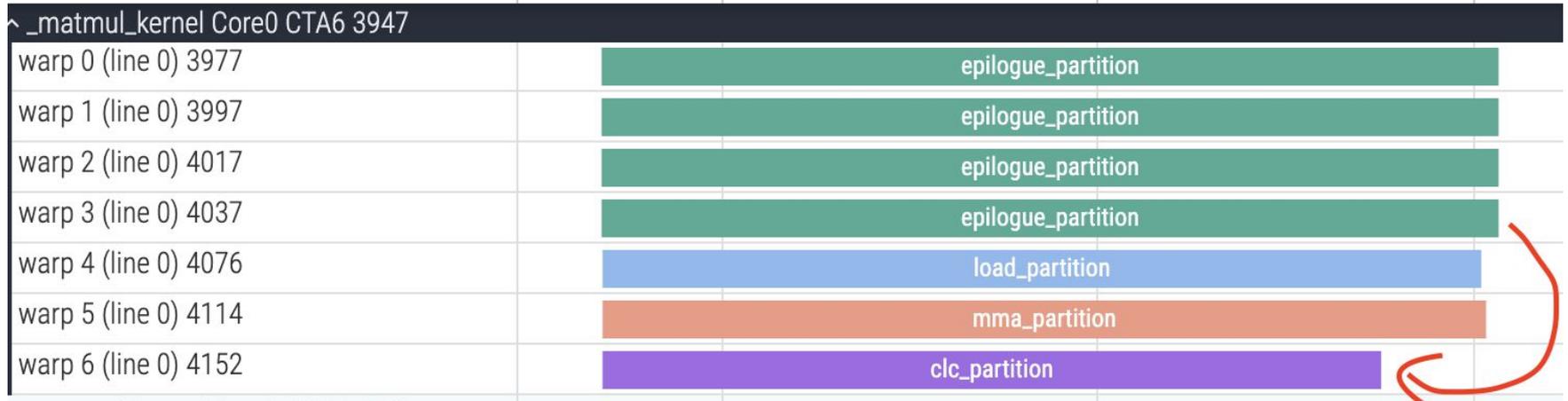
```
1  acc_layout: gl.constexpr = TensorMemoryLayout(
2      block=(cta_m, tile_n),
3      col_stride=1,
4      cga_layout=cga_layout,
5      two_ctas=True,
6  )
```

Matmul Performance

- Our warp specialized + 2 CTA matmul achieves similar performance with cuBLAS without excessive autotuning



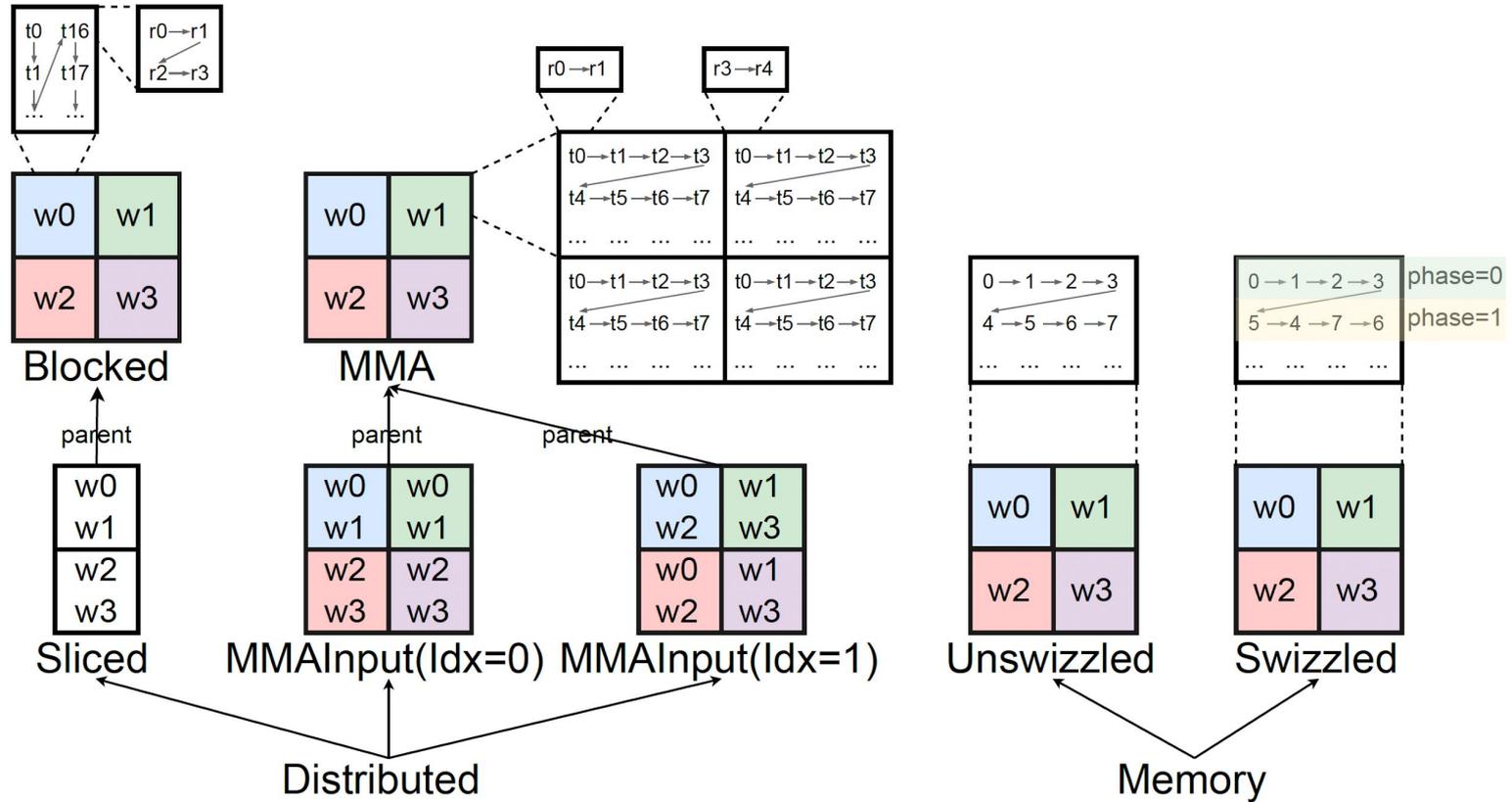
Tracing



Tensor Layouts

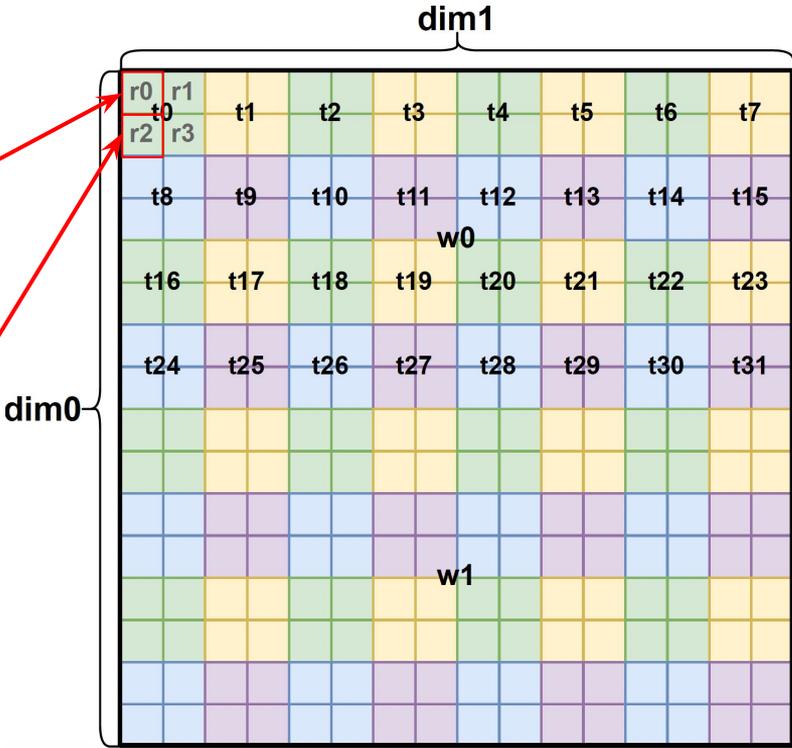
- Mapping between logic coordinates and hardware resources
- Coordinates: $x_0, x_1, \dots, x_{(n-1)}$ for a N-dimensional tile/tensor
- Execution units: thread, warp, and CTA indices
- Memory: shared memory, tensor memory

Layout Family

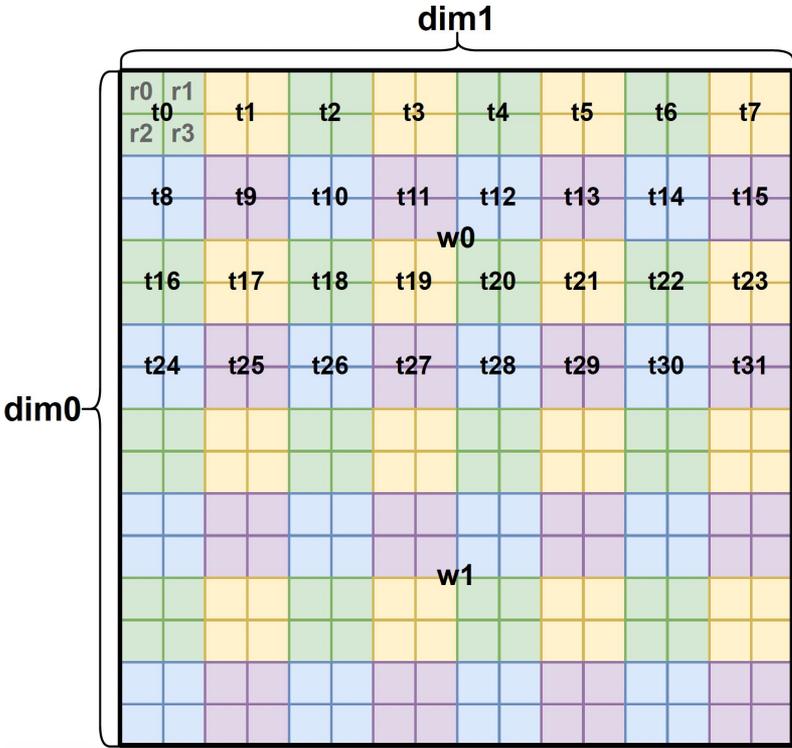
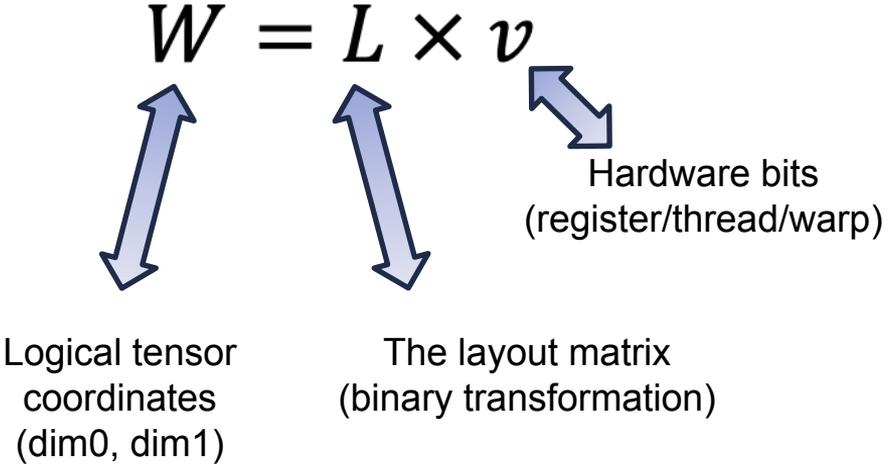


Linear Layouts

Location	Register	Thread	Warp
(0, 0) / (0b0, 0b0)	r_0 / 0b0	t_0 / 0b0	w_0 / 0b0
(0, 1) / (0b0, 0b1)	r_1 / 0b1	t_0 / 0b0	w_0 / 0b0
(0, 2) / (0b0, 0b10)	r_0 / 0b0	t_1 / 0b1	w_0 / 0b0
(0, 3) / (0b0, 0b11)	r_1 / 0b1	t_1 / 0b1	w_0 / 0b0
...
(1, 0) / (0b1, 0b0)	r_2 / 0b10	t_0 / 0b0	w_0 / 0b0
(1, 1) / (0b1, 0b1)	r_3 / 0b11	t_0 / 0b0	w_0 / 0b0
...
(2, 2) / (0b10, 0b10)	r_0 / 0b0	t_9 / 0b1001	w_0 / 0b0
(2, 3) / (0b10, 0b11)	r_1 / 0b1	t_9 / 0b1001	w_0 / 0b0
...
(3, 2) / (0b11, 0b10)	r_2 / 0b10	t_9 / 0b1001	w_0 / 0b0
(3, 3) / (0b11, 0b11)	r_3 / 0b11	t_9 / 0b1001	w_0 / 0b0
...



Defining Linear Layouts



Modeling Swizzling

- `swizzled_col = col XOR constant`
 - If the constant is less than `#columns`, the swizzled column indices will still be unique and less than `#columns`

Tile WITHOUT Swizzle			Tile WITH Swizzle					
	0	1	2	3	4	5	6	7
Row 0	0	1	2	3	4	5	6	7
Row 1	0	1	2	3	4	5	6	7
Row 2	0	1	2	3	4	5	6	7
Row 3	0	1	2	3	4	5	6	7
Row 4	0	1	2	3	4	5	6	7
Row 5	0	1	2	3	4	5	6	7
Row 6	0	1	2	3	4	5	6	7
Row 7	0	1	2	3	4	5	6	7

→

	0	1	2	3	4	5	6	7
Row 0	0	1	2	3	4	5	6	7
Row 1	1	0	3	2	5	4	7	6
Row 2	2	3	0	1	6	7	4	5
Row 3	3	2	1	0	7	6	5	4
Row 4	4	5	6	7	0	1	2	3
Row 5	5	4	7	6	1	0	3	2
Row 6	6	7	4	5	2	3	0	1
Row 7	7	6	5	4	3	2	1	0

$$\begin{bmatrix} I_n & C \\ 0 & I_m \end{bmatrix}.$$

where I_m and I_n denote identity matrices of size m and n accordingly. Each row c_i in C is given by

$$c_i = (\text{vec} \cdot (\frac{2^i}{\text{per_phase}} \bmod \text{max_phase})) \bmod 2^n.$$

Use Cases

- Software emulation of MXFP types
 - Linear layouts-based layout conversion
- Tile-based load/store (e.g., `ldmatrix/stmatrix`)
 - Linear layout division operations
- Contiguous elements
 - Largest u such that $L_{\text{reg}}^{-1}(i)=i$ for any $i \leq u$
- Generalized vectorization
 - Permute register values to get a larger number of contiguous elements

Rethink Tile-based Programming Models

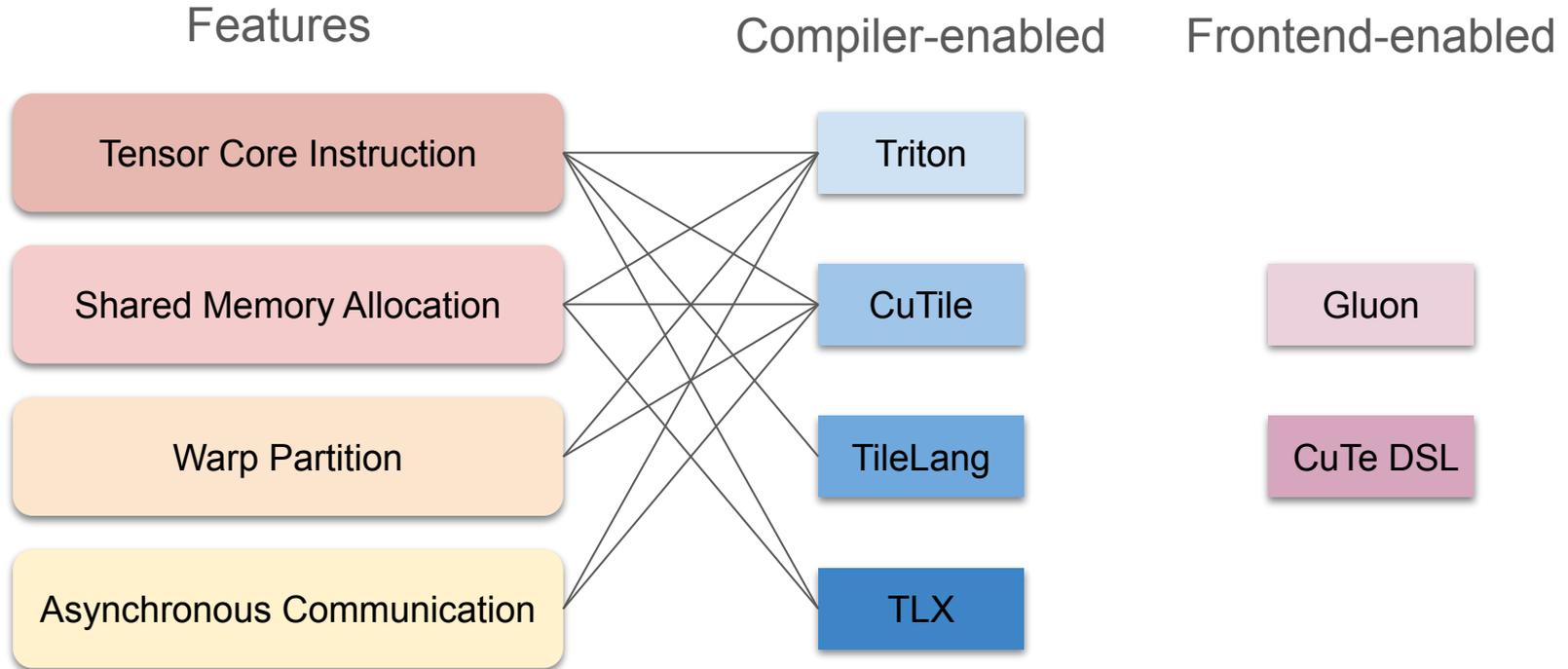
Tile-based Programming Models

- Programmers explicitly divide the work into tiles and write tiled code using APIs or frameworks
 - Gain fine-grained control
 - Write more flexible and sophisticated kernels
 - Allows for a relatively simple autotuner

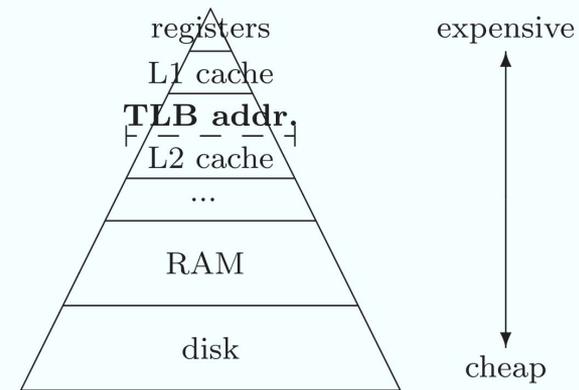
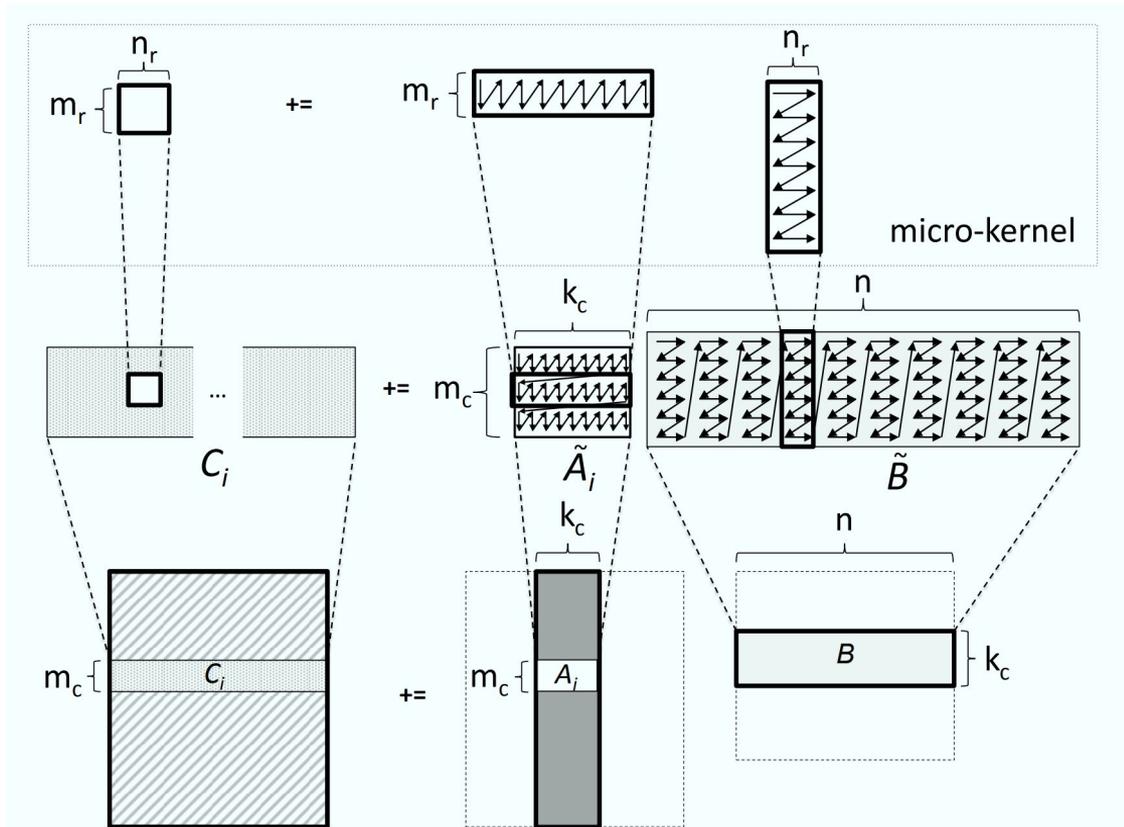
Tile-based Programming Models Comparison

- cuTile - MLIR
 - NVIDIA GPUs only
- NKI - MLIR
 - AWS Trainium and Inferentia
- TileLang - TVM
 - NVIDIA, AMD, and Huawei GPUs
- Triton - MLIR
 - AMD, NVIDIA, Intel GPUs, as well as a wide range of accelerators

Transformation Complexities

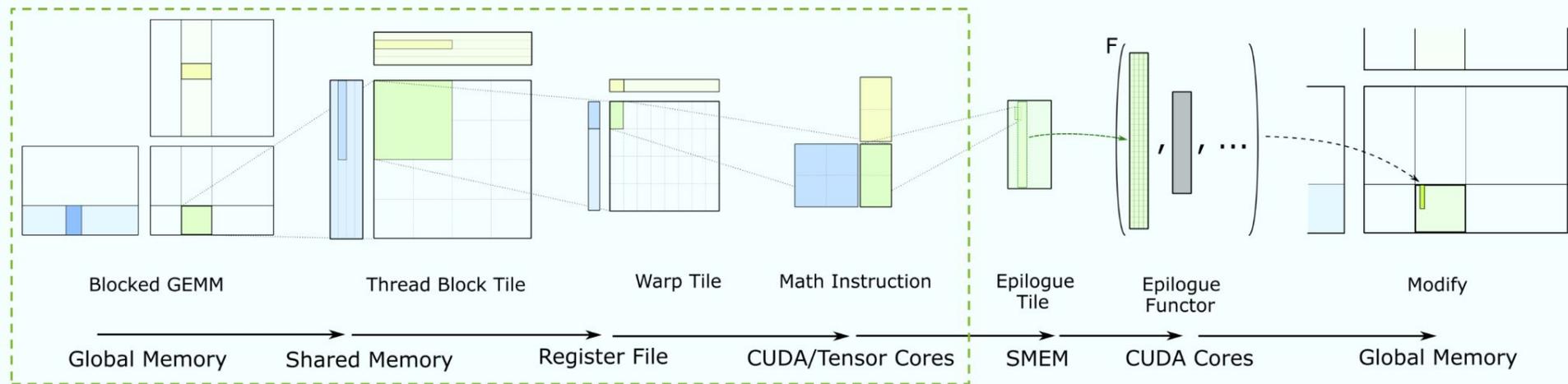


Sophisticated Tiling on CPUs



Refined Model

Sophisticated Tiling on GPUs



Tiled, hierarchical model: reuse data in Shared Memory and in Registers

Automated Transformation - Ideas

```
Matmul
1 # Multiply A[M][K] * B[K][N] = C[M][N]
2 for i in range(M):
3     for j in range(N):
4         C[i][j] = 0
5         for k in range(K):
6             C[i][j] += A[i][k] * B[k][j]
```



```
Matmul
1 # Assume tile size is TILE
2 for i0 in range(0, M, TILE):
3     for j0 in range(0, N, TILE):
4         for k0 in range(0, K, TILE):
5             for i in range(i0, min(i0 + TILE, M)):
6                 for j in range(j0, min(j0 + TILE, N)):
7                     for k in range(k0, min(k0 + TILE, K)):
8                         C[i][j] += A[i][k] * B[k][j]
```

Automated Transformation - Pros & Cons

Pros

- Developers write simple code; the compiler or transformation tool handles the complex optimization

Cons

- Requires sophisticated compilers (e.g., polyhedral frameworks). Not always available in all toolchains
- Customized code (e.g., fusion) may not be automatically optimized

Automatic Shared Memory Management

Pros

- Developers manage shared memory allocation and deallocation

Cons

- Manual shared memory management may lead to problems including
 - Memory leaks
 - Data races
 - Out of boundary accesses

Automatic Warp Specialization

Pros

- Developers determine the partition of the computation graph in each warp specialized region

Cons

- Requires a sophisticated performance model to estimate load balance among partitions involving
 - Instruction latency prediction
 - Cache hit prediction
 - Register pressure prediction

The Future Trend 1: Lightweight DSL/Compiler

- Building a high-level DSL compiler might be too heavy
 - Triton/CuTile requires a lot of engineering efforts
 - We need to abstract away low-level semantics
 - Shared memory
 - Registers
 - Barrier
- Research directions
 - Fast prototype custom, lightweight DSL co-designed with the underlying architecture
 - We don't always need one language for all architectures
 - Performance is still the first-class citizen

The Future Trend 2: LLM for Code Optimization

- A lot of research work has been done to generate code with higher performance than torch.compile on KernelBench
- MOKARA is an AI system for automated GPU kernel generation and optimization across CUDA, Triton, and other backends
 - You can try it out on <https://makora.com>
- Research directions
 - LLM for complex, end-to-end code optimization

The Future Trend 3: Profile-Guided Optimization

- Timing information is insufficient for LLMs to make wise and effective decisions about code transformation
 - Fine-grained performance metrics are necessary
- Wafer (<https://www.wafer.ai>) is a plugin for profiling, diagnosing, and optimizing your GPU code with LLM-based suggestions and fixes
- Research directions
 - Integrate not only hardware metrics but also custom metrics to LLMs
 - Modeling the code transformation effects to reduce the number of LLM tokens

Conclusions

Learning Resources

- [rkinas/triton-resources: A curated list of resources for learning and exploring Triton, OpenAI's programming language for writing efficient GPU code.](#)
- [gpu-mode/lectures: Material for gpu-mode lectures](#)
- <https://discord.gg/gpumode>

Triton Puzzles

- <https://github.com/srush/Triton-Puzzles>
- A set of questions for you to learn Triton from scratch
- You will start with trivial examples and build your way up to real algorithms like Flash Attention and Quantized neural networks
- These puzzles do not need to run on the GPU since they use the Triton interpreter

Visualization

- <https://deep-learning-profiling-tools.github.io/triton-viz/>
- One area that learners have trouble with is memory loading and storage which is critical for speed on low-level devices
- Triton-Viz is a visualizer that illustrates load/store and other triton operations using a user-friendly GUI
- It also provides simple a statistic summary about operations done by the triton kernel