# Tile-based Programming Models for Al

Keren Zhou kzhou6@gmu.edu



#### **Evolution of Al**











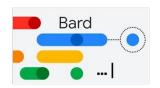


#### Al Software Stack



























## Hot LLM Topics

- Agent
  - Tool, Planner, Memory, Environment, Multi-Agent, ...
- Reasoning
  - o Chain-of-Thoughts, ReAct, Reflexion, Self-consistency, ...
- Reinforcement Learning
  - Tool-augmented, Instruction Tuning, Human Feedback, ...
- Systems
  - o KV Cache, Prefill, Decoding, Pipelining, Tensor Parallelism, Batching, ...
- Kernels
  - o MXFP, Quantization, MoE, Normalization, Batch GEMM, Linear/Flash/Paged/Attention, ...

We're going to focus on programming languages that empowers systems and kernels

## Outline

- Parallel Systems
- Tile-Based Programming
- Triton
- Triton-Puzzles

# Parallel Systems

## Overview

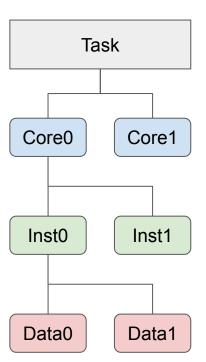
- Multi-core CPUs
- GPUs
- Accelerators

#### **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

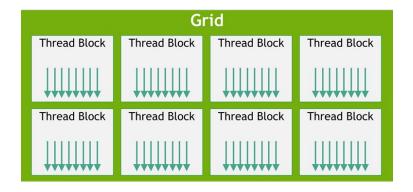
- Thread parallelism
- Instruction parallelism
- SIMD parallelism



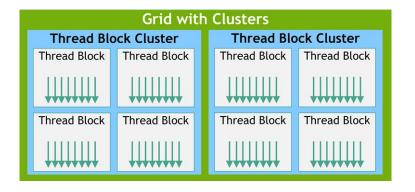
#### **GPUs**

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

#### **NVIDIA GPUs**



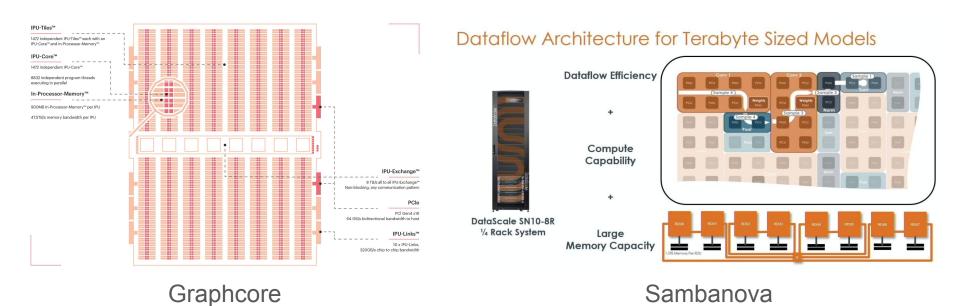
Before Hopper



Since Hopper

#### **Accelerators**

#### Different accelerators have different abstractions



# Tile-Based Programming

# Thread-Based Programming

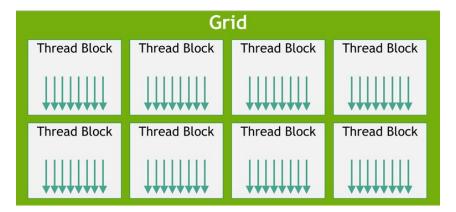
- Users specify the behavior of each thread
- The number of blocks and threads within each block is controlled on the host side

#### Behavior of Each Thread

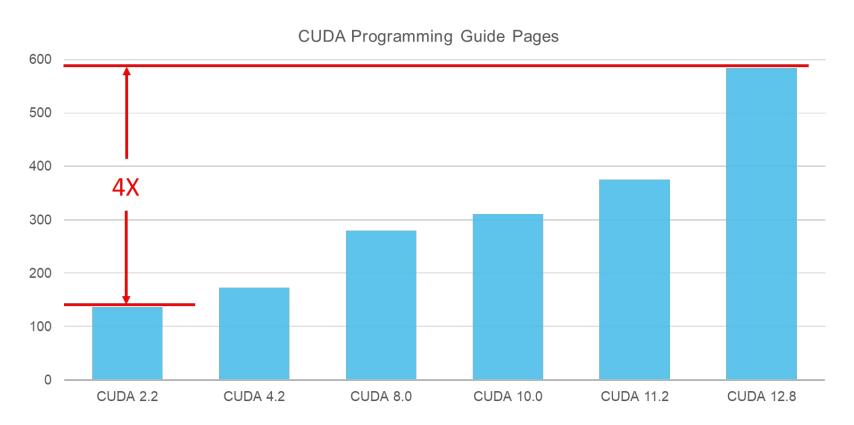
```
vecAdd
__global__ void vectorAdd(const float *A, const float *B, float *C, int numElements) {
    int i = blockDim.x * blockIdx.x + threadIdx.x;
    if (i < numElements) {</pre>
        C[i] = A[i] + B[i];
                                                                                       snappify.com
                                             Grid
                      Thread Block
                                    Thread Block
                                                 Thread Block
                                                               Thread Block
                      Thread Block
                                    Thread Block
                                                 Thread Block
                                                               Thread Block
```

#### Number of Blocks and Threads

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



# **CUDA Programming Guide**



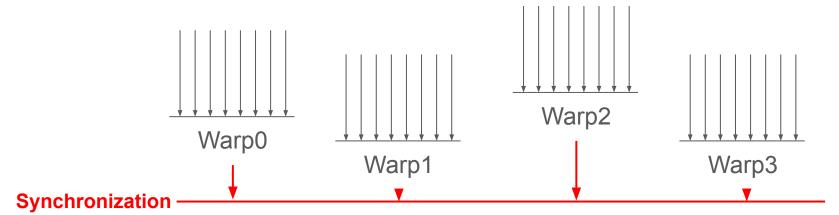
# The "Problems" with Full CUDA Specification

- 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"?

#### **Motivation**

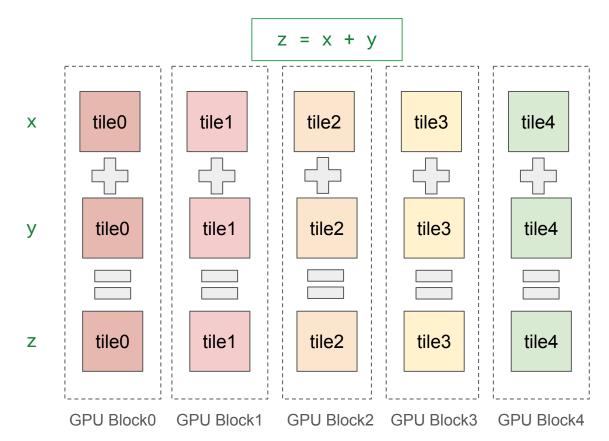
- Threads within a warp are executed in lock steps
  - This is the behavior of many simulators
  - Though not accurate with "independent thread scheduling"
- Warps within a thread block access the same programmable shared memory
- We can coarsen the execution unit as using warps or thread blocks



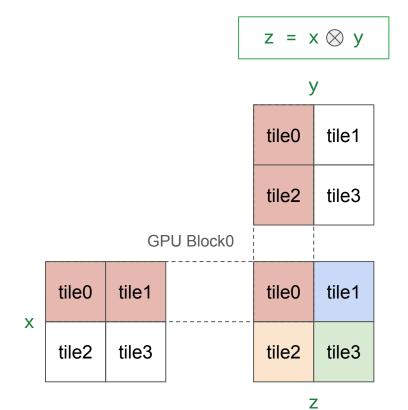
# Core Concepts - Tiling/Blocking

- 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

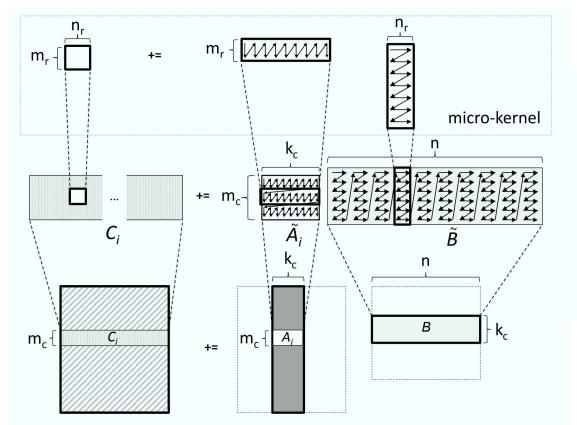
# **Example - Vector Addition**

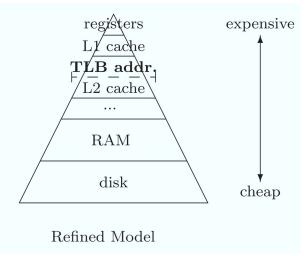


# **Example - Matrix Multiplication**



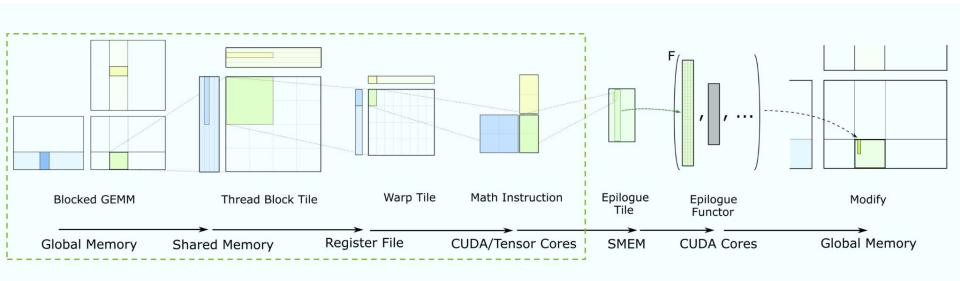
# Sophisticated Tiling on CPUs





Goto, Kazushige, and Robert A. van de Geijn. "Anatomy of high-performance matrix multiplication." ACM Transactions on Mathematical Software (TOMS) 34.3 (2008): 1-25.

# Sophisticated Tiling on GPUs



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

#### **Automated Transformation - Ideas**

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

# 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

- cuTile
  - NVIDIA GPUs
- NKI
  - AWS Trainium and Inferentia
- Triton
  - AMD, NVIDIA, Intel GPUs, as well as some accelerators
- ThunderKittens, TileLang
  - Research prototypes

# **Triton**

# Survey

- How many of you have heard of Triton before?
- How many of you know that Triton starts from a graduate student project at Harvard?

# Al System Software Stack



# Why Triton?









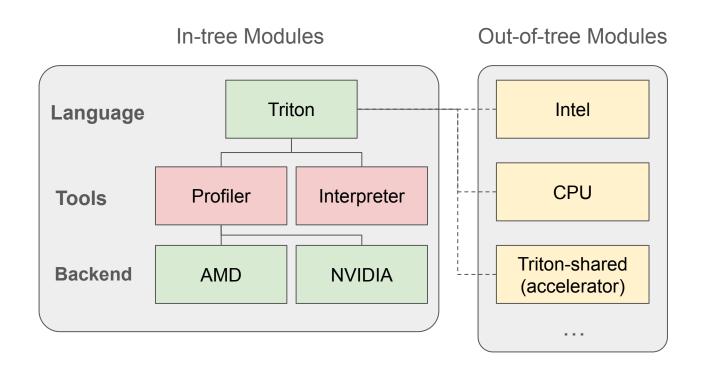








#### **Triton Modules**



# 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

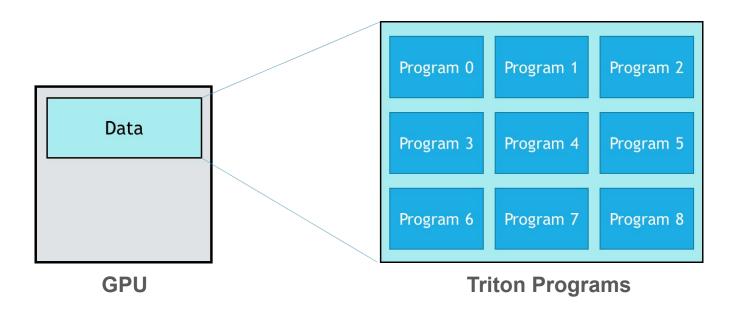
#### Triton vs CUDA

	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

#### SPMD Model

- Each program executes the "same" code
  - Only program ids are different



## A Simple Triton Program - Kernel Code

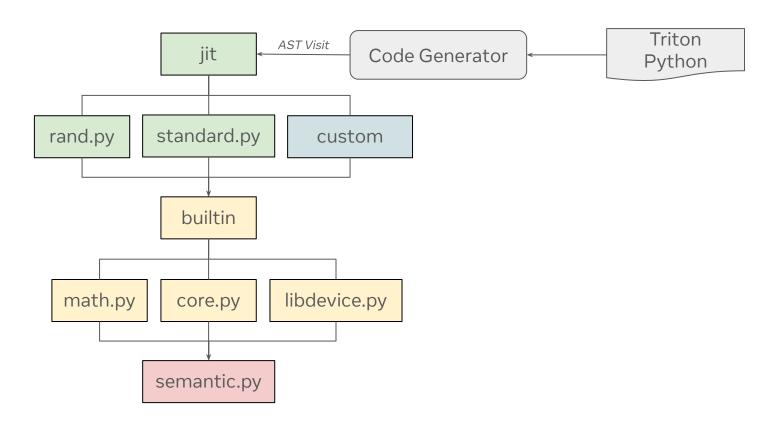
```
z: \dim 0 \times \dim 1 = x: \dim 0 \times \dim 1 + y: \dim 0 \times \dim 1
                                     vecAdd
   Kernel decorator
                              ລtriton.jit
                              def add_kernel(x_ptr, y_ptr, z_ptr, dim0, dim1,
                                              BLOCK_DIMO: tl.constexpr, BLOCK_DIM1: tl.constexpr):
                                pid_x = tl.program_id(axis=0)
Programming model
                                pid_y = tl.program_id(axis=1)
                                block_start = pid_x * BLOCK_DIMO * dim1 + pid_y * BLOCK_DIM1
                                offsets_dim0 = tl.arange(0, BLOCK_DIM0)[:.None]
       Creation ops
                                offsets_dim1 = tl.arange(0, BLOCK_DIM1)[None, :]
                                offsets = block_start + offsets_dim0 * dim1 + offsets_dim1
                                masks = (offsets_dim0 < dim0) & (offsets_dim1 < dim1)</pre>
                                x = tl.load(x_ptr + offsets, mask=masks)
       Memory ops
                                v = tl.load(y_ptr + offsets, mask=masks)
                                output = x + v
                                tl.store(z_ptr + offsets, output, mask=masks)
```

## A Simple Triton Program - Kernel Launch

```
z: \dim 0 \times \dim 1 = x: \dim 0 \times \dim 1 + y: \dim 0 \times \dim 1
```

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

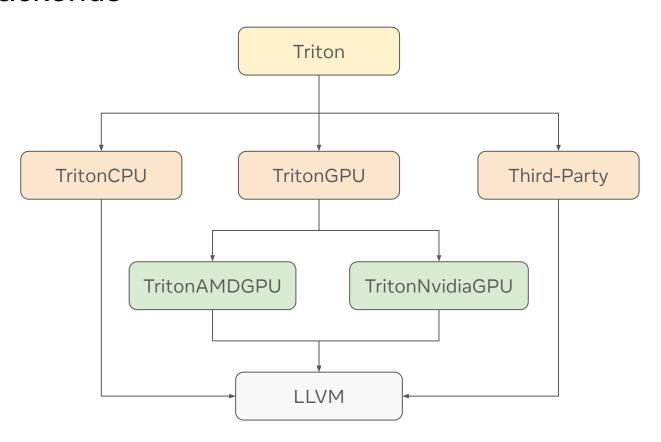
### **Triton Frontend**



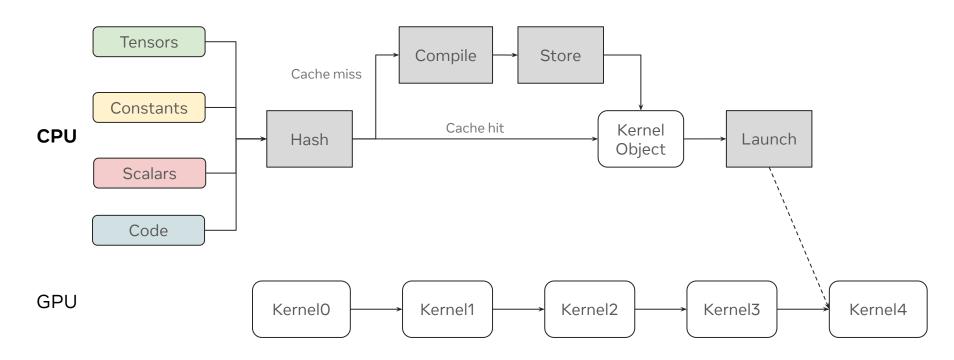
### MLIR: Multi-Level IR Compiler Framework

- 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)

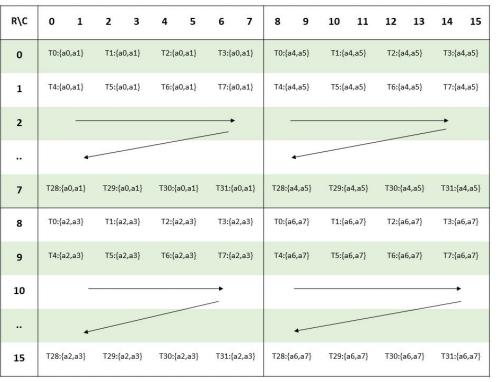
```
tt.func
                                     Mangled function name
@matmul kernel Pfp32 Pfp32 Pfp32 i32 i32 i32 i32 i32 i32 i32 i32 i32 12c64 13c64 14c64 15c8
(%arg0: !tt.ptr<f32> {tt.divisibility = 16 : i32}, ...) {
  %cst = arith.constant dense<true>: tensor<64x64xi1>
                                                                  Tensor
  %c64 = arith.constant 64 : i32
                                                                  Scalar
  %c0 = arith.constant 0 : i32
  \%0 = \text{tt.get program id } \times : i32
                                                              Triton operation
  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"

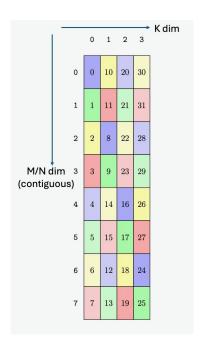
```
#blocked = #ttg.blocked<{sizePerThread = [8, 1], threadsPerWarp = [8, 4], warpsPerCTA = [1, 4], order = [0, 1]}>
#blocked1 = #ttg.blocked<{sizePerThread = [1, 8], threadsPerWarp = [4, 8], warpsPerCTA = [4, 1], order = [1, 0]}>
#mma = #ttg.nvidia mma<{versionMajor = 2, versionMinor = 0, warpsPerCTA = [4, 1], instrShape = [16, 8]}>
module attributes {"ttq.num-ctas" = 1 : i32, "ttq.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>
```

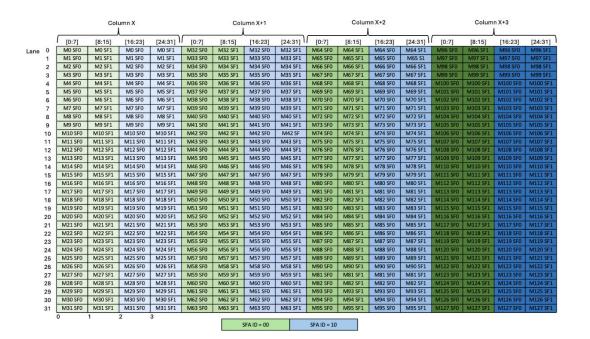
## Example Layouts - MMA



%laneid:{fragments}

## Example Layouts - tcgen05





## **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

### State-of-the-art GEMM Performance

FLOPS

Time

Numbers represents fp8 TFLOPS on Blackwell GB200

```
1690.079 3659.446 ROOT
└ nan 429.774 hvjet qqhsq 256x256 128x4 2x1 2cta v bz TNT

→ 1778.448 386.401 matmul kernel persistent [M=8192, N=8192, K=512]

→ 1972.040 348.469 matmul kernel tma persistent [M=8192, N=8192, K=512]

2359.637 291.229 matmul kernel tma persistent ws [M=8192, N=8192, K=512]
```

## Recent Advances in Tile-Based Programming Models

Paper	Conf	Target Problems
Rammer: Enabling Holistic Deep Learning Compiler Optimizations with rTasks	OSDI'20	Intra- and inter operator scheduling
Roller: Fast and Efficient Tensor Compilation for Deep Learning	OSDI'22	Tile-code generation with performance modeling
Welder: Scheduling Deep Learning Memory Access via Tile-graph	OSDľ23	Tile scheduling
Cocktailer: Analyzing and Optimizing Dynamic Control Flow in Deep Learning	OSDI'23	uTask scheduling
TensorIR: An Abstraction for Automatic Tensorized Program Optimization	ASPLOS'23	Block abstraction
Hidet: Task-mapping programming paradigm for deep learning tensor programs	ASPLOS'23	Tile with layouts
Graphene: An ir for optimized tensor computations on gpus	ASPLOS'23	Tile with layouts
Task-Based Tensor Computations on Modern GPUs	PLDI'25	Tile with warp specialization

# Triton-Puzzles

### **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

- 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

### **Get Started**

https://github.com/srush/Triton-Puzzles



## Learning Resources

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