



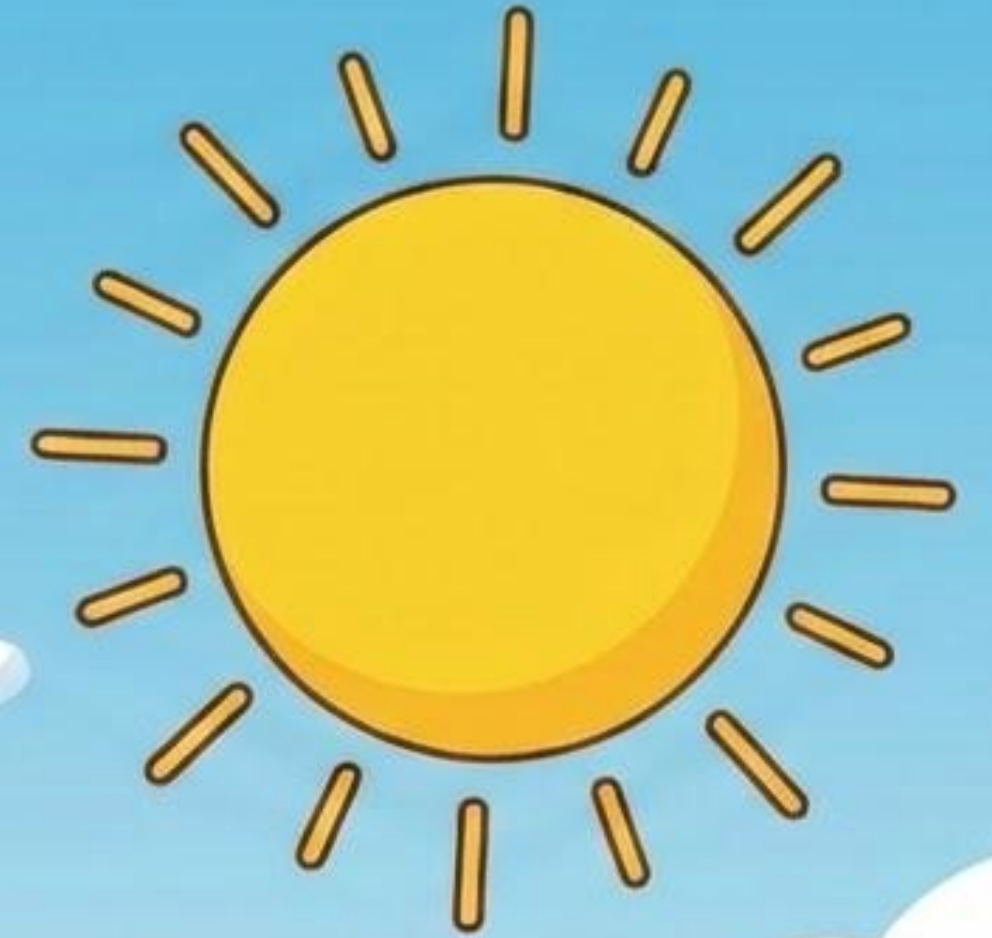
# CONDUCTING THE GPU SYMPHONY

WITH RUST  AND TRITON 

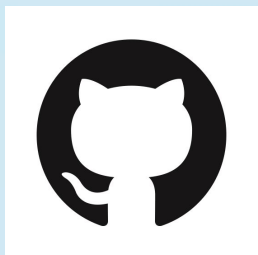
**KUBA SOŁTYS**



# **SOLUTIONS ARCHITECT**



**[HTTPS://WWW.YOUTUBE.COM/@QOوبا](https://www.youtube.com/@qooba)**



**[HTTPS://GITHUB.COM/QOوبا](https://github.com/qooba)**



**[HTTPS://WWW.LINKEDIN.COM/IN/QOوبا](https://www.linkedin.com/in/qooba)**



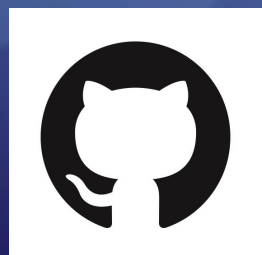


**BIELIK LM CONTRIBUTOR**

**WRITING GPU KERNELS  
& RUST CODE**



**[HTTPS://WWW.YOUTUBE.COM/@QOوبا](https://www.youtube.com/@qooba)**



**[HTTPS://GITHUB.COM/QOوبا](https://github.com/qooba)**



**[HTTPS://WWW.LINKEDIN.COM/IN/QOوبا](https://www.linkedin.com/in/qooba)**



# THE OVERTURE

**f FROM AGENTS TO KERNELS**

4/4 time signature. Treble clef. Dynamics: *f*. Includes a neural network icon in the top right.

**f TRITON VS CUDA**

4/4 time signature. Treble clef. Dynamics: *f*. Includes a trill ornament. Includes a neural network icon in the top right.

**f INSIDE BIELIK 1.5B**

4/4 time signature. Treble clef. Dynamics: *f*. Includes a neural network icon in the top right.

**espressivo, C++ MatMul, Softmax, RMSNorm**

4/4 time signature. Treble clef. Dynamics: *f*. Includes a neural network icon in the top right.

# THE GRAND FINALE

**mf AOT f JIT VS AOT TRITON**

4/4 time signature. Treble clef. Dynamics: *mf*, *f*. Includes a trill ornament. Includes a neural network icon in the top right.

**f DEEP DIVE: TRITON TO PTX**

2/2 time signature. Treble clef. Dynamics: *f*. Includes a trill ornament. Includes a neural network icon in the top right.

**mf RUST BINDINGS**

4/4 time signature. Bass clef. Dynamics: *mf*. Includes a Rust logo. Includes a neural network icon in the top right.

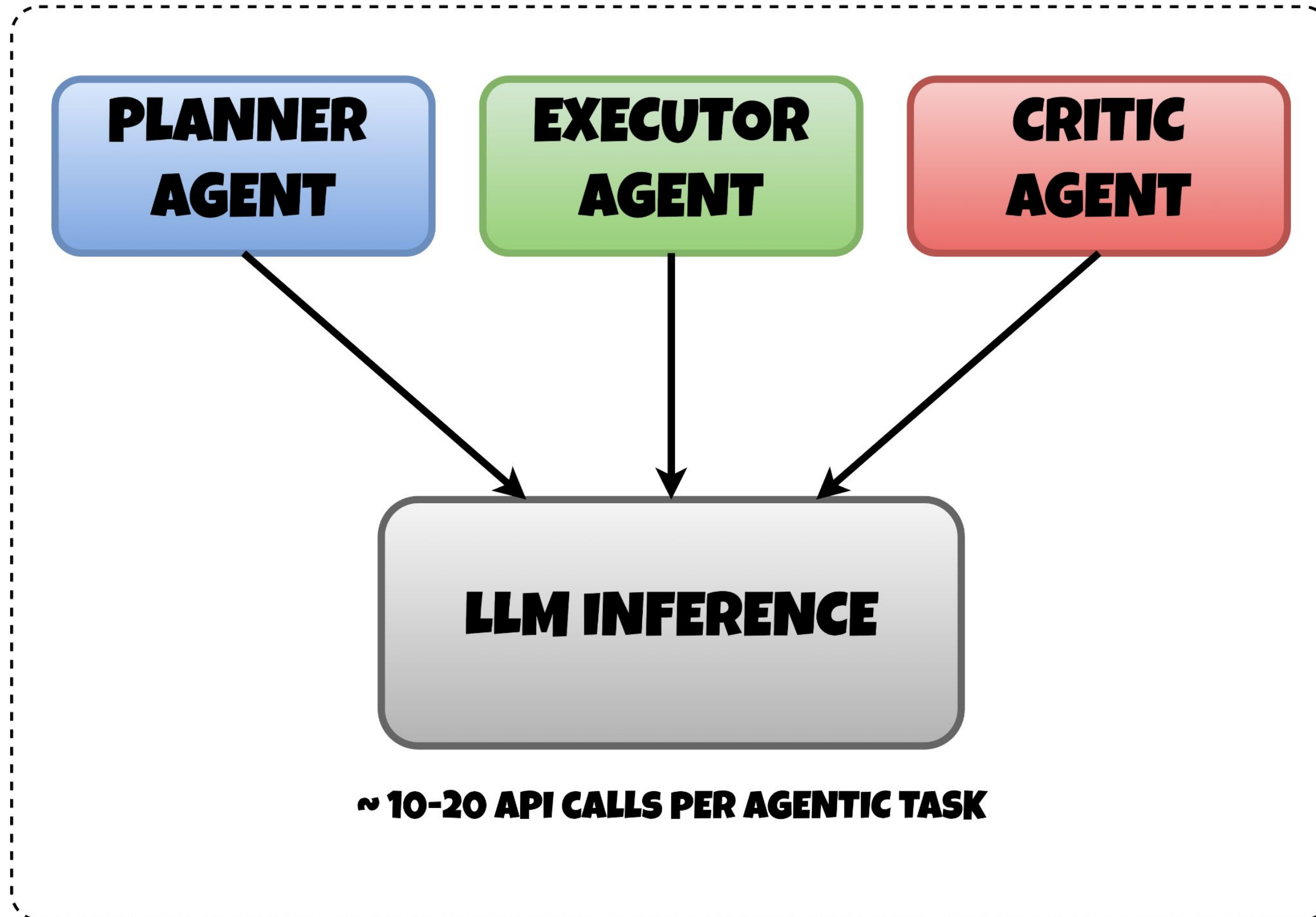
**BENCHMARK: RUST VS PYTHON**

4/4 time signature. Treble clef. Dynamics: *ff*, *pp*, *ff*. Includes a Rust logo and Python logo. Includes a bar chart comparing Rust and Python speed. Includes a neural network icon in the top right.

# WHY GPU KERNELS?



## AGENTIC LOOP



**INPUT TOKENS**

**EMBEDDING**

**LAYER 1**

**LAYER 2**

**LAYER 3**

...

**LAYER 32**

**RMS NORM**

**LM HEAD**

**LOGITS**

**ONE LAYER (EXPANDED)**

**RMS NORM**

**QKV PROJECTION**

**CAUSAL SOFTMAX**

**OUTPUT PROJ.**

**MLP (SWIGLU)**

# GPU EXECUTION

**WORK DISTRIBUTION**

**TILE PARTITIONING / SM ALLOCATION**

**MEMORY COALESCING (HBM -> SRAM)**

**FUSED CAUSAL CALCULATION**

**GLOBAL MEMORY COMMIT (VRAM)**

**FASTER KERNEL -> FASTER TOKEN -> FASTER CALL -> FASTER TASK**

# CUDA VS TRITON: ABSTRACTION LEVEL

## RAW CUDA C++

## TRITON

```
1 // You manage every thread
2
3 __global__ void softmax_kernel(
4     float* x, float* out, int n
5 ) {
6     int tid = blockIdx.x * blockDim.x
7         + threadIdx.x;
8
9     // manual warp sync
10    __syncwarp();
11
12    // manual shared mem
13    __shared__ float shmem[1024];
14
15    // ... 200 more lines
16
17 }
```

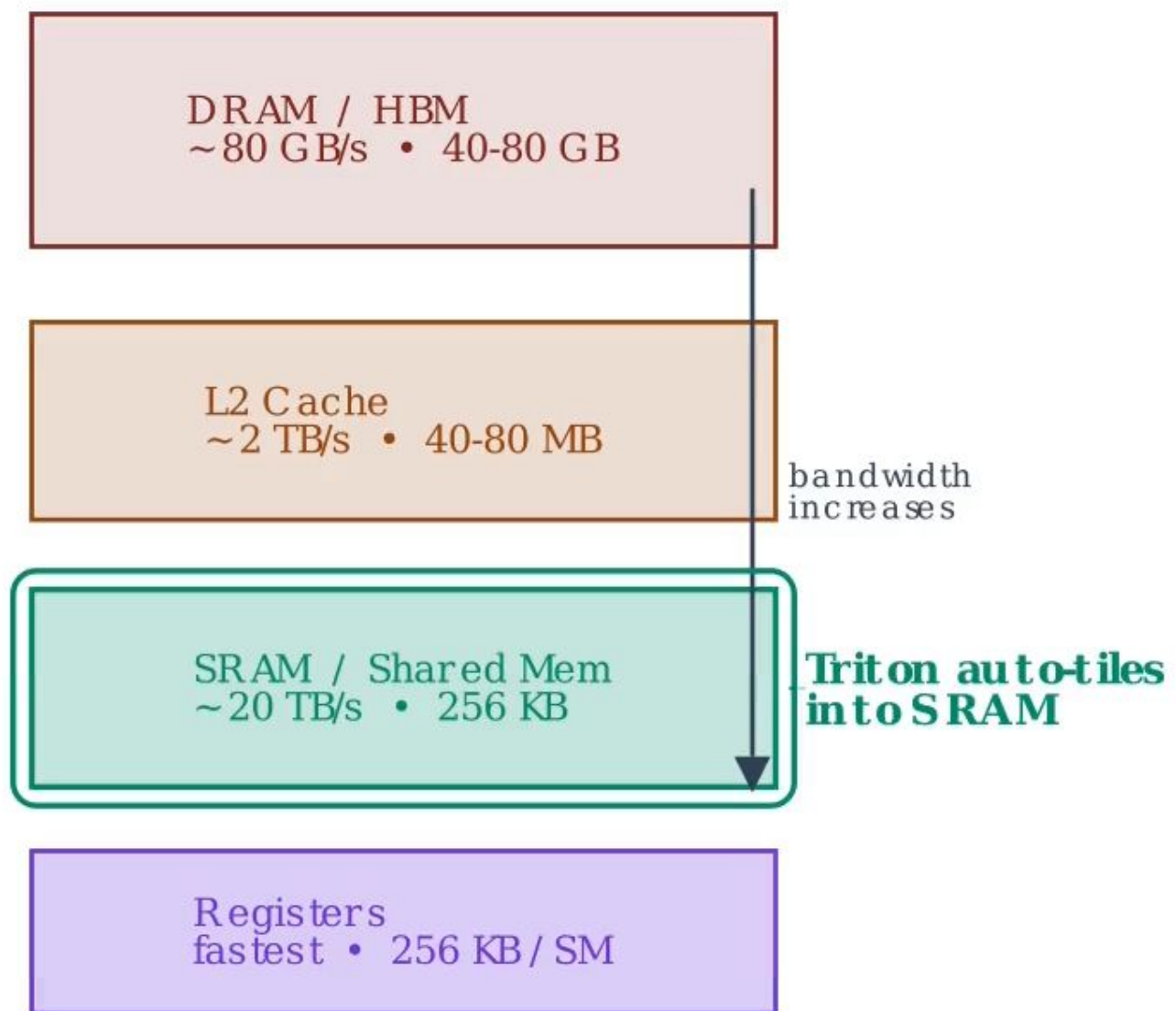
```
1 @triton.jit
2 def softmax_kernel(
3     x_ptr, out_ptr, n_rows, n_cols,
4     BLOCK_SIZE: tl.constexpr
5 ):
6     row = tl.program_id(0) # block idx
7     ofs = tl.arange(0, BLOCK_SIZE)
8     x = tl.load(x_ptr + row * n_cols
9                + ofs,
10                mask=ofs < n_cols,
11                other=-float('inf'))
12     x -= tl.max(x, axis=0)
13     e = tl.exp(x)
14     out = e / tl.sum(e, axis=0)
15     tl.store(out_ptr + row * n_cols
16            + ofs, out,
17            mask=ofs < n_cols)
```

**THREAD-LEVEL CONTROL. WARP DIVERGENCE,  
BANK CONFLICTS, COALESCING - ALL YOURS.**

**BLOCK-LEVEL. TRITON HANDLES SRAM,  
COALESCING & SYNC.**

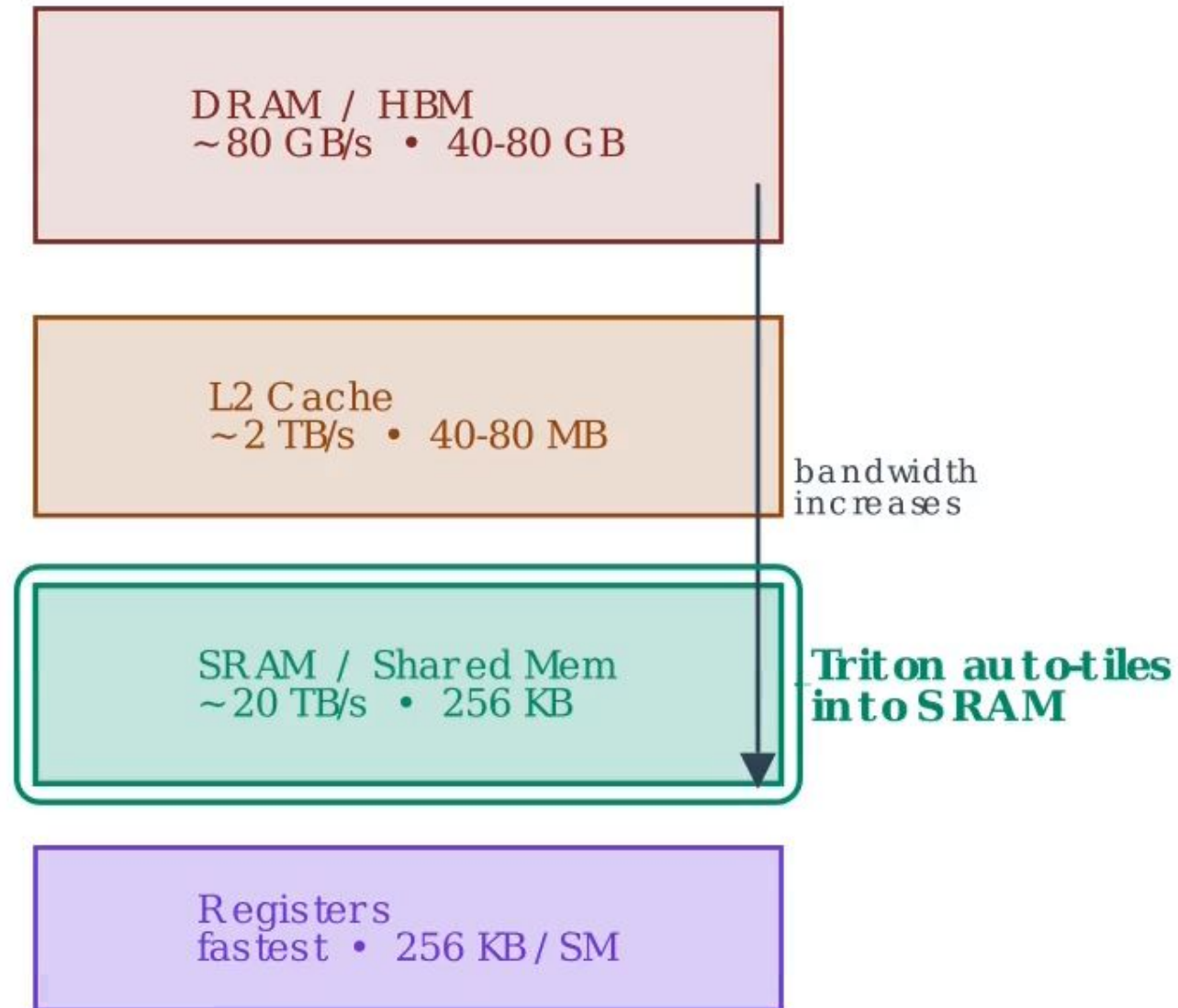
# GPU MEMORY HIERARCHY & TRITON'S ABSTRACTION

## Memory Hierarchy



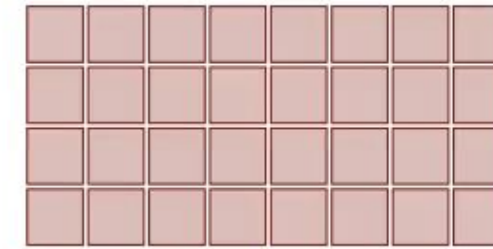
# GPU MEMORY HIERARCHY & TRITON'S ABSTRACTION

## Memory Hierarchy



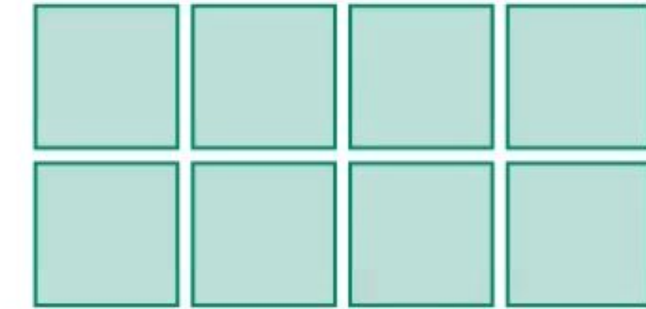
## Triton vs CUDA Programming

### CUDA: thread-level



Each thread → 1 element  
You manage tiling, sync, banks

### Triton: block-level



Each program → 1 tile (BLOCK\_SIZE)  
Auto SRAM, coalescing, sync

### Triton kernel — block-level thinking

```
pid = tl.program_id(0)      # block index
offs = pid * BLOCK_SIZE + tl.arange(0, BLOCK_SIZE)
x = tl.load(ptr + offs)    # load tile → SRAM
y = tl.exp(x - tl.max(x))  # compute in regs
```

**Triton hides warp/thread management & SRAM scheduling**

# WHAT TRITON GIVES US

- ✓ **Automatic SRAM management** - tiles loaded to shared memory without explicit `__shared__` arrays
- ✓ **Global memory coalescing** - Triton ensures consecutive threads access consecutive addresses
- ✓ **Latency hiding** - instruction-level pipelining between memory ops and arithmetic
- ✓ **Tensor cores for free** - `tl.dot(a, b, allow_tf32=True)` emits `mma.sync` instructions automatically; no manual warp-level intrinsics
- ✓ **Auto-tuning** - `@triton.autotune` tries multiple `BLOCK_SIZE` values, picks best

**Triton's insight:** optimal GPU kernels differ only in *tile size*, not structure.  
Write the structure once; let the compiler tune the rest.

# MATMUL, SOFTMAX, RMSNORM



# MATMUL

```
1 import torch
2 import triton
3 import triton.language as tl
4
5 @triton.jit
6 def matmul_tiled_kernel(
7     a_ptr, b_ptr, c_ptr,
8     M, N, K,
9     BLOCK_M: tl.constexpr, BLOCK_N: tl.constexpr, BLOCK_K: tl.constexpr,
10 ):
11     pid_m = tl.program_id(axis=0)
12     pid_n = tl.program_id(axis=1)
13
14     offs_m = pid_m * BLOCK_M + tl.arange(0, BLOCK_M)
15     offs_n = pid_n * BLOCK_N + tl.arange(0, BLOCK_N)
16
17     acc = tl.zeros((BLOCK_M, BLOCK_N), dtype=tl.float32)
18
19     for k_start in range(0, K, BLOCK_K):
20         offs_k = k_start + tl.arange(0, BLOCK_K)
21
22         a_tile = tl.load(a_ptr + offs_m[:, None] * K + offs_k[None, :],
23                         mask=(offs_m[:, None] < M) & (offs_k[None, :] < K),
24                         other=0.0,
25                         )
26
27         b_tile = tl.load(
28             b_ptr + offs_k[:, None] * N + offs_n[None, :],
29             mask=(offs_k[:, None] < K) & (offs_n[None, :] < N),
30             other=0.0,
31             )
32
33         acc += tl.dot(a_tile, b_tile, allow_tf32=True)
34
35     out_mask = (offs_m[:, None] < M) & (offs_n[None, :] < N)
36     tl.store(
37         c_ptr + offs_m[:, None] * N + offs_n[None, :],
38         acc,
39         mask=out_mask,
40     )
```

# CUDA CORE VS TENSOR CORE

## CUDA CORE

$$c = a \times b + c$$

Example:

$$3.5 \times 2.1 + 0 = 7.35$$

Performance:

~19 TFLOPS (FP32)

op1

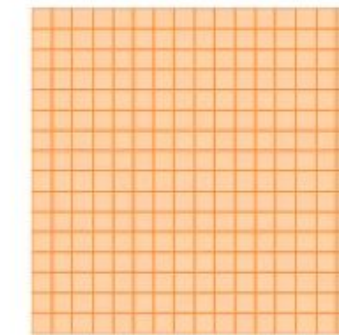
op2

op3

op4

## TENSOR CORE

$$D_{16 \times 16} = A_{16 \times 16} @ B_{16 \times 16} + C_{16 \times 16}$$



Performance:

~312 TFLOPS (FP16)

**16\* FASTER!**



```

1 import triton
2 import triton.language as tl
3
4 @triton.jit
5 def softmax_causal_simple_kernel(
6     input_ptr, output_ptr, n_rows, n_cols, seq_len,
7     input_row_stride, output_row_stride,
8     BLOCK_SIZE: tl.constexpr,
9 ):
10     row_idx = tl.program_id(axis=0)
11
12     if row_idx >= n_rows:
13         return
14
15     position = row_idx % seq_len
16
17     input_row_start = input_ptr + row_idx * input_row_stride
18     output_row_start = output_ptr + row_idx * output_row_stride
19
20     col_offsets = tl.arange(0, BLOCK_SIZE)
21     mask = col_offsets < n_cols
22
23     row_vals = tl.load(input_row_start + col_offsets, mask=mask, other=-float('inf'))
24     is_future = col_offsets > position
25     row_vals = tl.where(is_future, -float('inf'), row_vals)
26
27     row_max = tl.max(row_vals, axis=0)
28     numerator = tl.exp(row_vals - row_max)
29     denominator = tl.sum(numerator, axis=0)
30     softmax_output = numerator / denominator
31
32     tl.store(output_row_start + col_offsets, softmax_output, mask=mask)

```

**CAUSAL MASK**

## DECODER

**RMS NORM**

**QKV PROJECTION**

**CAUSAL SOFTMAX**

**OUTPUT PROJ.**

**MLP (SWIGLU)**

$$\text{softmax}(\mathbf{z})_i = \frac{e^{z_i - \max(\mathbf{z})}}{\sum_{j=1}^N e^{z_j - \max(\mathbf{z})}}$$

$$\text{attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$



```
1 import triton
2 import triton.language as tl
3
4
5 @triton.jit
6 def rms_norm_simple_kernel(
7     X, W, Y,
8     stride_x_row, stride_y_row,
9     n_cols, eps,
10    BLOCK_SIZE: tl.constexpr,
11 ):
12     row_idx = tl.program_id(axis=0)
13
14     X += row_idx * stride_x_row
15     Y += row_idx * stride_y_row
16
17     col_offsets = tl.arange(0, BLOCK_SIZE)
18     mask = col_offsets < n_cols
19
20     x = tl.load(X + col_offsets, mask=mask, other=0.0)
21     w = tl.load(W + col_offsets, mask=mask, other=1.0)
22
23     x_squared = x * x
24     mean_x_squared = tl.sum(x_squared) / n_cols
25     rms = tl.sqrt(mean_x_squared + eps)
26     x_normalized = x / rms
27     y = x_normalized * w
28
29     tl.store(Y + col_offsets, y, mask=mask)
```

## DECODER

**RMS NORM**

**QKV PROJECTION**

**CAUSAL SOFTMAX**

**OUTPUT PROJ.**

**MLP (SWIGLU)**

$$\text{RMSNorm}(\mathbf{x}) = \frac{\mathbf{x}}{\text{RMS}(\mathbf{x}) + \epsilon} \odot \mathbf{w}$$

$$\text{RMS}(\mathbf{x}) = \sqrt{\frac{1}{d} \sum_{i=1}^d x_i^2}$$

# JIT → AOT: THE MINDSET SHIFT

## JIT (PyTorch standard)

- ✓ Convenient - kernel compiled on first call
- ✓ BLOCK\_SIZE tuned per input shape
- ⚠ Requires Python + Triton at runtime
- ⚠ First-call latency ("cold start")
- ⚠ 500+ MB Python dependency tree

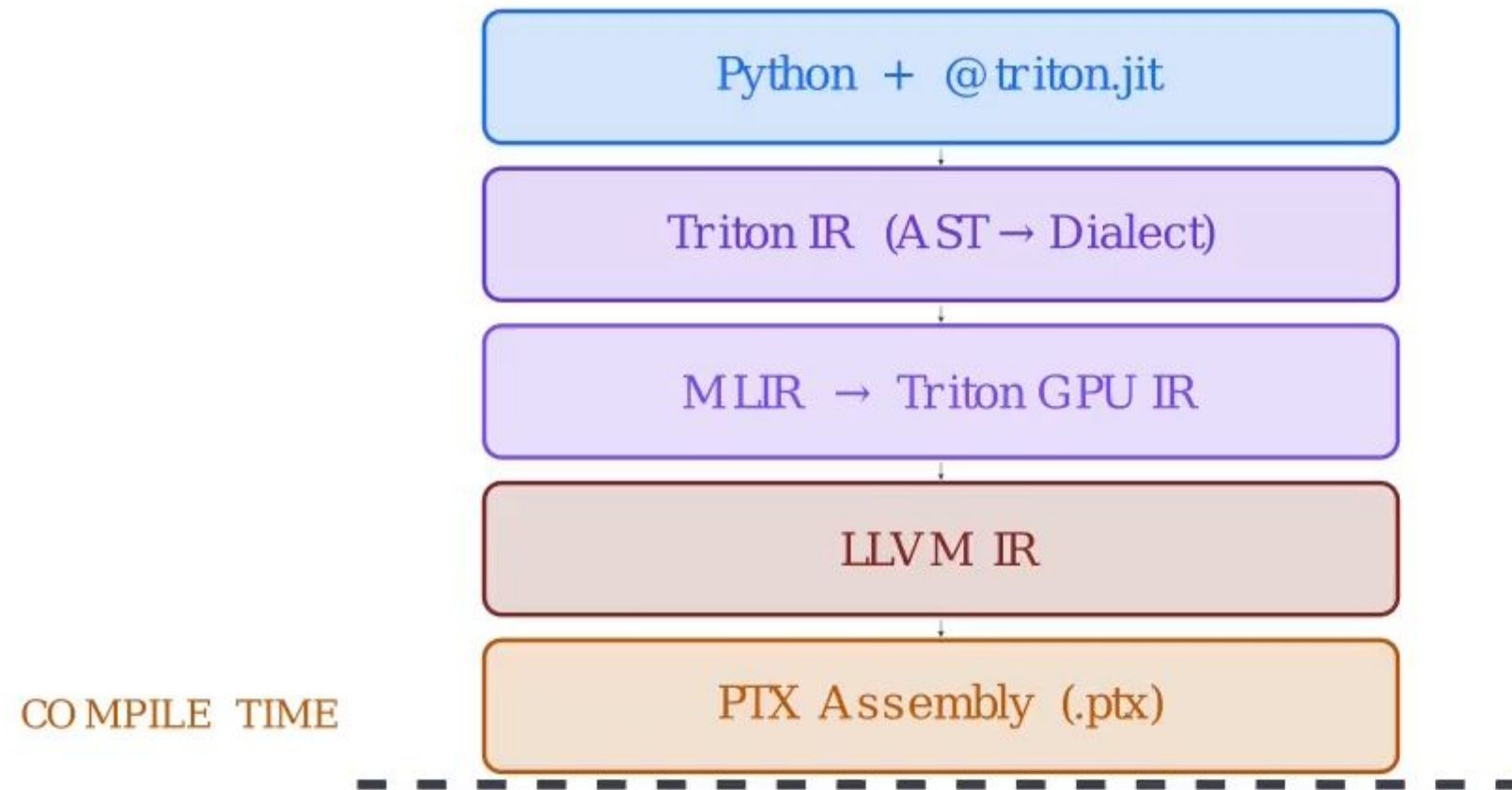
## AOT (our approach)

- ✓ PTX compiled *once*, embedded in binary
- ✓ Zero Python at runtime
- ✓ Instant load - no compilation step
- ⚠ BLOCK\_SIZE *frozen* at compile time
- ⚠ Must target specific GPU (SM arch)

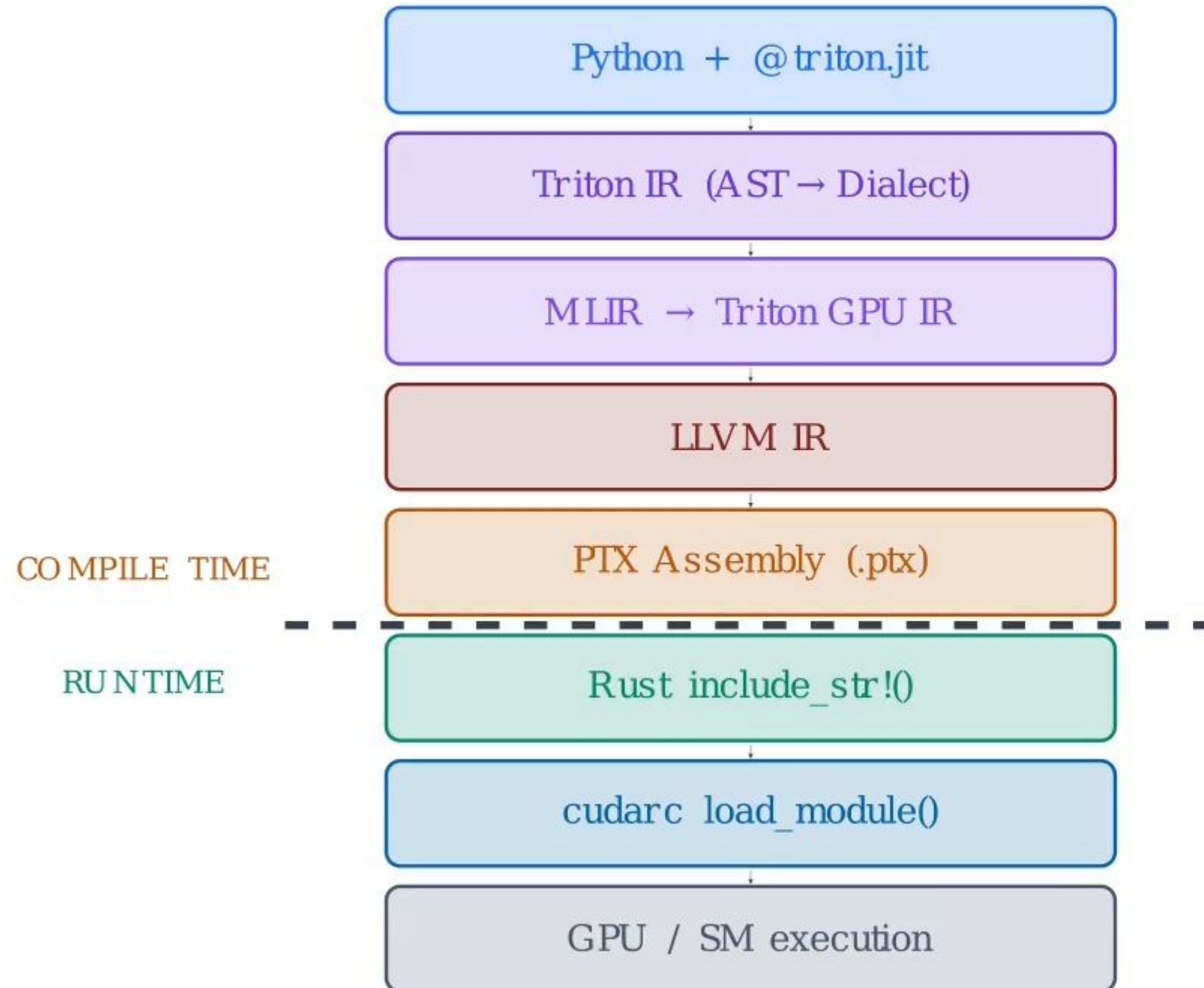
# TRITON COMPILATION DEEP DIVE



# AOT COMPILATION: PYTHON → GPU ASSEMBLY

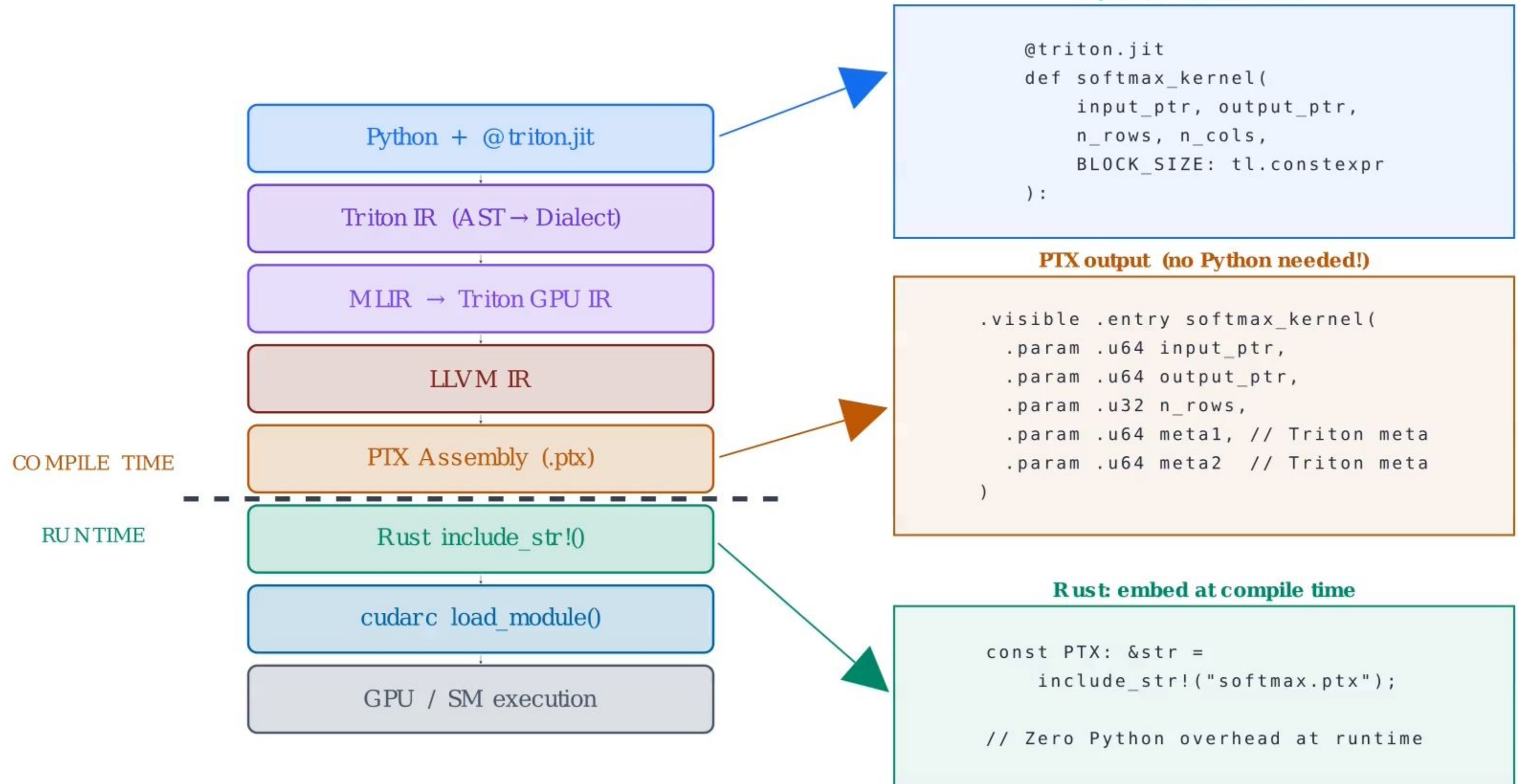


# AOT COMPILATION: PYTHON → GPU ASSEMBLY



# AOT COMPILATION: PYTHON → GPU ASSEMBLY

Python / Triton source



```
@triton.jit
def softmax_kernel(
    input_ptr, output_ptr,
    n_rows, n_cols,
    BLOCK_SIZE: tl.constexpr
):
```

**PTX output (no Python needed!)**

```
.visible .entry softmax_kernel(
    .param .u64 input_ptr,
    .param .u64 output_ptr,
    .param .u32 n_rows,
    .param .u64 meta1, // Triton meta
    .param .u64 meta2 // Triton meta
)
```

**Rust: embed at compile time**

```
const PTX: &str =
    include_str!("softmax.ptx");

// Zero Python overhead at runtime
```

```

1 import triton
2 import triton.language as tl
3 import triton.compiler as tc
4
5 @triton.jit
6 def add_kernel(x_ptr, y_ptr, ouptut_ptr, n_elements, BLOCK_SIZE:
  tl.constexpr):
7     pid = tl.program_id(0)
8     block_start = pid * BLOCK_SIZE
9
10    offset = block_start + tl.arange(0, BLOCK_SIZE)
11    mask = offset < n_elements
12
13    x = tl.load(x_ptr + offset, mask=mask, other=0.0)
14    y = tl.load(y_ptr + offset, mask=mask, other=0.0)
15
16    output = x + y
17
18    tl.store(ouptut_ptr + offset, output, mask=mask)
19

```



```

1 target = triton.runtime.driver.active.get_current_target()
2
3 src = tc.ASTSource(
4     fn=add_kernel,
5     signature={"x_ptr": "*fp32", "y_ptr": "*fp32",
6               "ouptut_ptr": "*fp32", "n_elements": "i32"},
7     constexprs={"BLOCK_SIZE": 256}
8 )
9
10 compiled_kernel = triton.compile(src, target=target)
11
12 print("# TRITON IR")
13 print("---")
14 print(compiled_kernel.asm["ttir"])
15 print("# TRITON-GPU IR")
16 print("---")
17 print(compiled_kernel.asm["ttgir"])
18 print("# LLVM IR")
19 print("---")
20 print(compiled_kernel.asm["llir"])
21 print("# PTX")
22 print("---")
23 print(compiled_kernel.asm["ptx"])
24 print("# SASS")
25 print("---")
26 print(compiled_kernel.asm["sass"])

```

# STEP 1: TRITON IR

**GPU-AGNOSTIC TILE OPERATIONS.** PARSSES PYTHON AST, RESOLVES TL.CONSTEXPR (BLOCK\_SIZE BAKED IN). NO THREADS, NO WARPS YET.

```
1 #loc = loc("/home/qooba/add_kernel_compile.py":6:0)
2 #loc15 = loc("x_ptr"(#loc))
3 #loc16 = loc("y_ptr"(#loc))
4 #loc17 = loc("ouptut_ptr"(#loc))
5 #loc18 = loc("n_elements"(#loc))
6 module {
7   tt.func public @add_kernel(%x_ptr: !tt.ptr<f32> loc("x_ptr"(#loc)), %y_ptr: !tt.ptr<f32>
  loc("y_ptr"(#loc)), %ouptut_ptr: !tt.ptr<f32> loc("ouptut_ptr"(#loc)), %n_elements: i32
  loc("n_elements"(#loc))) attributes {noinline = false} {
8     %cst = arith.constant dense<0.000000e+00> : tensor<256xf32> loc(#loc1)
9     %c256_i32 = arith.constant 256 : i32 loc(#loc1)
10    %pid = tt.get_program_id x : i32 loc(#loc19)
11    %block_start = arith.muli %pid, %c256_i32 : i32 loc(#loc20)
12    %offset = tt.make_range {end = 256 : i32, start = 0 : i32} : tensor<256xi32> loc(#loc21)
13    %offset_0 = tt.splat %block_start : i32 -> tensor<256xi32> loc(#loc22)
14    %offset_1 = arith.addi %offset_0, %offset : tensor<256xi32> loc(#loc22)
15    %mask = tt.splat %n_elements : i32 -> tensor<256xi32> loc(#loc23)
16    %mask_2 = arith.cmpi slt, %offset_1, %mask : tensor<256xi32> loc(#loc23)
17    %x = tt.splat %x_ptr : !tt.ptr<f32> -> tensor<256x!tt.ptr<f32>> loc(#loc24)
18    %x_3 = tt.addptr %x, %offset_1 : tensor<256x!tt.ptr<f32>>, tensor<256xi32> loc(#loc24)
19    %x_4 = tt.load %x_3, %mask_2, %cst : tensor<256x!tt.ptr<f32>> loc(#loc25)
20    %y = tt.splat %y_ptr : !tt.ptr<f32> -> tensor<256x!tt.ptr<f32>> loc(#loc26)
21    %y_5 = tt.addptr %y, %offset_1 : tensor<256x!tt.ptr<f32>>, tensor<256xi32> loc(#loc26)
22    %y_6 = tt.load %y_5, %mask_2, %cst : tensor<256x!tt.ptr<f32>> loc(#loc27)
23    %output = arith.addf %x_4, %y_6 : tensor<256xf32> loc(#loc28)
24    %0 = tt.splat %ouptut_ptr : !tt.ptr<f32> -> tensor<256x!tt.ptr<f32>> loc(#loc12)
25    %1 = tt.addptr %0, %offset_1 : tensor<256x!tt.ptr<f32>>, tensor<256xi32> loc(#loc12)
26    tt.store %1, %output, %mask_2 : tensor<256x!tt.ptr<f32>> loc(#loc13)
27    tt.return loc(#loc14)
28  } loc(#loc)
```

# STEP 2: TRITONGPU IR

**TILES MAPPED TO HARDWARE. DECIDES WARPS, THREADS, SHARED MEMORY. EACH TENSOR GETS A LAYOUT ANNOTATION.**

```
1 #blocked = #ttg.blocked<{sizePerThread = [1], threadsPerWarp = [32], warpsPerCTA = [4], order = [0]}>
2 #loc = loc("/home/qooba/add_kernel_compile.py":6:0)
3 #loc15 = loc("x_ptr"(#loc))
4 #loc16 = loc("y_ptr"(#loc))
5 #loc17 = loc("ouptut_ptr"(#loc))
6 #loc18 = loc("n_elements"(#loc))
7 module attributes {"ttg.num-ctas" = 1 : i32, "ttg.num-warps" = 4 : i32, ttg.target = "cuda:89", "ttg.threads-per-warp" =
  32 : i32} {
8   tt.func public @add_kernel(%x_ptr: !tt.ptr<f32> loc("x_ptr"(#loc)), %y_ptr: !tt.ptr<f32> loc("y_ptr"(#loc)),
  %ouptut_ptr: !tt.ptr<f32> loc("ouptut_ptr"(#loc)), %n_elements: i32 loc("n_elements"(#loc))) attributes {noinline = false}
  {
9     %c256_i32 = arith.constant 256 : i32 loc(#loc1)
10    %cst = arith.constant dense<0.000000e+00> : tensor<256xf32, #blocked> loc(#loc1)
11    %pid = tt.get_program_id x : i32 loc(#loc19)
12    %block_start = arith.muli %pid, %c256_i32 : i32 loc(#loc20)
13    %offset = tt.make_range {end = 256 : i32, start = 0 : i32} : tensor<256xi32, #blocked> loc(#loc21)
14    %offset_0 = tt.splat %block_start : i32 -> tensor<256xi32, #blocked> loc(#loc22)
15    %offset_1 = arith.addi %offset_0, %offset : tensor<256xi32, #blocked> loc(#loc22)
16    %mask = tt.splat %n_elements : i32 -> tensor<256xi32, #blocked> loc(#loc23)
17    %mask_2 = arith.cmpi slt, %offset_1, %mask : tensor<256xi32, #blocked> loc(#loc23)
18    %x = tt.splat %x_ptr : !tt.ptr<f32> -> tensor<256x!tt.ptr<f32>, #blocked> loc(#loc24)
19    %x_3 = tt.addptr %x, %offset_1 : tensor<256x!tt.ptr<f32>, #blocked>, tensor<256xi32, #blocked> loc(#loc24)
20    %x_4 = tt.load %x_3, %mask_2, %cst : tensor<256x!tt.ptr<f32>, #blocked> loc(#loc25)
21    %y = tt.splat %y_ptr : !tt.ptr<f32> -> tensor<256x!tt.ptr<f32>, #blocked> loc(#loc26)
22    %y_5 = tt.addptr %y, %offset_1 : tensor<256x!tt.ptr<f32>, #blocked>, tensor<256xi32, #blocked> loc(#loc26)
23    %y_6 = tt.load %y_5, %mask_2, %cst : tensor<256x!tt.ptr<f32>, #blocked> loc(#loc27)
24    %output = arith.addf %x_4, %y_6 : tensor<256xf32, #blocked> loc(#loc28)
25    %0 = tt.splat %ouptut_ptr : !tt.ptr<f32> -> tensor<256x!tt.ptr<f32>, #blocked> loc(#loc12)
26    %1 = tt.addptr %0, %offset_1 : tensor<256x!tt.ptr<f32>, #blocked>, tensor<256xi32, #blocked> loc(#loc12)
27    tt.store %1, %output, %mask_2 : tensor<256x!tt.ptr<f32>, #blocked> loc(#loc13)
28    tt.return loc(#loc14)
29  } loc(#loc)
```

# STEP 3: LLVM IR

**PER-THREAD CODE. TILE ABSTRACTION GONE. REGISTER ALLOCATION, LOOP UNROLLING, INLINE PTX FOR MASKED LOADS.**

```
1 define ptx_kernel void @add_kernel(ptr addrspace(1) %0, ptr addrspace(1) %1, ptr addrspace(1) %2, i32 %3, ptr addrspace(1) readnone captures(none) %4, ptr addrspace(1)
  readnone captures(none) %5) local_unnamed_addr #0 !dbg !4 {
2   %7 = tail call i32 @llvm.nvvm.read.ptx.sreg.ctaid.x(), !dbg !7
3   %8 = shl i32 %7, 8, !dbg !8
4   %9 = tail call i32 @llvm.nvvm.read.ptx.sreg.tid.x(), !dbg !9
5   %10 = and i32 %9, 127, !dbg !9
6   %11 = or disjoint i32 %8, %10, !dbg !10
7   %12 = or disjoint i32 %11, 128, !dbg !10
8   %13 = icmp slt i32 %11, %3, !dbg !11
9   %14 = icmp slt i32 %12, %3, !dbg !11
10  %15 = sext i32 %11 to i64, !dbg !12
11  %16 = getelementptr float, ptr addrspace(1) %0, i64 %15, !dbg !12
12  %17 = sext i32 %12 to i64, !dbg !12
13  %18 = getelementptr float, ptr addrspace(1) %0, i64 %17, !dbg !12
14  %19 = tail call i32 asm sideeffect "mov.u32 $0, $1;\0A\09@\3 ld.global.b32 { $0 }, [ $2 + 0 ];", "=r,r,l,b"(i32 0, ptr addrspace(1) %16, i1 %13) #2, !dbg !13
15  %20 = bitcast i32 %19 to float, !dbg !13
16  %21 = tail call i32 asm sideeffect "mov.u32 $0, $1;\0A\09@\3 ld.global.b32 { $0 }, [ $2 + 0 ];", "=r,r,l,b"(i32 0, ptr addrspace(1) %18, i1 %14) #2, !dbg !13
17  %22 = bitcast i32 %21 to float, !dbg !13
18  %23 = getelementptr float, ptr addrspace(1) %1, i64 %15, !dbg !14
19  %24 = getelementptr float, ptr addrspace(1) %1, i64 %17, !dbg !14
20  %25 = tail call i32 asm sideeffect "mov.u32 $0, $1;\0A\09@\3 ld.global.b32 { $0 }, [ $2 + 0 ];", "=r,r,l,b"(i32 0, ptr addrspace(1) %23, i1 %13) #2, !dbg !15
21  %26 = bitcast i32 %25 to float, !dbg !15
22  %27 = tail call i32 asm sideeffect "mov.u32 $0, $1;\0A\09@\3 ld.global.b32 { $0 }, [ $2 + 0 ];", "=r,r,l,b"(i32 0, ptr addrspace(1) %24, i1 %14) #2, !dbg !15
23  %28 = bitcast i32 %27 to float, !dbg !15
24  %29 = fadd float %20, %26, !dbg !16
25  %30 = fadd float %22, %28, !dbg !16
26  %31 = getelementptr float, ptr addrspace(1) %2, i64 %15, !dbg !17
27  %32 = getelementptr float, ptr addrspace(1) %2, i64 %17, !dbg !17
28  %33 = bitcast float %29 to i32, !dbg !18
29  tail call void asm sideeffect "@$2 st.global.b32 [ $1 + 0 ], { $0 };", "r,l,b"(i32 %33, ptr addrspace(1) %31, i1 %13) #2, !dbg !18
30  %34 = bitcast float %30 to i32, !dbg !18
31  tail call void asm sideeffect "@$2 st.global.b32 [ $1 + 0 ], { $0 };", "r,l,b"(i32 %34, ptr addrspace(1) %32, i1 %14) #2, !dbg !18
32  ret void, !dbg !19
33 }
34
35 ; Function Attrs: mustprogress nocallback nofree nosync nounwind speculatable willreturn memory(none)
36 declare noundef range(i32 0, 2147483647) i32 @llvm.nvvm.read.ptx.sreg.ctaid.x() #1
37
38 ; Function Attrs: mustprogress nocallback nofree nosync nounwind speculatable willreturn memory(none)
39 declare noundef range(i32 0, 1024) i32 @llvm.nvvm.read.ptx.sreg.tid.x() #1
40
41 attributes #0 = { nounwind "nvvm.reqntid"="128" }
42 attributes #1 = { mustprogress nocallback nofree nosync nounwind speculatable willreturn memory(none) }
43 attributes #2 = { nounwind }
```

# STEP 4: PTX — GPU

**PORTABLE ASSEMBLY.** THIS IS WHAT WE SAVE AND EMBED IN RUST.

**NOTE: 6 PARAMS, NOT 4 — TRITON ADDS 2 METADATA POINTERS!**

```
1 .version 8.7
2 .target sm_89
3 .address_size 64
4
5     // .globl  add_kernel          // -- Begin function add_kernel
6                                     // @add_kernel
7 .visible .entry add_kernel(
8     .param .u64 .ptr .global .align 1 add_kernel_param_0,
9     .param .u64 .ptr .global .align 1 add_kernel_param_1,
10    .param .u64 .ptr .global .align 1 add_kernel_param_2,
11    .param .u32 add_kernel_param_3,
12    .param .u64 .ptr .global .align 1 add_kernel_param_4,
13    .param .u64 .ptr .global .align 1 add_kernel_param_5
14 )
15 .reqntid 128
16 {
17     .reg .pred  %p<7>;
18     .reg .b32  %r<18>;
19     .reg .b64  %rd<11>;
20     .loc      1 6 0                // add_kernel_compile.py:6:0
21 $L__func_begin0:
22     .loc      1 6 0                // add_kernel_compile.py:6:0
23
24 // %bb.0:
25     ld.param.b64  %rd7, [add_kernel_param_0];
26     ld.param.b64  %rd8, [add_kernel_param_1];
27 $L__tmp0:
28     .loc      1 7 23                // add_kernel_compile.py:7:23
29
30     ...
```

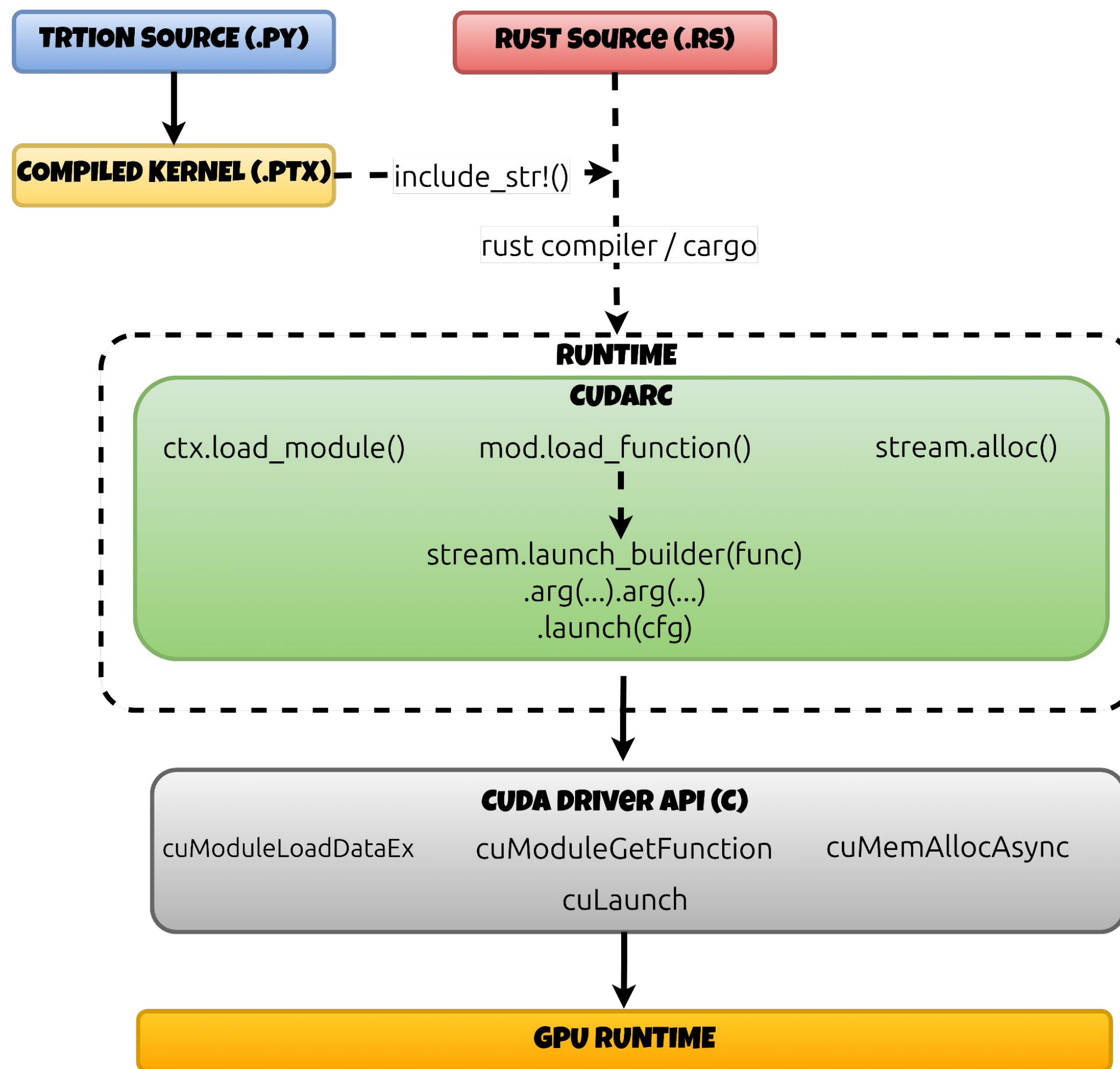
# STEP 5: SASS — MACHINE CODE

**GPU-SPECIFIC BINARY.** CUDA DRIVER COMPILES PTX -> SASS AT RUNTIME ON THE TARGET GPU. WE NEVER SHIP SASS, ONLY PTX.

```
1 ---
2 Function:add_kernel
3 ---:---:2 MOV R1, c[0x0][0x28];
4 ---:0:--:1 S2R R0, SR_TID.X;
5 ---:---:1 MOV R11, 0x4;
6 ---:---:1 CS2R R6, SRZ;
7 ---:---:1 ULDC.64 UR4, c[0x0][0x118];
8 ---:1:--:1 S2R R3, SR_CTAID.X;
9 ---:---:1 CS2R R8, SRZ;
10 01:--:Y:4 LOP3.LUT R0, R0, 0x7f, RZ, 0xc0, !PT;
11 02:--:Y:4 LEA R0, R3, R0, 0x8;
12 ---:---:1 ISETP.GE.AND P0, PT, R0.reuse, c[0x0][0x178], PT;
13 ---:---:1 IMAD.WIDE R4, R0.reuse, R11, c[0x0][0x168];
14 ---:---:Y:4 IADD3 R2, R0, 0x80, RZ;
15 ---:---:1 ISETP.GE.AND P1, PT, R2, c[0x0][0x178], PT;
16 ---:---:Y:6 IMAD.WIDE R2, R0, R11, c[0x0][0x160];
17 ---:2:--:4 @!P0 LDG.E R6, [R2.64];
18 ---:2:--:4 @!P0 LDG.E R7, [R4.64];
19 ---:3:--:4 @!P1 LDG.E R8, [R2.64+0x200];
20 ---:3:--:1 @!P1 LDG.E R9, [R4.64+0x200];
21 04:--:--:1 FADD R13, R7, R6;
22 ---:---:Y:5 IMAD.WIDE R6, R0, R11, c[0x0][0x170];
23 ---:0:--:1 @!P0 STG.E [R6.64], R13;
24 08:--:--:1 FADD R9, R9, R8;
25 ---:---:6 @P1 EXIT;
26 ---:---:1 STG.E [R6.64+0x200], R9;
27 ---:---:5 EXIT;
28 ...
```

# TRITON + RUST





# CUDARC: THE GPU BRIDGE

## CUDA Runtime API (libcudart)

- ✗ Requires nvcc compiler
- ✗ C++ ABI complexity in Rust FFI
- ✗ Runtime & driver must match

## CUDA Driver API (libcuda) via cudarc

- ✓ Load any .ptx at runtime
- ✓ No nvcc needed at deploy time
- ✓ Pure Rust, no unsafe in caller code
- ✓ Lazy module caching built-in

```
// CudaDevice wraps the driver handle
let device = CudaDevice::new(0)?; // GPU 0

// Alloc a typed buffer on the device
let buf: CudaSlice<f32> = unsafe { device.alloc(n)? };

// Load a PTX module and get a function handle (cached)
let func = device.get_or_load_custom_func(
    "kernel_name", "module_name", PTX_SOURCE
)?;
```

# OPTION 1: PURE CUDARC

Raw GPU buffers — minimal dependencies, full control:

```
1 const PTX: &str = include_str!("../kernels/matmul/compiled/matmul_tiled_kernel.ptx");
2
3 pub struct PureMatmulKernel {
4     stream: Arc<CudaStream>,
5     func:   CudaFunction,
6 }
7
8 impl PureMatmulKernel {
9     pub fn new(stream: Arc<CudaStream>) -> Result<Self> {
10         let func = stream.context().get_or_load_custom_func(
11             "matmul_tiled_kernel", "matmul_module", PTX)?;
12         Ok(Self { stream, func })
13     }
14
15     pub fn forward(&self, a: &CudaSlice<f32>, b: &CudaSlice<f32>,
16                 m: usize, k: usize, n: usize) -> Result<CudaSlice<f32>> {
17         let mut c = unsafe { self.stream.alloc:::<f32>(m * n)? };
18         let cfg = LaunchConfig {
19             grid_dim:      (m.div_ceil(32) as u32, n.div_ceil(32) as u32, 1),
20             block_dim:     (128, 1, 1), // matches .reqntid 128 in PTX
21             shared_mem_bytes: 16384, // 2 * (BLOCK_MxBLOCK_K + BLOCK_KxBLOCK_N) * 4
22         };
23         unsafe {
24             self.stream.launch_builder(&self.func)
25                 .arg(a).arg(b).arg(&mut c)
26                 .arg(&(m as u32)).arg(&(n as u32)).arg(&(k as u32))
27                 .arg(&0u64).arg(&0u64) // ← Triton metadata (must not be omitted!)
28                 .launch(cfg)?;
29         }
30         Ok(c)
31     }
32 }
```

# OPTION 2: CANDLE CUSTOMOP

Type-safe tensors — shape tracking, layout offsets, device binding handled automatically:

```
1 pub struct MatmulKernel;
2
3 impl MatmulKernel {
4     pub fn forward(&self, a: &Tensor, b: &Tensor) -> Result<Tensor> {
5         a.apply_op2_no_bwd(b, self)
6     }
7 }
8
9 impl CustomOp2 for MatmulKernel {
10     fn name(&self) -> &'static str { "matmul_tiled" }
11
12     fn cuda_fwd(&self,
13         a_s: &CudaStorage, a_l: &Layout,
14         b_s: &CudaStorage, b_l: &Layout,
15     ) -> Result<(CudaStorage, Shape)> {
16         let dev = a_s.device();
17         let (m, k) = (a_l.dim(0)?, a_l.dim(1)?);
18         let n = b_l.dim(1)?;
19
20         let a = f32::as_cuda_slice(a_s)?.slice(a_l.start_offset()..);
21         let b = f32::as_cuda_slice(b_s)?.slice(b_l.start_offset()..);
22         let mut c = unsafe { dev.alloc::<f32>(m * n)? };
23
24         let func = dev.get_or_load_custom_func(
25             "matmul_tiled_kernel", "matmul_module", PTX?);
26         let cfg = LaunchConfig {
27             grid_dim:      (m.div_ceil(32) as u32, n.div_ceil(32) as u32, 1),
28             block_dim:     (128, 1, 1),
29             shared_mem_bytes: 16384,
30         };
31         unsafe {
32             dev.launch_builder(&func)
33                 .arg(&a).arg(&b).arg(&mut c)
34                 .arg(&(m as u32)).arg(&(n as u32)).arg(&(k as u32))
35                 .arg(&0u64).arg(&0u64)
36                 .launch(cfg)?;
37         }
38         Ok((f32::wrap_cuda_slice(c, dev.clone()), Shape::from(&[m, n])))
39     }
40 }
```

# PURE CUDARC VS CANDLE CUSTOMOP

## Pure cudarc

- ✓ Minimal dependencies
- ✓ Direct buffer control
- ✓ Works outside Candle ecosystem
- ⚠ Manual shape tracking
- ⚠ No layout offset handling
- ⚠ Returns `CudaSlice`, not `Tensor`

## Candle CustomOp

- ✓ Automatic shape tracking
- ✓ `start_offset()` handles sliced tensors
- ✓ Device binding & grad graph
- ✓ Integrates with Candle model API
- ⚠ Candle dependency required

# BENCHMARK



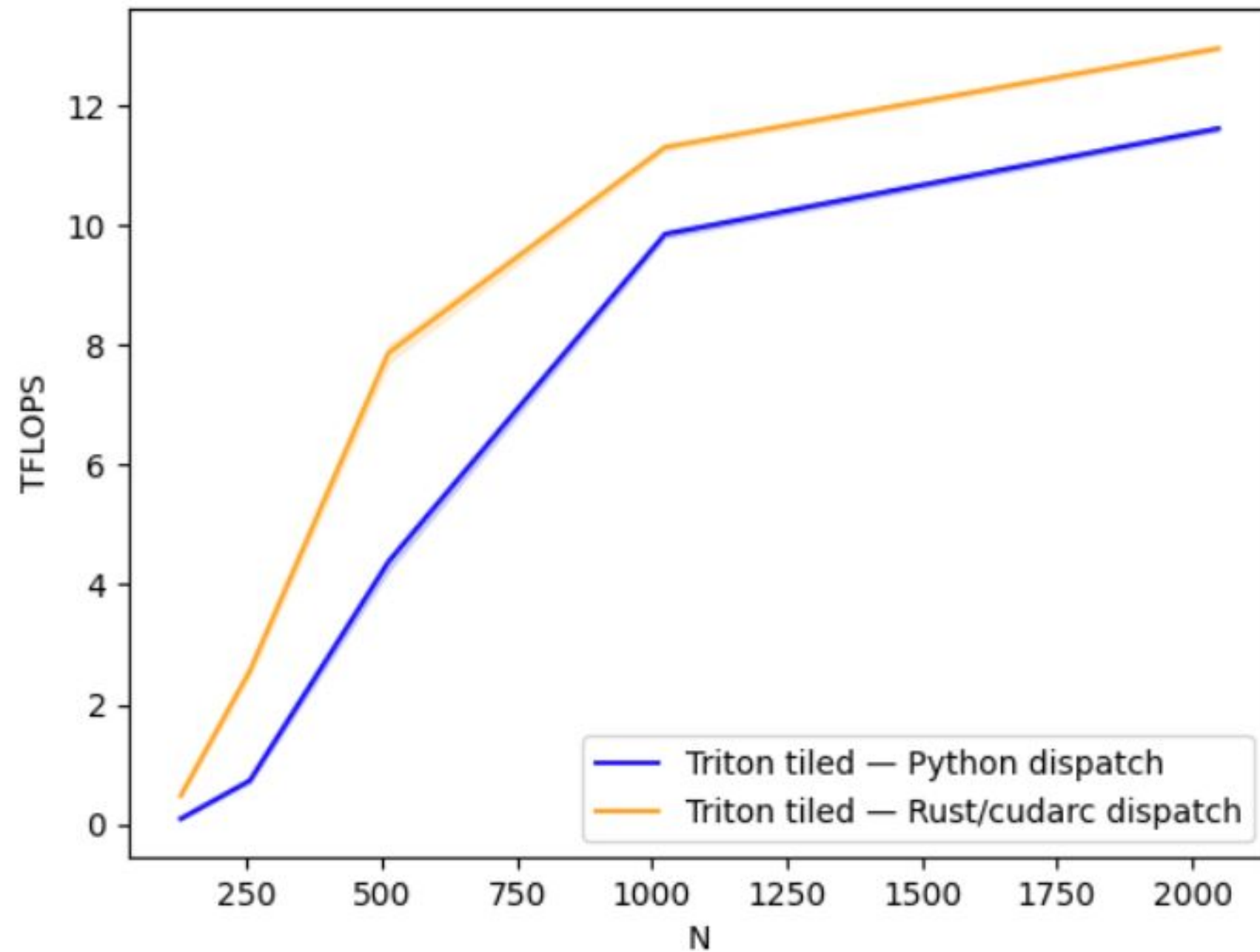
TRITON + **PYTHON** (RUBATO)



TRITON + **RUST** (PRESTO CON FUOCO)

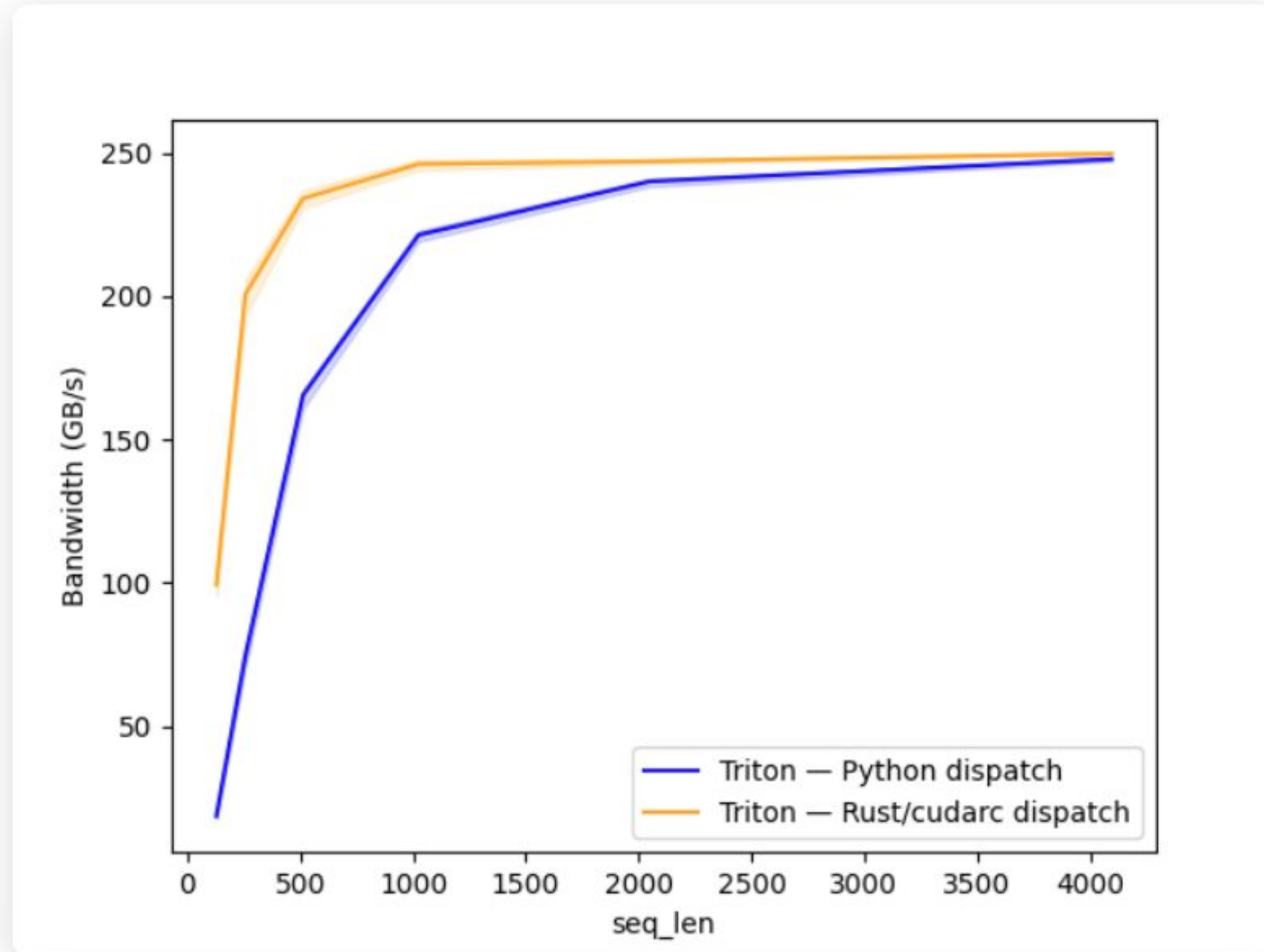
# BENCHMARK: TILED MATMUL

TFLOPS (fp32) — Python vs Rust · RTX 4060 Ti



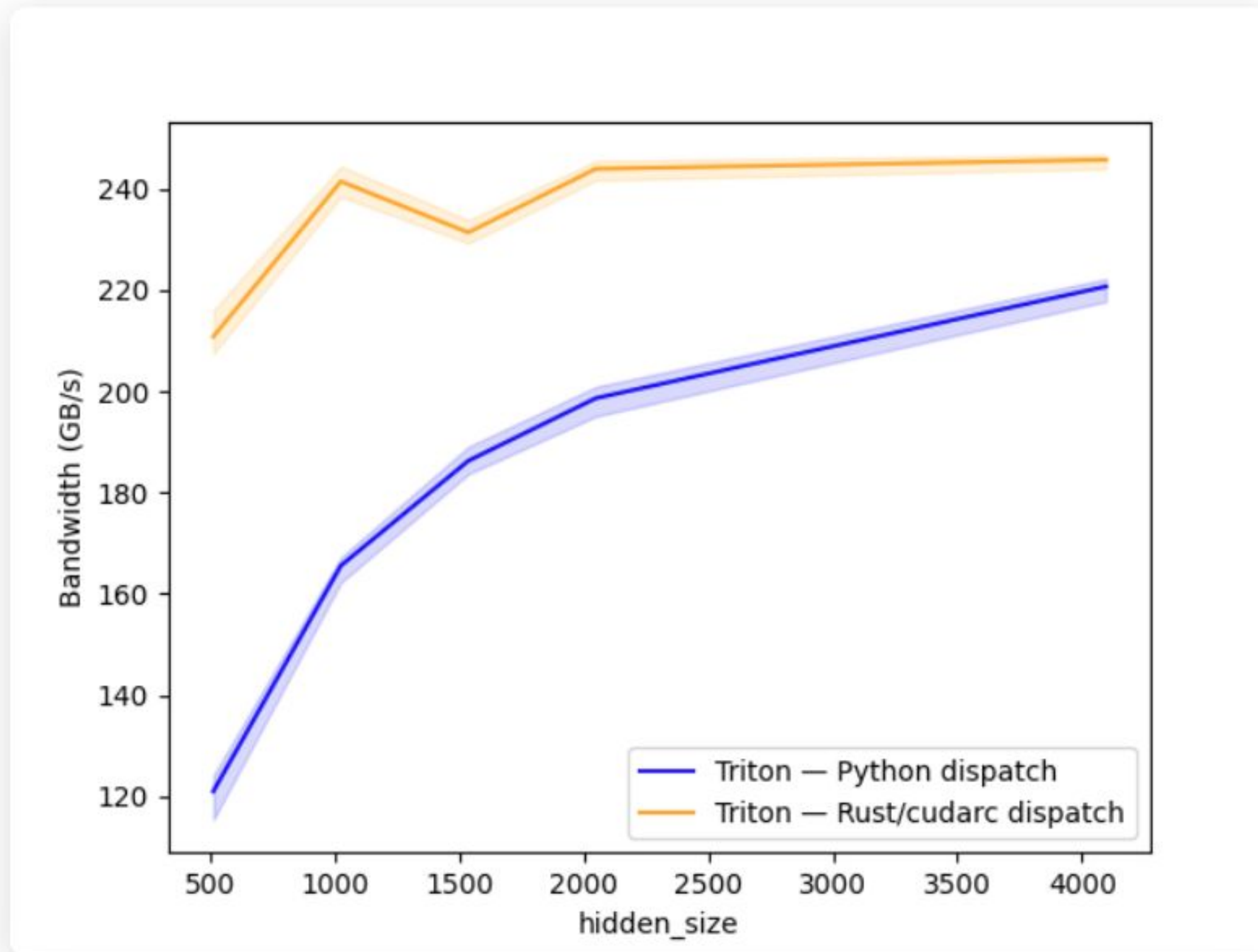
# BENCHMARK: CAUSAL SOFTMAX

Cold-cache bandwidth (GB/s) — Python vs Rust · RTX 4060 Ti



# BENCHMARK: RMSNORM

Cold-cache bandwidth (GB/s) — Python vs Rust · RTX 4060 Ti



# BENCHMARKS (RTX 4060 TI)

Kernel	Size	Python dispatch	Rust/cudarc dispatch	Notes
Causal Softmax	seq_len=128	18.6 GB/s	99.3 GB/s (5.3×)	Dispatch overhead dominates
Causal Softmax	seq_len=2048	240 GB/s	247 GB/s (~1.0×)	DRAM-bound — both hit ~250 GB/s ceiling
RMSNorm	hidden=1536 (Bielik)	186 GB/s	231 GB/s	<b>Dip: next_pow2(1536)=2048 → 25% wasted BW</b>
RMSNorm	hidden=2048	199 GB/s	244 GB/s	No padding waste — full bandwidth
Tiled MatMul	N=128	0.09 TFLOPS	0.48 TFLOPS (5.3×)	Dispatch overhead dominates at small N
Tiled MatMul	N=2048	11.6 TFLOPS	13.0 TFLOPS (1.12×)	TF32 tensor cores — compute-bound

**Small inputs:** Rust dispatch 5× faster  
(saves 35–80 μs Python overhead)

**Large inputs:** kernels converge —  
DRAM/compute is the ceiling, not dispatch

# WHEN IS THIS STACK THE RIGHT CHOICE?

## Use it when...

- ✓ Serving a *fixed* model on *known* hardware
- ✓ Latency SLAs in the single-digit millisecond range
- ✓ Memory-bound decoding (every microsecond matters)
- ✓ Agentic loops — many sequential kernel launches
- ✓ Edge / embedded deployment (no Python runtime)
- ✓ You need custom ops Candle doesn't have

## Skip it when...

- ⚠️ Rapid architecture experimentation (use PyTorch)
- ⚠️ Model changes frequently (recompile PTX each time)
- ⚠️ Team unfamiliar with Rust + GPU programming
- ⚠️ Standard cuBLAS/NCCL kernels are sufficient
- ⚠️ Many different GPU SKUs in production fleet

# SUMMARY



**PRESENTATION**

**THE END**

