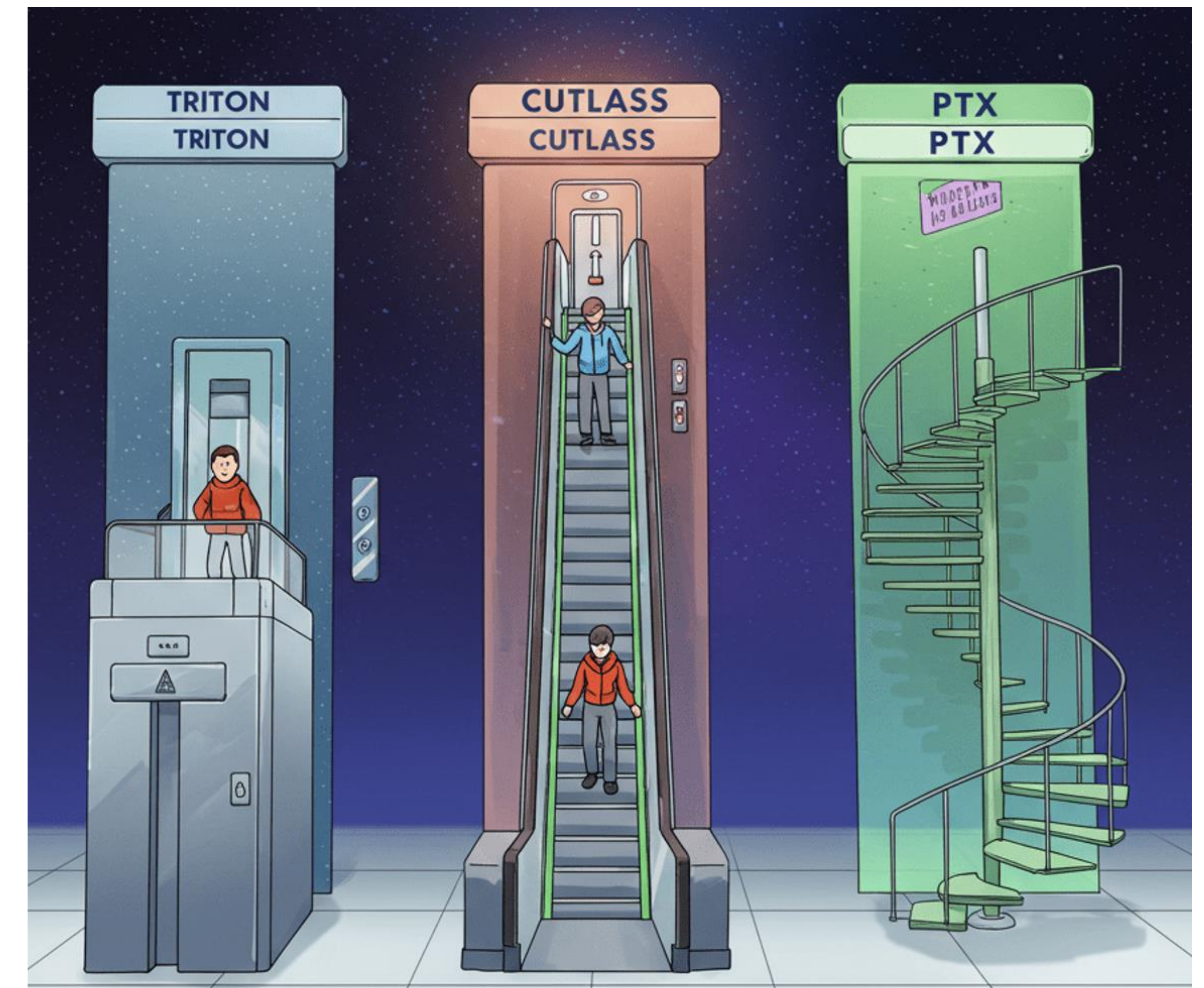**NVIDIA**

# Enable Tensor Core Programming in Python with CUTLASS 4.0

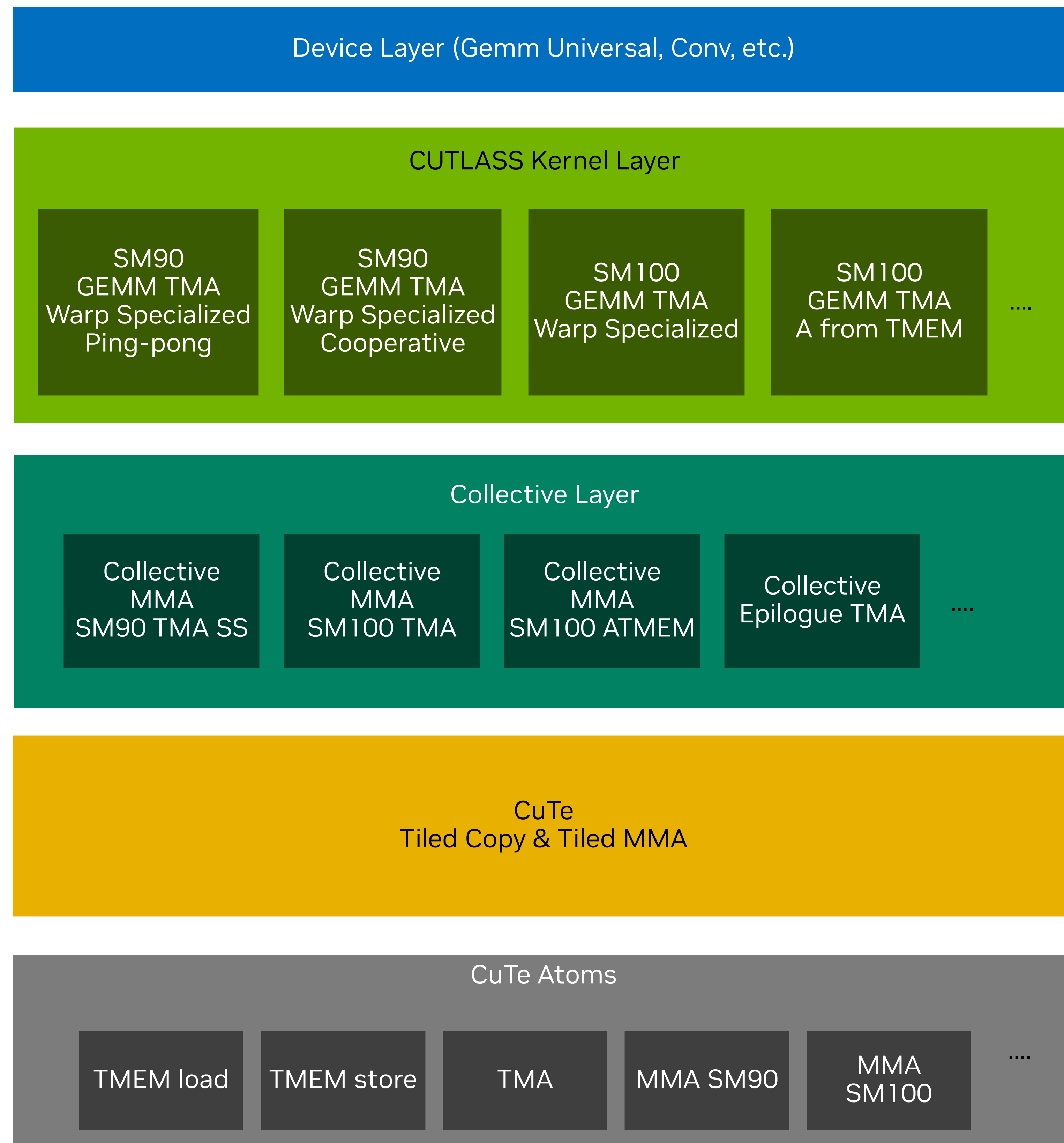Kihiro Bando, Brandon Sun | 2025-03-21

# Why CUTLASS?

## Enabling innovation for SOL performance

- High-level generators leveraging compilers are popular
  - Hide and automate a lot of details
  - Get excellent performance on common use cases, but...
  - Algorithmic innovations require lower-level **abstractions**
  - Advanced features like PDL require fine grain control

- With CUTLASS
  - Available on day 0 with full control!
  - Expressive abstractions for performance in all cases
  - Modular and extensible design robust throughout GPU generations

# CUTLASS

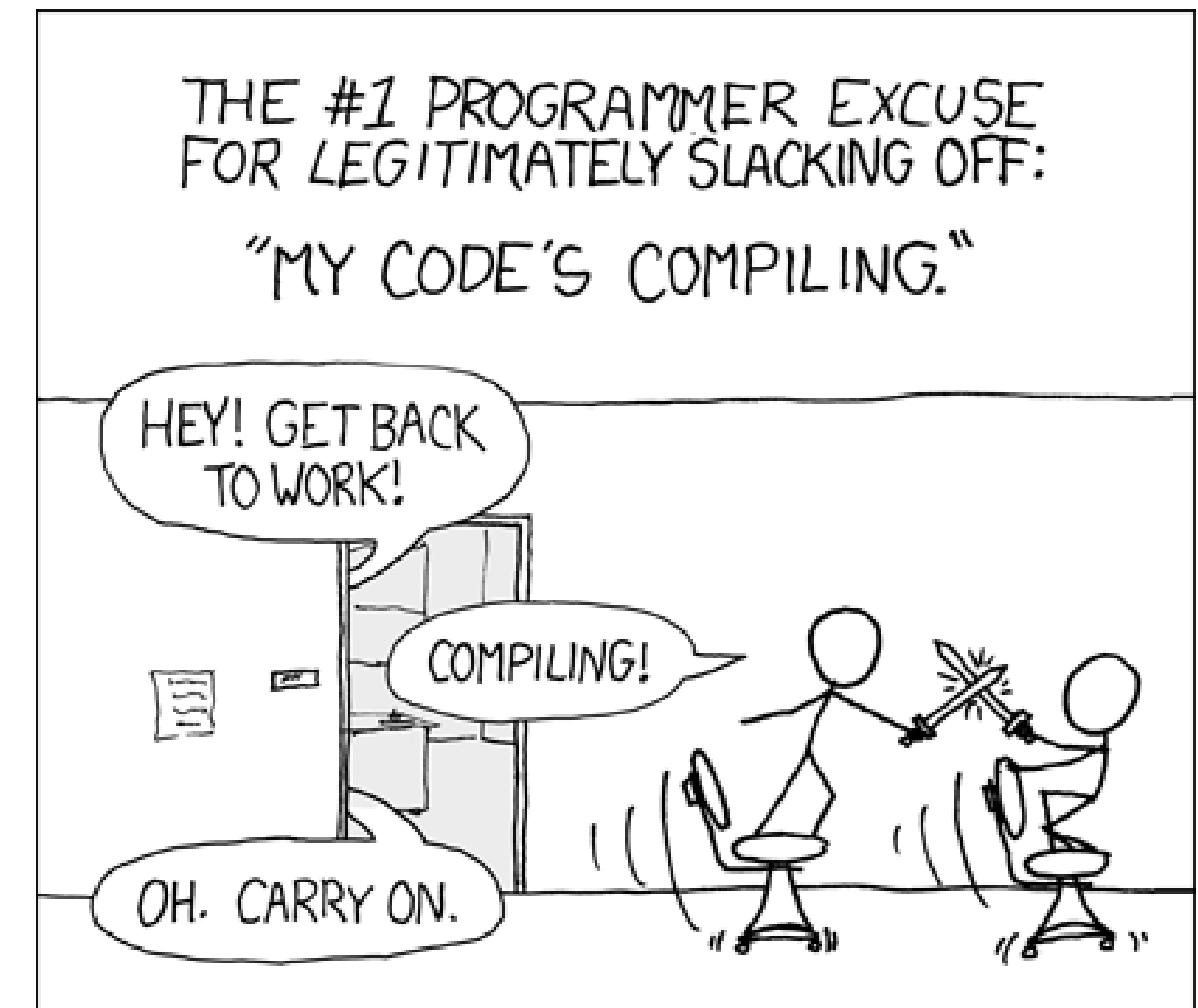A set of useful abstractions for productivity and performance at all scopes and scales

**Device Layer (Gemm Universal, Conv, etc.)**

**CUTLASS Kernel Layer**

| SM90 GEMM TMA Warp Specialized Ping-pong | SM90 GEMM TMA Warp Specialized Cooperative | SM100 GEMM TMA Warp Specialized | SM100 GEMM TMA A from TMEM | .... |

**Collective Layer**

| Collective MMA SM90 TMA SS | Collective MMA SM100 TMA | Collective MMA SM100 ATMEM | Collective Epilogue TMA | .... |

**CuTe Tiled Copy & Tiled MMA**

**CuTe Atoms**

| TMEM load | TMEM store | TMA | MMA SM90 | MMA SM100 | .... |

**More pre-tuned recipes**

**More Control**

- Open source https://github.com/NVIDIA/cutlass
- Presented: GTC'18, GTC'19, GTC'20, GTC'21, GTC'22, GTC'22, GTC'23, GTC'24
- Multiple entry points depending on your needs
- More details this year in *Programming Blackwell Tensor Cores with CUTLASS* [S72720]

# Major pain points with C++

C++ templates and unfortunate consequences

- C++ templates suffer from slow compilation time
  - Front-end too generic for our purposes
  - Prevents fast iteration
  - Prohibits JIT-ting at scale and brute force auto-tuning
- C++ templates are inconvenient
  - Additional mental load when writing compile-time logic
  - Error messages are longer than novels
- The DL space fully embraces the python ecosystem regardless
  - Everybody hates writing binding code
  - Dependency on `nvcc`
- LLMs are likely better at generating python programs



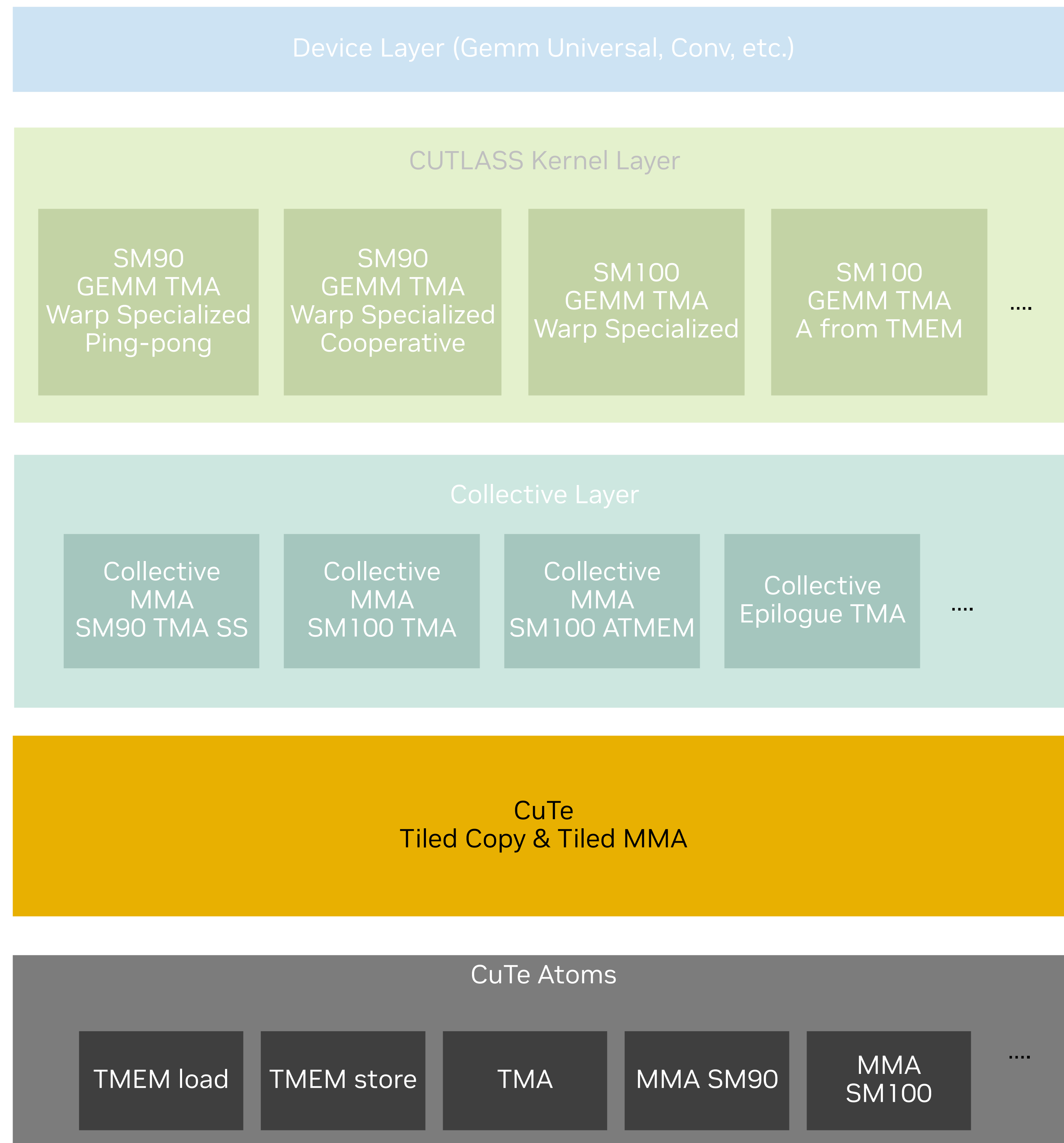Do I have to tolerate all of this to use CUTLASS?

# Introducing with CUTLASS 4.0
## Tensor core programming in Python

# CUTLASS in Python

## Initial launch to include CuTe

**Device Layer (Gemm Universal, Conv, etc.)**

**CUTLASS Kernel Layer**

| SM90 GEMM TMA Warp Specialized Ping-pong | SM90 GEMM TMA Warp Specialized Cooperative | SM100 GEMM TMA Warp Specialized | SM100 GEMM TMA A from TMEM | .... |

**Collective Layer**

| Collective MMA SM90 TMA SS | Collective MMA SM100 TMA | Collective MMA SM100 ATMEM | Collective Epilogue TMA | .... |

**CuTe**
**Tiled Copy & Tiled MMA**

**CuTe Atoms**

| TMEM load | TMEM store | TMA | MMA SM90 | MMA SM100 | .... |

More pre-tuned recipes

More Control

This first release makes available a *mature* low-level tensor programming model, giving access to tensor cores with full control.

More to come later...

- `make_shape(Int<1>{}, Int<2>{}, x) → (1,2,x)`
- `make_layout`
- `make_identity_tensor`
- `zipped_divide`
- `local_tile`
- `tiled_mma.get_slice`
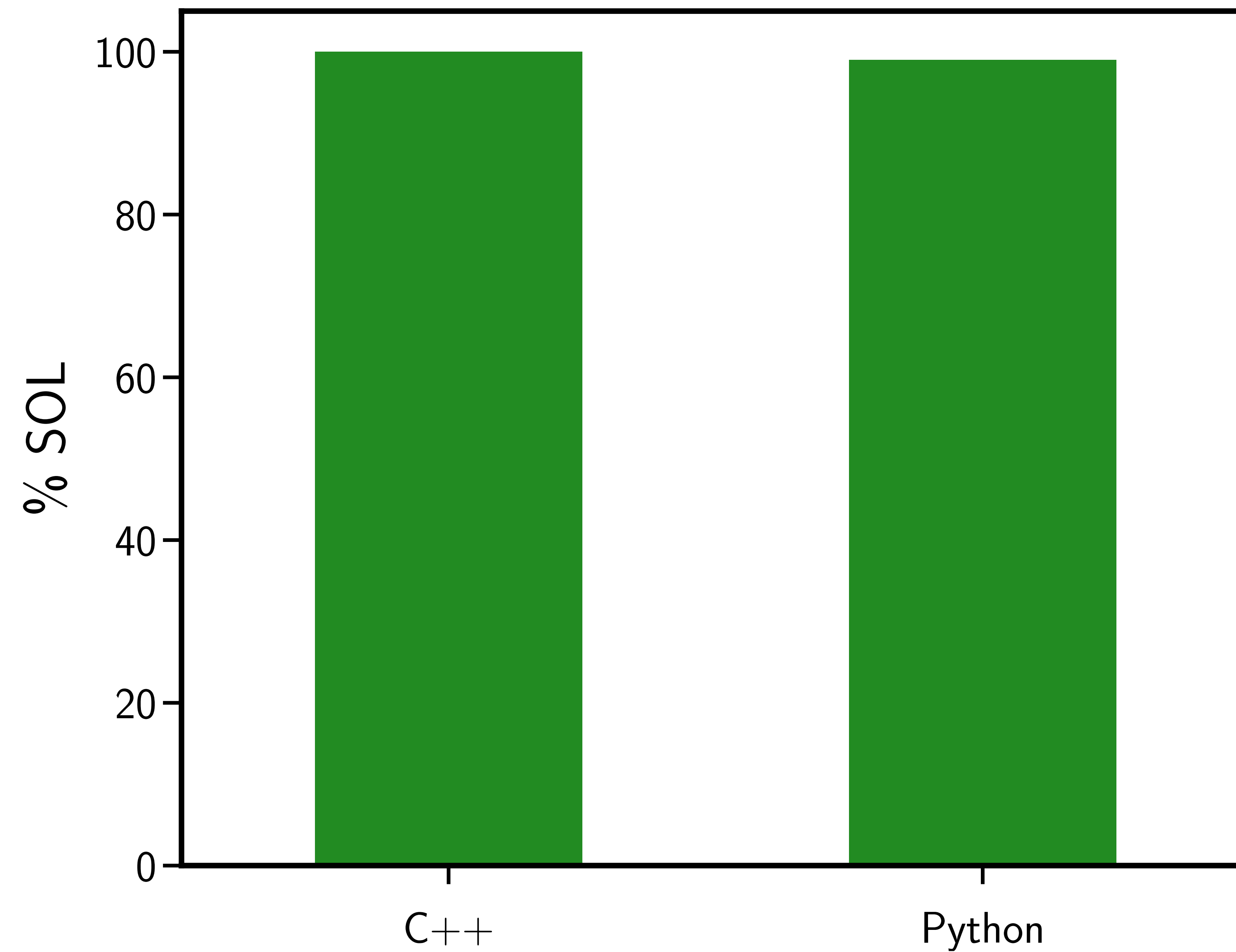- `thr_mma.partition_A`
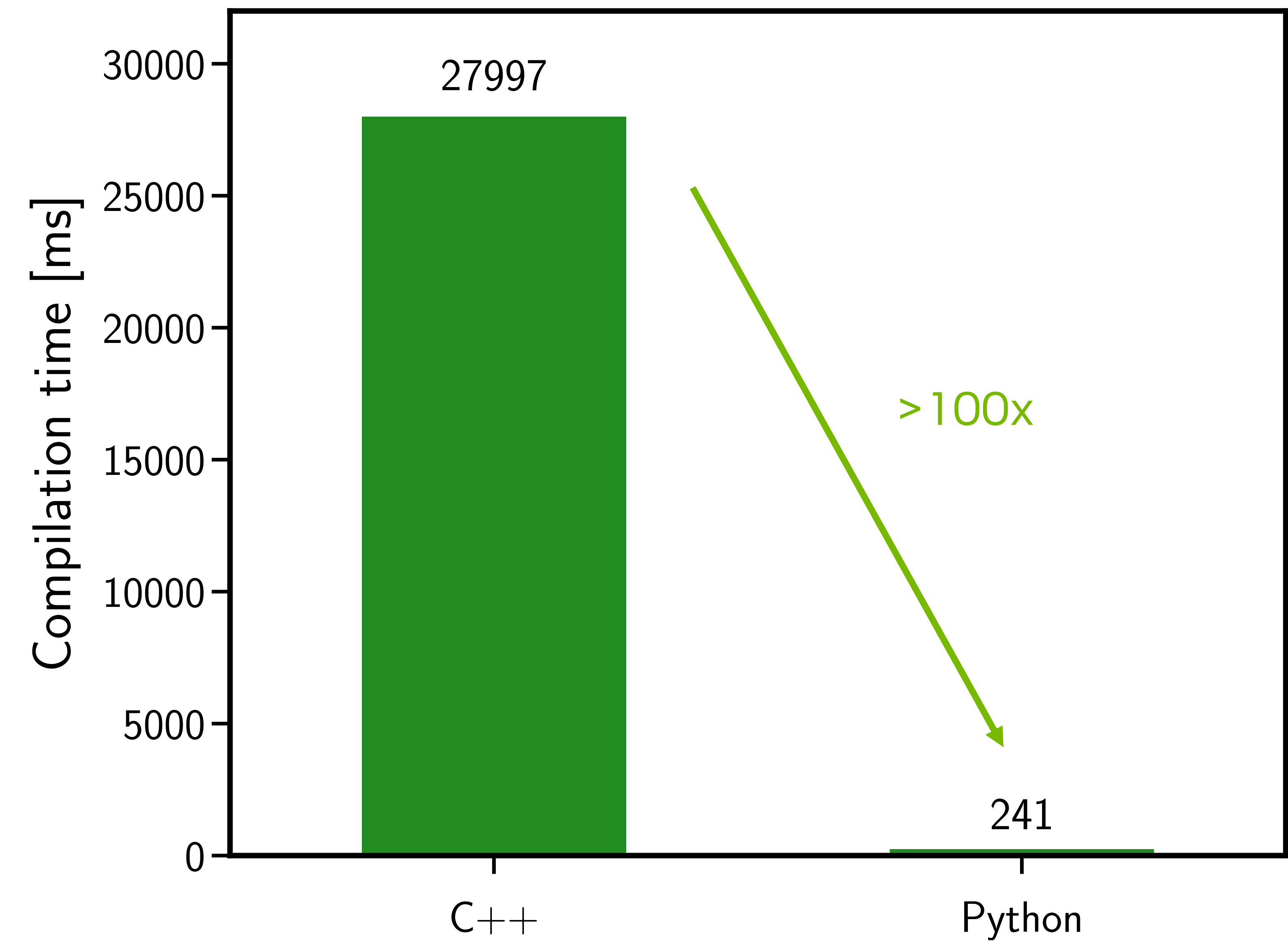- `thr_copy.partition_S`

# CUTLASS in Python

What do you get?

**8kx8kx8k GEMM**

**Peak Performance!**

**Blazing Fast Compilation Time!**



>100x

Higher is better

Lower is better

# CUTLASS in Python

## How will you get started?

`pip install nvidia-cutlass-dsl`

```python
import cutlass
import cutlass.cute as cute

@cute.kernel
def kernel():
  tidx, _, _ = cutlass.nvvm.thread_idx()
  if tidx == 0:
    cute.print_("Hello world")

@cute.jit
def host():
  kernel(config=cutlass.LaunchConfig(
    grid=(1, 1, 1), block=(32, 1, 1)))

host()
```

`python3 hello_world.py`

# Agenda

- Introduction and Motivations
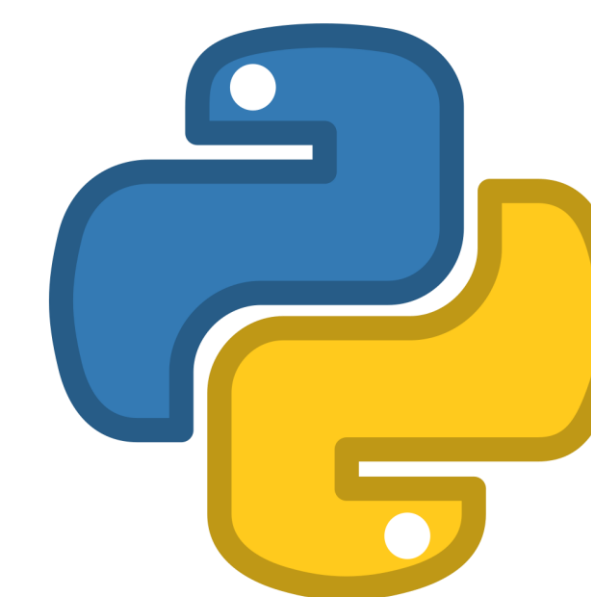
---

- **The DSL Infrastructure**

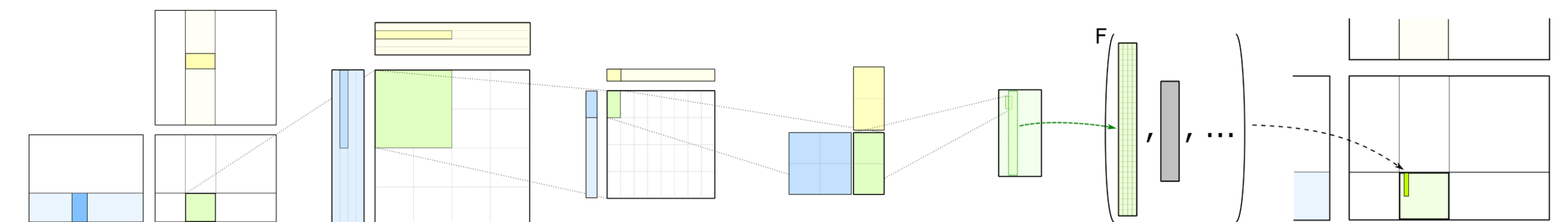---

- Kernel Authoring in Python with CuTe

---

- Runtime Performance
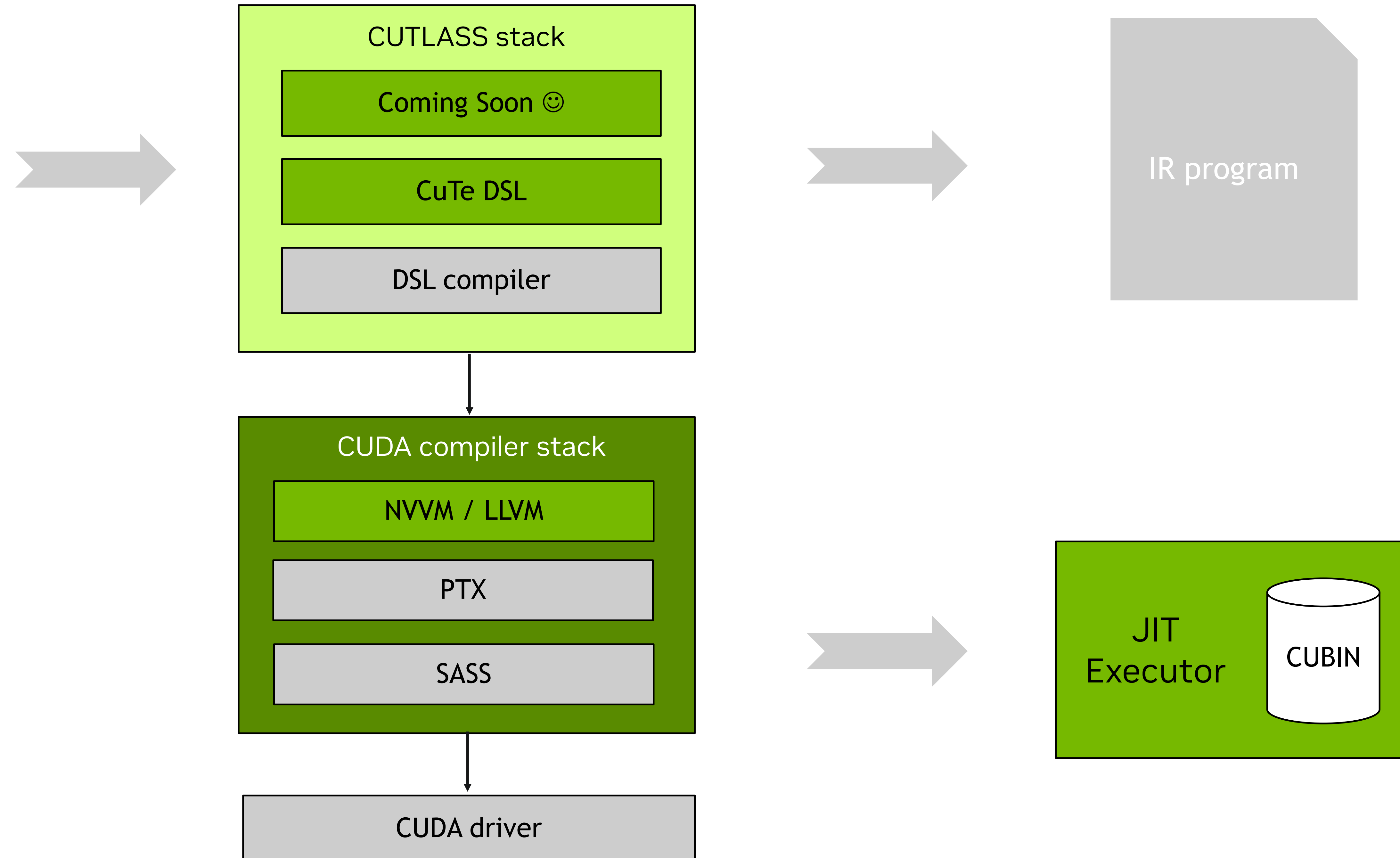
---

- Conclusion

# Soul of CuTe DSL

- **A Python based programming language to program Tensor Cores with CuTe semantics for best possible performance**

- Enables kernel authoring in Python rather than just accessing CUTLASS kernels

- Empowered by CuTe abstractions

- Easy integration with popular Python frameworks, like Pytorch

- Model hardware accurately for full control of performance

- Based on MLIR framework to leverage the power of MLIR ecosystem
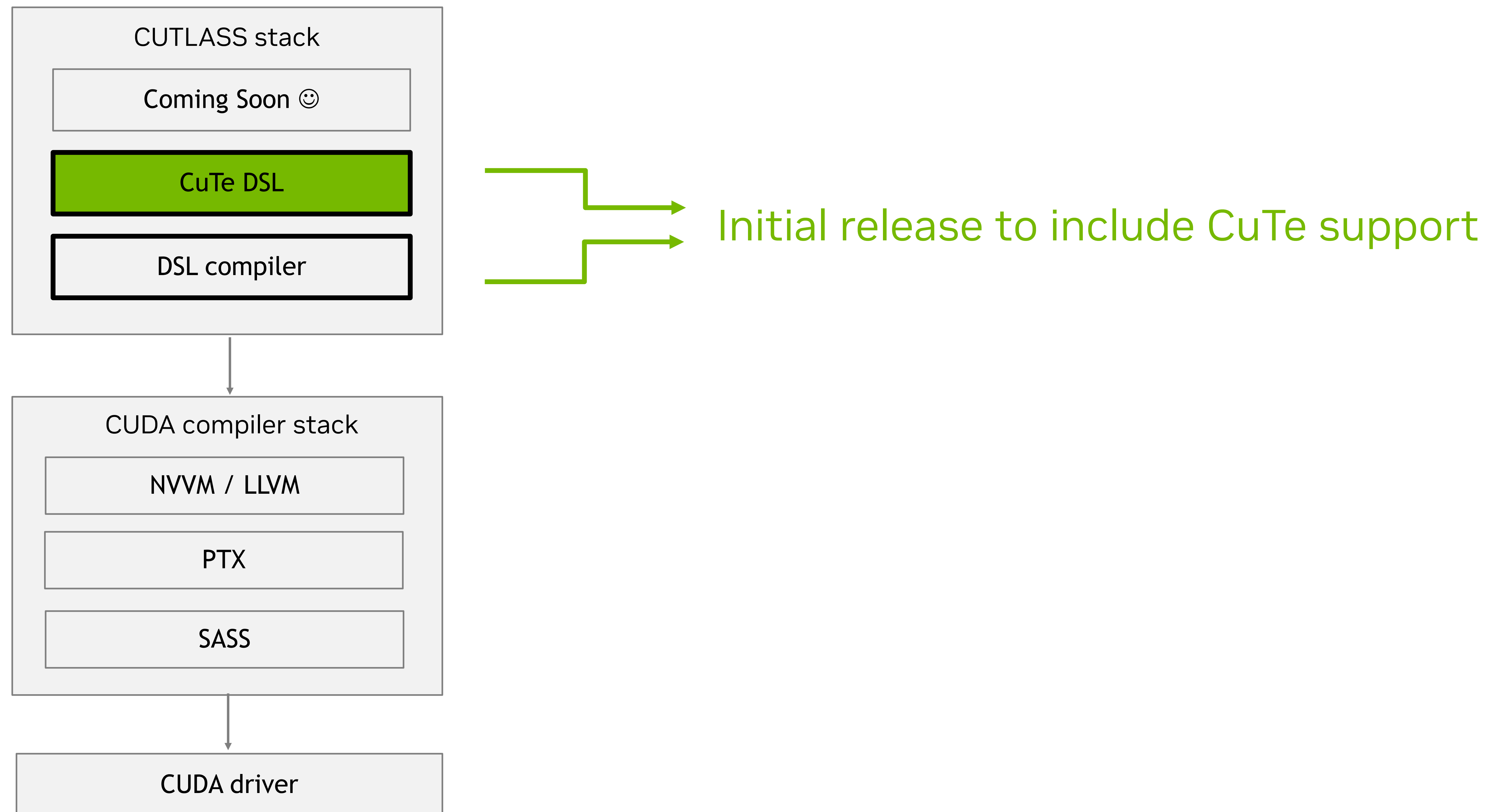
# CUTLASS Python architecture

```
@cute.kernel
def my_kernel(A: cute.Tensor, ...):
    ...
    atom = cute.make_copy_atom(...)
    ...
    cute.make_tiled_copy_tv(...)
    ...
```

**CUTLASS stack**

Coming Soon ☺

CuTe DSL

DSL compiler

IR program

**CUDA compiler stack**

NVVM / LLVM

PTX

SASS

CUDA driver

JIT Executor  CUBIN

# Tensor Core programming in Python

```
CUTLASS stack
  ┌─────────────────────┐
  │   Coming Soon ☺     │
  └─────────────────────┘
  ┌─────────────────────┐
  │      CuTe DSL       │
  └─────────────────────┘
  ┌─────────────────────┐
  │    DSL compiler     │
  └─────────────────────┘
```

Initial release to include CuTe support

```
CUDA compiler stack
  ┌─────────────────────┐
  │     NVVM / LLVM     │
  └─────────────────────┘
  ┌─────────────────────┐
  │         PTX         │
  └─────────────────────┘
  ┌─────────────────────┐
  │        SASS         │
  └─────────────────────┘
```

```
  ┌─────────────────────┐
  │     CUDA driver     │
  └─────────────────────┘
```

# Writing a kernel in Python
## @cute.jit / @cute.kernel

```cpp
template <class ProblemShape, class CtaTiler,
          class TA, class SmemLayoutA, class TmaA,
          class TB, class SmemLayoutB, class TmaB,
          class TC, class CStride, class TiledMma,
          class Alpha, class Beta>
__global__ static
__launch_bounds__(decltype(size(TiledMma{}))::value)
void
gemm_device(ProblemShape shape_MNK, CtaTiler cta_tiler,
            TA const* A, CUTLASS_GRID_CONSTANT TmaA const tma_a,
            TB const* B, CUTLASS_GRID_CONSTANT TmaB const tma_b,
            TC      * C, CStride dC, TiledMma mma,
            Alpha alpha, Beta beta) {

  ...


// Kernel Launch
cutlass::Status status = cutlass::launch_kernel_on_cluster(
                                    params, kernel_ptr,
                                    prob_shape, cta_tiler,
                                    A, tmaA,
                                    B, tmaB,
                                    C, dC, tiled_mma,
                                    alpha, beta);
```

```python
import cutlass
import cutlass.cute as cute

@cute.kernel
def kernel(self, tma_atom_a: cute.CopyAtom, mA_mkl: cute.Tensor,
           tma_atom_b: cute.CopyAtom, mB_nkl: cute.Tensor,
           mC_mnl: cute.Tensor,
           cluster_layout_vmnk: cute.Layout,
           a_smem_layout_staged: cute.Layout,
           b_smem_layout_staged: cute.Layout,
           epilogue_op: cutlass.Constexpr = lambda x: x,
):
    ...


@cute.jit
def __call__(self, mA: cute.Tensor, mB: cute.Tensor, mC: cute.Tensor,
             epilogue_op=lambda x: x,
):
    ...

    # Launch the kernel
    self.kernel(...)
```

# Easy integration: Interop with Pytorch

Support DLPack protocol

```python
import cutlass
import cutlass.cute as cute
from cutlass.cute.runtime import from_dlpack
import torch


@cute.kernel
def jit_kernel(A: cute.Tensor):
    ...


@cute.jit
def jit_func(A: cute.Tensor):
    jit_kernel(
        A, config=cutlass.LaunchConfig(grid=[1, 1, 1], block=[1, 1,
1])
    )
```

```python
A_tensor = torch.tensor([0, 0, 0], dtype=torch.int32).cuda()
jit_func(A_tensor)
# Or jit_func(from_dlpack(A_tensor).mark_layout_dynamic())
```

- Take torch.tensor as input seamlessly
- Finer grained control with explicit call

# Easy integration: Interop with Pytorch
## Support DLPack protocol

```cpp
// Create instantiation for device reference gemm kernel
cutlass::reference::device::Gemm<ElementInputA,
                                 LayoutInputA,
                                 ElementInputB,
                                 LayoutInputB,
                                 ElementOutput,
                                 LayoutOutput,
                                 ElementComputeEpilogue,
                                 ElementComputeEpilogue>
    gemm_device;

// Launch device reference gemm kernel
gemm_device(problem_size,
            alpha,
            tensor_a.device_ref(),
            tensor_b.device_ref(),
            beta,
            tensor_c.device_ref(),
            tensor_ref_d.device_ref());

// Wait for kernels to finish
cudaDeviceSynchronize();

// Copy output data from CUTLASS and reference kernel to host for comparison
tensor_d.sync_host();
tensor_ref_d.sync_host();

// Check if output from CUTLASS kernel and reference kernel are equal or not
bool passed = cutlass::reference::host::TensorEquals(
  tensor_d.host_view(),
  tensor_ref_d.host_view());

std::cout << (passed ? "Passed" : "Failed") << std::endl;
```

```python
ref_c = (torch.einsum("mkl,nkl->mnl", a_ref, b_ref)).cpu()
torch.testing.assert_close(gpu_c, ref_c, atol=tolerance, rtol=1e-05)
```

# Easy integration: Interop with Pytorch

From static layout to dynamic layout

```python
@cute.kernel
def jit_kernel(A: cute.Tensor, x: cutlass.Int32, y: cutlass.Int32):
    A[x] = y


@cute.jit
def jit_func(A: cute.Tensor):
    x = 0
    y = 3
    jit_kernel(
        A, x, y, config=cutlass.LaunchConfig(grid=[1, 1, 1],
                                             block=[1, 1, 1])

    )


A_tensor = torch.tensor([0, 0, 0], dtype=torch.int32).cuda()
jit_func(from_dlpack(A_tensor))
```

- cute.Tensor type would have a layout (3:1) that follows the size of input A_tensor

# Easy integration: Interop with Pytorch

From static layout to dynamic layout

```python
@cute.kernel
def jit_kernel(A: cute.Tensor, x: cutlass.Int32, y: cutlass.Int32):
    A[x] = y


@cute.jit
def jit_func(A: cute.Tensor):
    x = 0
    y = 3
    jit_kernel(
        A, x, y, config=cutlass.LaunchConfig(grid=[1, 1, 1],
                                             block=[1, 1, 1])
    )


A_tensor = torch.tensor([0, 0, 0], dtype=torch.int32).cuda()
jit_func(from_dlpack(A_tensor))

B_tensor = torch.tensor([0, 0, 0, 0, 0], dtype=torch.int32).cuda()
jit_func(from_dlpack(B_tensor))
```

- cute.Tensor type would have a layout (5:1) that follows the size of input B_tensor

- Two sets of JIT functions are compiled:
  - One with cute.Tensor type of layout (3:1)
  - One with cute.Tensor type of layout (5:1)

- Static layout results in distinct codes

17  NVIDIA.

# Easy integration: Interop with Pytorch

From static layout to dynamic layout

```python
@cute.kernel
def jit_kernel(A: cute.Tensor, x: cutlass.Int32, y: cutlass.Int32):
    A[x] = y


@cute.jit
def jit_func(A: cute.Tensor):
    x = 0
    y = 3
    jit_kernel(
        A, x, y, config=cutlass.LaunchConfig(grid=[1, 1, 1],
                                             block=[1, 1, 1])
    )


A_tensor = torch.tensor([0, 0, 0], dtype=torch.int32).cuda()
jit_func(from_dlpack(A_tensor).mark_layout_dynamic(mode=[0]))

B_tensor = torch.tensor([0, 0, 0, 0, 0], dtype=torch.int32).cuda()
jit_func(B_tensor)
```

- cute.Tensor type would have a dynamic layout (?:1)

- Only one set of JIT functions compiled for both cases

- Use dynamic layout to allow generalized code generation

18

# Easy integration: Interop with Pytorch
## Integration with LLaMA 8b

- Wire up the gate/up/down projection layer with the customized linear module

```python
class LlamaMLP(nn.Module):
    def __init__(self, config):
        super().__init__()
        self.config = config
        self.hidden_size = config.hidden_size
        self.intermediate_size = config.intermediate_size
        self.gate_proj = nn.Linear(
            self.hidden_size,
            self.intermediate_size,
            bias=config.mlp_bias
        )
        self.up_proj = nn.Linear(
            self.hidden_size,
            self.intermediate_size,
            bias=config.mlp_bias
        )
        self.down_proj = nn.Linear(
            self.intermediate_size,
            self.hidden_size,
            bias=config.mlp_bias
        )
        self.act_fn = ACT2FN[config.hidden_act]
```

```python
class LlamaMLP(nn.Module):
    def __init__(self, config):
        super().__init__()
        self.config = config
        self.hidden_size = config.hidden_size
        self.intermediate_size = config.intermediate_size
        self.gate_proj = MyCutlassLinear(
            self.hidden_size,
            self.intermediate_size,
            bias=config.mlp_bias
        )
        self.up_proj = MyCutlassLinear(
            self.hidden_size,
            self.intermediate_size,
            bias=config.mlp_bias
        )
        self.down_proj = MyCutlassLinear(
            self.intermediate_size,
            self.hidden_size,
            bias=config.mlp_bias
        )
        self.act_fn = ACT2FN[config.hidden_act]
```

# Easy integration: Interop with Pytorch

Integration with LLaMA 8b

```python
class MyCutlassLinear(nn.Module):
    def __init__(self, in_features, out_features, bias=False):
        ...
        from blackwell.gemm import MyGemmKernel

        self.gemm = MyGemmKernel(
            cutlass.Float16,
            cutlass.Float32,
            cutlass.Float16,
            False,              # 2-CTA optimization
            (128, 128),         # MMA tile shape
            (2, 1, 1),          # cluster shape
            True)               # Use TMA


    def forward(self, input, bias=None):
        batch_size, seq_len, hidden_size = input.shape

        try:
            input = input.reshape(batch_size * seq_len, hidden_size)

            output = torch.empty(
                input.size(0), self.out_features, device=input.device,
                                                  dtype=input.dtype)

            ...

            self.gemm(
                input.detach().contiguous(),
                weight.detach().contiguous(),
                output.contiguous(),
                stream)
            ...
```

- Setup the linear module with your customized kernel implemented by CUTLASS Python APIs

- Invoke your kernel for the forward pass
- Leverage implicit from_dlpack conversion

NVIDIA.

# Metaprogramming with an imperative style
## Python to Python code generation: dynamic if

- Pythonic way to write "meta-kernel" that automatically fits for dynamic layout
  - Conditional execution based on dynamic expression

```python
@cute.kernel
def jit_kernel(A: cute.Tensor, x: cutlass.Int32, y: cutlass.Int32):

    # Conditional on dynamic value
    if x < cute.size(A):
        A[x] = y
    else:
        ...


@cute.jit
def jit_func(A: cute.Tensor):
    x = 0
    y = 3
    jit_kernel(
        A, x, y, config=cutlass.LaunchConfig(grid=[1, 1, 1],
                                             block=[1, 1, 1])
    )

A_tensor = torch.tensor([0, 0, 0], dtype=torch.int32).cuda()
jit_func(A_tensor)
```

This will be converted to a dynamic if automatically which will be executed at runtime

# Metaprogramming with an imperative style
## Python to Python code generation: dynamic loop

- Pythonic way to write "meta-kernel" that automatically fits for dynamic layout
  - Loop on dynamic expression

```python
@cute.kernel
def jit_kernel(A: cute.Tensor, x: cutlass.Int32, y: cutlass.Int32):

    # Loop on dynamic value
    for i in range(cute.size(A)):
        ...

    # Loop on dynamic value with loop unrolling
    for i in range_dynamic(cute.size(A), unroll=1):
        ...


@cute.jit
def jit_func(A: cute.Tensor):
    x = 0
    y = 3
    jit_kernel(
        A, x, y, config=cutlass.LaunchConfig(grid=[1, 1, 1],
                                             block=[1, 1, 1])
    )

A_tensor = torch.tensor([0, 0, 0], dtype=torch.int32).cuda()
jit_func(A_tensor)
```

Both will be converted to dynamic loop automatically which will be executed at runtime

# Metaprogramming with an imperative style

Compile-time constants as Constexpr

- Pythonic way to write "meta-kernel" that automatically fits for dynamic shape
  - Use Constexpr for constants known at compile-time

```python
@cute.kernel
def jit_kernel(A: cute.Tensor, x: cutlass.Int32, y: cutlass.Int32 , z: cutlass.Constexpr[int]):

    # Loop on compile-time constant
    for i in range(z):
        A[x + i] = y


@cute.jit
def jit_func(A: cute.Tensor):
    x = 0
    y = 3
    z = 3
    jit_kernel(
        A, x, y, z, config=cutlass.LaunchConfig(grid=[1, 1, 1],
                                                 block=[1, 1, 1]))

A_tensor = torch.tensor([0, 0, 0], dtype=torch.int32).cuda()
jit_func(A_tensor)
```

No conversion for compile-time constants

# Metaprogramming with an imperative style

TiledCopy and SMEM layout as kernel parameters

```cpp
using CollectiveOp = cutlass::gemm::collective::CollectiveMma<
    DispatchPolicy,
    TileShape_MNK,
    ElementA,
    cute::tuple<cutlass::gemm::TagToStrideA_t<GmemLayoutATag>,
                cutlass::gemm::TagToStrideA_t<GmemLayoutSFATag>>,
    ElementB,
    cute::tuple<cutlass::gemm::TagToStrideB_t<GmemLayoutBTag>,
                cutlass::gemm::TagToStrideB_t<GmemLayoutSFBTag>>,
    TiledMma,
    GmemTiledCopyA,
    SmemLayoutAtomA,
    void,
    cute::identity,
    GmemTiledCopyB,
    SmemLayoutAtomB,
    void,
    cute::identity
  >;
```

```python
@cute.kernel
def kernel(
    self,
    mA: cute.Tensor,
    mB: cute.Tensor,
    mC: cute.Tensor,
    sA_layout: cute.Layout,
    sB_layout: cute.Layout,
    tiled_copy_A: cute.TiledCopy,
    tiled_copy_B: cute.TiledCopy,
    tiled_mma: cute.TiledMma,
    epilogue_op: cutlass.Constexpr = lambda x: x,
):
    ...
```

# Data types modeling to generate fast code

```
DslType
   ├── Numeric
   │      ├── Integer
   │      │      ├── Uint8
   │      │      ├── Int32
   │      │      └── ...
   │      └── Float
   │             ├── Float8E5M2
   │             ├── Float6E3M2FN
   │             ├── Float4E2M1FN
   │             └── ...
   ├── Tensor
   └── Pointer
```

- Primitive types
  - Full support of CUTLASS data types
  - One type for all contexts
    - Host/device JIT function
    - Non-JIT context
  - Easy type access and conversions with CuTe DSL type
    - torch.dtype
    - numpy.dtype

- Compound types
  - Tensor
    - Modeling CuTe tensor concept
  - Pointer
    - Modeling the raw pointer to memory location

# Better code expressiveness and readability

Operator overloading for arithmetic, comparison and bitwise

```python
@cute.jit
def jit_func(x: cutlass.Int32, y: cutlass.Int32):
    # basic arithmetic
    x + y
    x - y
    x * y
    x // y
    x ** y
    x % y

    x >= y
    x < y

    x & y
    x | y
    x ^ y
    x << y
    x >> y
    ~x

    ...
```

- Pythonic way to deal with DslTyped values

  ❌ Avoid tedious spelling like **arith.muli(a, arith.constant(a.type, 4))**

  ✅ Instead, simply write **a * 4** and let the DSL take care of it for you

# Better code expressiveness and readability

## Operator overloading with vectorization

```python
@cute.kernel
def sgemm_kernel(
    mA: cute.Tensor,
    mB: cute.Tensor,
    mC: cute.Tensor,
    ...
):
    ...
    # Activation function fusion
    tCrC = cute.make_tensor_like(tAcc)
    for i in range_dynamic(cute.size(tAcc)):
        a = tAcc[i]
        tCrC[i] = if a > 0 a else a.dtype(0)
    ...
```

```python
@cute.kernel
def sgemm_kernel(
    ...
    epilogue_op: cutlass.Constexpr = lambda x: x,
):
    ...
    # Activation function fusion
    tCrC.store(epilogue_op(tCrC.load()))
    ...

# Fused GEMM and ReLU
gemm(input.detach(),
     weight.detach(),
     output,
     epilogue_op=lambda x: cute.where(x > 0, x, cute.full_like(x, 0)),
)
```

- **TensorSSA**: thread local data modeling for CuTe Tensor in value semantics and immutable
  - Vector based with nested CuTe shape support
  - Load tensor elements as vector / store vector data into tensor
  - Operator overloading for vectorized operations

# Customized C struct like data types
@cute.struct

- cute.struct decorator
  - Transform a Python class into a memory-mapped structure with precise control over memory layout, alignment and offsets
  - Support scalar, MemRange or nested struct as data members
  - Allow customized data alignment which is essential for better performance

```python
@cute.struct
class complex:
    real: cutlass.Float32
    imag: cutlass.Float32

@cute.struct
class MyStorage:
    x: cutlass.Float32
    y: cutlass.Int32
    nested: cute.struct.align(complex, 16)
    mem: cute.struct.align(cute.struct.MemRange(
        cutlass.Float32, cute.cosize(layout_a)), 1024)
```

| offset | MyStorage | alignment | |
|--------|-----------|-----------|---|
| 0 | x | 4 | Natural alignment per dtype |
| 4 | y | 4 | |
| 16 | nested | 16 | User specified alignment |
| 1024 | mem | 1024 | |

# Kernel writing in an OOP manner

## @cute.struct

```python
@cute.kernel
def kernel(...):
    ...
    smem = cutlass.utils.SmemAllocator()
    ab_full_mbar_ptr = smem.allocate_array(cutlass.Int64,
                                           self.ab_stage)

    ab_empty_mbar_ptr = smem.allocate_array(cutlass.Int64,
                                            self.ab_stage)

    buffer_align_bytes = 1024
    sa = smem.allocate_tensor(
        self.a_dtype,
        a_smem_layout_staged.outer,
        buffer_align_bytes,
        swizzle=a_smem_layout_staged.inner)
    sb = ...
    ...
def _compute_smem(cta_tile_shape_mnk, a_dtype, b_dtype, ab_stage, c_dtype):
        a_shape = cute.slice_(cta_tile_shape_mnk, (None, 0, None))
        b_shape = cute.slice_(cta_tile_shape_mnk, (0, None, None))
        ab_bytes_per_stage = (
            cute.size(a_shape) * a_dtype.width // 8
            + cute.size(b_shape) * b_dtype.width // 8)

        mbar_helpers_bytes = 1024
        num_smem_bytes = ab_bytes_per_stage * ab_stage + mbar_helpers_bytes
        return num_smem_bytes
self.kernel(...,
            config=cutlass.LaunchConfig(
                ...
                smem=self._compute_smem(...),
                async_deps=[stream]))
```

```python
@cute.kernel
def kernel(...):
    ...
    smem = cutlass.utils.SmemAllocator()
    storage = smem.allocate_struct(self.shared_storage)
    # access ab_full_mbar_ptr through struct data members
    ab_full_mbar_ptr = storage.ab_full_mbar_ptr.data_ptr()
    ...

@cute.struct
class SharedStorage:
    ab_full_mbar_ptr: cute.struct.MemRange(cutlass.Int64,
                                           self.ab_stage)
    ab_empty_mbar_ptr: cute.struct.MemRange(cutlass.Int64,
                                            self.ab_stage)
    sa: ...
    sb: ...

self.kernel(...,
            config=cutlass.LaunchConfig(
                ...
                smem=SharedStorage.size_in_bytes(),
                async_deps=[stream]))
```

NVIDIA.

# Reduced kernel launching latency with caching

## Significant overhead without caching

```python
class MyGemmKernel:
    @cute.jit
    def __call__(
        self,
        a: cute.Tensor,
        b: cute.Tensor,
        c: cute.Tensor,
        stream: cutlass.Stream,
        epilogue_op: cutlass.Constexpr = lambda x: x,
    ):
        ...

def forward(self, input, bias=None):
    batch_size, seq_len, hidden_size = input.shape

    try:
        # Reshape input to prepare for GEMM
        input = input.reshape(batch_size * seq_len, hidden_size)

        output = torch.empty(
            input.size(0), self.out_features, device=input.device,
                                              dtype=input.dtype)

        ...

        self.gemm(
            input.detach().contiguous(),
            weight.detach().contiguous(),
            output.contiguous(),
            stream)
        ...
```

- In each forward pass, the method will be called with kernel JIT compilation which would cause a significant runtime overhead

NVIDIA.

# Reduced kernel launching latency with caching

Zero Compile: JIT Executor with CUBIN cached

```python
def forward(self, input, bias=None):
    batch_size, seq_len, hidden_size = input.shape

    try:
        # Reshape input to prepare for GEMM
        input = input.reshape(batch_size * seq_len, hidden_size)

        output = torch.empty(
            input.size(0), self.out_features, device=input.device, dtype=input.dtype)
        ...


        input_tensor = from_dlpack(input.detach().contiguous()).mark_layout_dynamic()
        weight_tensor = from_dlpack(weight.detach().contiguous()).mark_layout_dynamic()
        output_tensor = from_dlpack(output.contiguous()).mark_layout_dynamic()
```

```python
        key = input.shape
        if key not in self.cached_kernels:
            self.cached_kernels[key] = cute.compile(
                self.gemm,
                input_tensor,
                weight_tensor,
                output_tensor,
                cutlass_stream,
            )
        self.cached_kernels[key](
            input_tensor, weight_tensor, output_tensor, cutlass_stream
        )
        ...
```

- Customized keys for kernel caching

- JIT Executor with CUBIN cached

- JIT Executor supports serialization to file and deserialization from file

NVIDIA.

# Reduced kernel launching latency with caching

## Zero Compile: JIT Executor with CUBIN cached

```
@cute.kernel
def my_kernel(A: cute.Tensor, ...):
    ...
    atom = cute.make_copy_atom(...)
    ...
    cute.make_tiled_copy_tv(...)
    ...
```

Skip generation & compilation to
access cached JIT Executor(s) directly

JIT
Executor    CUBIN

JIT
Executor    CUBIN

# Generation, Compilation, and Launch overhead
### Blackwell B100 FP16 GEMM: M=N=K=8K

- **Without cache**
  - Goes all the way down to kernel launching
  - Includes IR generation, compilation and kernel launch

- **Stable cache**
  - By default ON
  - Always generates codes for functional correctness
  - Skips compilation if generated codes are identical

- **Zero Compile**
  - Launches kernel directly through JIT Executor managed by users
  - No IR generation and compilation
  - Minimized kernel launch overhead at ~4.6 us

# Agenda

# What is CuTe and why should you care about it

## Program GPUs at peak throughout architectural generations

- What's needed for performance
  - Complex tiling and partitioning patterns
  - Arch-specific instructions with their own set of requirements

- Layout book-keeping is hard, error-prone, and takes away time

- CuTe at your rescue without taking away control!
  - A single Layout concept to express all layouts of interest and more based on a hierarchical representation
  - A formal algebra
  - A programming model maintaining logical consistency
  - A consistent set of idioms applicable throughout GPU generations
  - A safer low-level programming
    - Detect illegal patterns by inspecting layouts at compile-time



Figure 3b: Maxwell thread assignments for LDS.U matching

Maxwell A access addr[tid] = addr[tid^1]

Maxwell B access addr[tid] = addr[tid^2]



LDG.128 register quad ($R_{n..n+3}$)

Shared memory after STS.128 w/swizzled thread id = (tid[1:0] << 3) | (tid & 4) | (tid[4:3] ^ tid[1:0])

# CuTe exposure

If you already use CUTLASS-C++, you will feel at home

- CuTe Layouts and Tensors in all their flavors
  - The algebra is available in its entirety
  - Mixed static/dynamic Layouts are fully supported

- Robust tensor programming model that users of CUTLASS-C++ are already familiar with

**MMA/Copy Atoms**
Encapsulate the PTX with metadata encoded using Layouts

**Tiled MMA/Copy**
Robust generic partitioning interface

**cute.{gemm|copy} algorithms**
Automatically dispatch to the requested PTX instruction

**CuTe Layouts and Algebra**

# Programming with CuTe in Python

## Examples

1. Tiling MMAs
2. Ampere warp-level MMA
3. Blackwell 2CTA MMA



Data Layouts
`logical coord` → `data offset`

$\texttt{shape:stride} = (8, 8){:}(8, 1)$

Thread-Value (TV) Layouts
`(thr,val)` → `logical coord`

# Programming with CuTe in Python

Example 1: Building complex tiling of MMA Atoms

N

0

K    0    T0
          V0

0

M    0    T0        T0
          V0        V0

MMA Atom MxNxK=1x1x1

1x1x1 Universal MMA based on FMA

# Programming with CuTe in Python

## Example 1: Building complex tiling of MMA Atoms



Tile this atom in an 8x8 fashion

- threads 0, 8, 16,... own in their registers the same entry of A
- threads [0,7] own in their registers the same entry of B

Description tracked by Tiled MMA with TV-Layouts

# Programming with CuTe in Python

Example 1: Building complex tiling of MMA Atoms



Tiled MMA 16x16x1

- Tile further across values
  - Each thread computes a 2x2 accumulator fragment
- Permute the tiling to enjoy vectorized loads from SMEM to RMEM for M-major A and N-major B

# Programming with CuTe in Python

## Example 1: Building complex tiling of MMA Atoms



Tiled MMA 16x16x1

A programming pattern applicable to *any* MMA!

```python
# Layout of Atoms: tiling across threads
atom_layout = cute.make_layout((8, 8, 1))
# Permutation Tiler:
# – tiling across values
# – permutation of the overall tiling
perm_tiler = (
    cute.make_layout((8, 2), stride=(2, 1)),
    cute.make_layout((8, 2), stride=(2, 1)),
    None,
)
tiled_mma = cute.make_tiled_mma(
    cute.nvgpu.MmaUniversalOp(),
    atom_layout,
    perm_tiler,
)
```

# Programming with CuTe in Python

Example 1: Building complex tiling of MMA Atoms



Tiled MMA 16x16x1

Partitioning without ever worrying about the physical stride

```
thr_mma = tiled_mma.get_slice(tidx)
tCsA = thr_mma.partition_A(sA)
```

M-major (column-major)
MxK=32x8 SMEM A

$$\mathtt{base\_ptr} \circ (32, 8){:}(1, 32)$$

$$\mathtt{offset\_ptr} \circ (1, (2, 2), 8){:}(0, (1, 16), 32)$$

- The MMA mode for a single atom

- # repetitions along M to cover the full 32x8 tile

- # repetitions along K to cover the full 32x8 tile

Partitioning at a high level:

$$\mathtt{t} = \mathrm{Data} \circ \mathrm{TV} : (\mathtt{thr\_id}, \mathtt{val\_id}) \to \mathtt{data\ offset}$$

$$\mathtt{thr\_part} = \mathtt{t}(\mathtt{my\_thr}, \_)$$

42

# Programming with CuTe in Python

Example 2: GMEM → SMEM → HMMA

- Consider the A operand of a half-precision F16 warp-level MMA (HMMA)

| K-major A in GMEM | → | 16x64 SMEM | → | HMMA |

# Programming with CuTe in Python

## Example 2: Resolving SMEM bank conflicts

128B = 32 x 32b = 64 x 16b

|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| 1 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
| 2 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 |
| 3 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
| 4 | 32 | 33 | 34 | 35 | 36 | 37 | 38 | 39 |
| 5 | 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 |
| 6 | 48 | 49 | 50 | 51 | 52 | 53 | 54 | 55 |
| 7 | 56 | 57 | 58 | 59 | 60 | 61 | 62 | 63 |

The first 8 rows of the 16x64 SMEM

16x64 SMEM

$(8, 8):(8, 1)$

- SMEM is organized into 32 x 32b banks
- Starting with a simple row-major SMEM layout
  - Represented in unit of 128b elements
  - Each color is a set of 4 banks
- Row-major A in GMEM → use LDG.128 + STS.128
- 128b writes across a row are free of bank conflicts… but
- 128b reads across a column all hit the same 4 banks
- Solution: swizzle

# Programming with CuTe in Python

Example 2: Resolving SMEM bank conflicts



| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| 1 | 9 | 8 | 11 | 10 | 13 | 12 | 15 | 14 |
| 2 | 18 | 19 | 16 | 17 | 22 | 23 | 20 | 21 |
| 3 | 27 | 26 | 25 | 24 | 31 | 30 | 29 | 28 |
| 4 | 36 | 37 | 38 | 39 | 32 | 33 | 34 | 35 |
| 5 | 45 | 44 | 47 | 46 | 41 | 40 | 43 | 42 |
| 6 | 54 | 55 | 52 | 53 | 50 | 51 | 48 | 49 |
| 7 | 63 | 62 | 61 | 60 | 59 | 58 | 57 | 56 |

- Reads/writes across rows/columns are now free of bank conflicts!

- This layout is exactly a swizzle layout

$$S<3,0,3> \circ 0 \ \circ (8,8){:}(8,1)$$

- ```
  l = cute.make_composed_layout(
      cute.swizzle(3,0,3),
      0,
      cute.make_layout((8,8), stride=(8,1)),
  )
  assert l( (3,6) ) == 29
  ```

  No indexing math!

- Such layout can be partitioned just like **any** other layout

# Programming with CuTe in Python

128b = 8 x 16b



• • •    64

- A's TV-layout for a 16x8x8 HMMA

- How do I prepare my register fragments *efficiently*?

- Each `ldmatrix` instruction loads 128b from 8 rows and places the values in the registers of threads according to HMMA

- This is exactly one column in the previous figure, thus free of bank conflicts

- This pattern can be repeated along K for larger tiles and remains free of bank conflicts

PTX doc for ldmatrix
"The eight addresses required for each matrix are provided by eight threads [...] Each address corresponds to the start of a matrix row."

Threads 0-7      ???

addr0–addr7

# Programming with CuTe in Python

## Example 2: SMEM → RMEM

Construct the Tiled MMA and partition the (swizzled!) SMEM tensor according to the MMA

```python
# Construct a 16x8x8 MMA with F16/F32 inputs/output
# Partition SMEM tensors according to the MMA
mma_op = warp.MmaF16BF16Op(
    cutlass.Float16, cutlass.Float32, (16, 8, 8))
tiled_mma = cute.make_tiled_mma(mma_op)
thr_mma = tiled_mma.get_slice(tidx)
tCsA = thr_mma.partition_A(sA)
tCrA = thr_mma.make_fragment_A(tCsA)
```

Construct a Tiled Copy of `ldmatrix`
***based off the MMA***

```python
# Construct a tiled Copy based off the MMA
# This takes care of any repetition (x2 along M here)
copy_op = warp.LdMatrix8x8x16bOp(transpose=False)
copy_atom = cute.make_copy_atom(copy_op, cutlass.Float16)
tiled_copy = cute.make_tiled_copy_A(copy_atom, tiled_mma)
```

Repartition the source SMEM tensor now according to the Copy

```python
# Repartition SMEM source tensors according to the Copy
# Retile RMEM destination tensors according to the Copy
thr_copy = tiled_copy.get_slice(tidx)
tCsA_copy = thr_copy.partition_S(sA)
tCrA_copy = thr_copy.retile(tCrA)
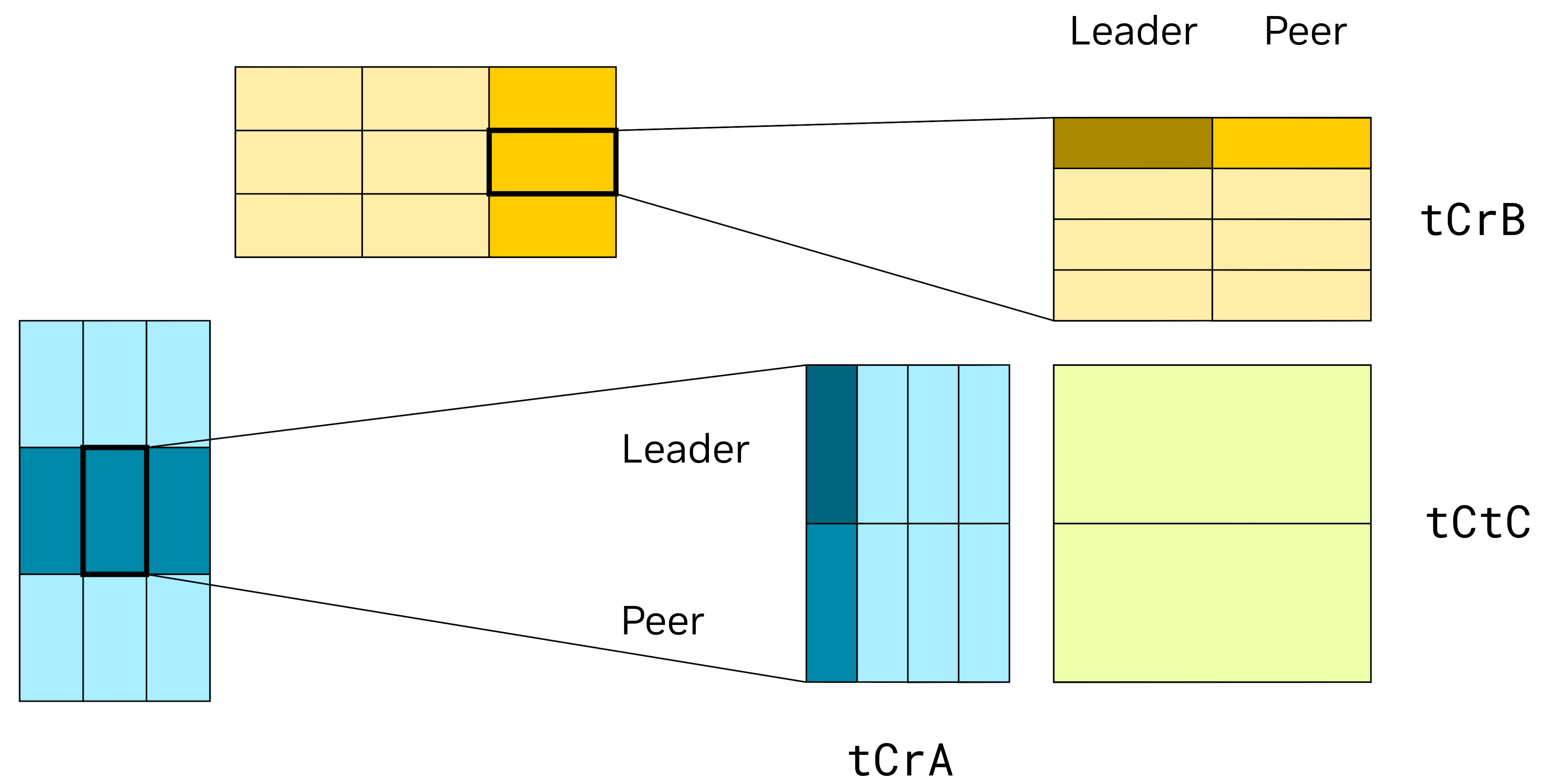```

# Programming with CuTe in Python

## Example 2: SMEM → RMEM



Copy's src view conforming with `ldmatrix`

Copy's dst view = MMA view

```python
# Construct a 16x8x8 MMA with F16/F32 inputs/output
# Partition SMEM tensors according to the MMA
mma_op = warp.MmaF16BF16Op(
    cutlass.Float16, cutlass.Float32, (16, 8, 8))
tiled_mma = cute.make_tiled_mma(mma_op)
thr_mma = tiled_mma.get_slice(tidx)
tCsA = thr_mma.partition_A(sA)
tCrA = thr_mma.make_fragment_A(tCsA)


# Construct a tiled Copy based off the MMA
# This takes care of any repetition (x2 along M here)
copy_op = warp.LdMatrix8x8x16bOp(transpose=False)
copy_atom = cute.make_copy_atom(copy_op, cutlass.Float16)
tiled_copy = cute.make_tiled_copy_A(copy_atom, tiled_mma)


# Repartition SMEM source tensors according to the Copy
# Retile RMEM destination tensors according to the Copy
thr_copy = tiled_copy.get_slice(tidx)
tCsA_copy = thr_copy.partition_S(sA)
tCrA_copy = thr_copy.retile(tCrA)
```

# Programming with CuTe in Python

Example 3: 2CTA Blackwell MMA

```python
tCrA = tiled_mma.make_fragment_A(sA)    # Create a tensor of SMEM descriptors for the A operand
tCgC = thr_mma.partition_C(gC)
tCtC = tiled_mma.make_fragment_C(tCgC) # Create a TMEM tensor for the accumulator


# Issue a GEMM
cute.gemm(
    tiled_mma,
    tCtC,
    tCrA[(None, None, k_block_idx)],
    tCrB[(None, None, k_block_idx)],
    tCtC,
)
```

The 2CTA MMA consumes the SMEM tensors of both CTAs and stores half of the result in each CTA's TMEM allocation.



Leader   Peer

tCrB

tCrA

Leader

Peer

tCtC

# What if something is missing?

## Seamless integration of raw Op builders

```python
from cutlass._mlir.dialects import llvm, nvvm


def thread_idx(*, loc=None, ip=None):
    return (
        nvvm.read_ptx_sreg_tid_x(T.i32(), loc=loc, ip=ip),
        nvvm.read_ptx_sreg_tid_y(T.i32(), loc=loc, ip=ip),
        nvvm.read_ptx_sreg_tid_z(T.i32(), loc=loc, ip=ip),
    )


def fence_tma_desc_acquire(tma_desc_ptr_i64, *, loc=None, ip=None):
    llvm.inline_asm(
        None,
        [tma_desc_ptr_i64],
        "fence.proxy.tensormap::generic.acquire.gpu [$0], 128;",
        "l",
        has_side_effects=True,
        is_align_stack=False,
        asm_dialect=llvm.AsmDialect.AD_ATT,
        loc=loc,
        ip=ip,
    )
```

- LLVM + NVVM Op builders exposed
- Direct access to NVVM Ops
- *Inline PTX* is straightforward
- Subject to upstream breaking changes

# What else?

## Refreshed tutorials, examples, and documentation

- Initial support for schedulers and persistency

- Initial support for pipeline abstractions

- New centralized documentation

- Educational notebooks from getting started to SOL

# Agenda

# Blackwell Performance: Python vs. C++

## FP16 I/O GEMM with M=N=8192



- 8192Mx8192N Python Math SOL%
- 8192Mx8192N C++ Math SOL%

Testing spec: B100 180GB HBM3e, 148SM GPC-800MHz/DRAM 4GHz 850W

128x256x64, Warp-specialized, 2x1 Cluster shape

53

# Blackwell Performance: Python vs. C++

## FP16 I/O GEMM with M=N=2048



Testing spec: B100 180GB HBM3e, 148SM GPC-800MHz/DRAM 4GHz 850W

256x128x64, Warp-specialized, 2x1 Cluster shape

# Hopper Performance: Python vs. C++

## FP16 I/O GEMM with M=N=8192



- Perf gap

Python example does not include persistent optimization thus expose more overhead among CTA waves for small GEMM-K cases

*Testing spec: H100 80GB HBM3, 132SM GPC-1500MHz/DRAM 2619MHz 700W*

*128x256x64 cooperative size, Swizzle size = 8*

# Hopper Performance: Python vs. C++

## FP16 I/O GEMM with M=N=2048



- Perf gap

For small GEMM-K case, it's single wave hence C++ exposed more overhead due to the tile scheduling calculation cost with its warp-specialized persistent scheduler which is not used in Python example.

*Testing spec: H100 80GB HBM3, 132SM GPC-1500MHz/DRAM 2619MHz 700W*

*128x256x64 cooperative size, Swizzle size = 8*

# Blackwell Performance: Python vs. C++

## FP16 I/O Group GEMM



- Perf gap

More optimized JIT codes b/c of

1) static cluster size
2) No insts related to residual if beta=0
3) Simplified prologue and epilogue

*Testing spec: B100 180GB HBM3e, 148SM GPC-800MHz/DRAM 4GHz 850W*

*128x256/128x32/64x32/128x128/128x32 CTA Tilesize, Warp-specialized,*

*Cluster shape (C++) 2x1/2x1/2x1/1x1/1x1 vs. (Python)2x1/1x1/1x1/1x1/1x1*

# CUTLASS Python

CUTLASS comes to Python and is first class citizen with full support!

## Initial beta release

- Release date: Q2'25

- Focus on GEMMs, including grouped GEMM

- Blackwell B200: support for major features
  - TMA, 2CTA MMA, TMEM

- Ada/Ampere/Hopper: experimental

- Schedulers and pipelines

- Jupyter notebook examples
  - Tutorial series on how to use the DSL
  - *Back to the basics* SGEMM tutorial series
  - Ampere TC GEMM and FA2
  - Hopper WGMMA GEMM
  - Blackwell GEMM tutorial series
  - Blackwell grouped GEMM
  - Blackwell FA2

## What will come next

- GeForce RTX 50 Series

- Feature complete support for Blackwell including
  - Cluster Launch Control
  - Programming Dependent Launch

- Narrow-precision data types support and block-scaled MMAs

- EVT

- More examples of advanced fusion

- Convolutions

- Fully fledged Ahead-of-time compilation support

- Higher level abstractions & primitives

- More comprehensive documentation

- Graduate Ada, Ampere, Hopper from experimental

# Acknowledgements

A great collaboration among various teams made it possible!

THANK YOU !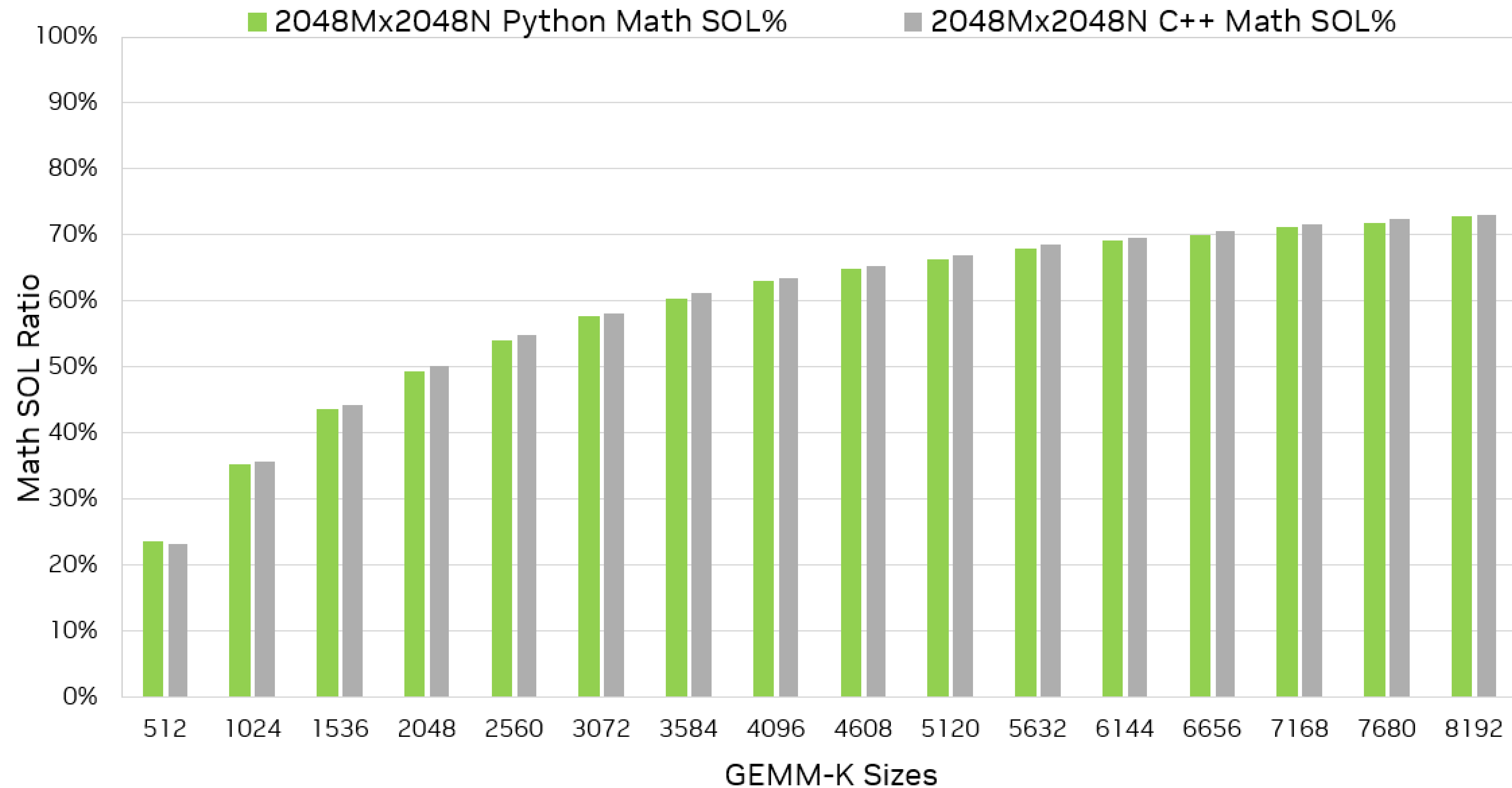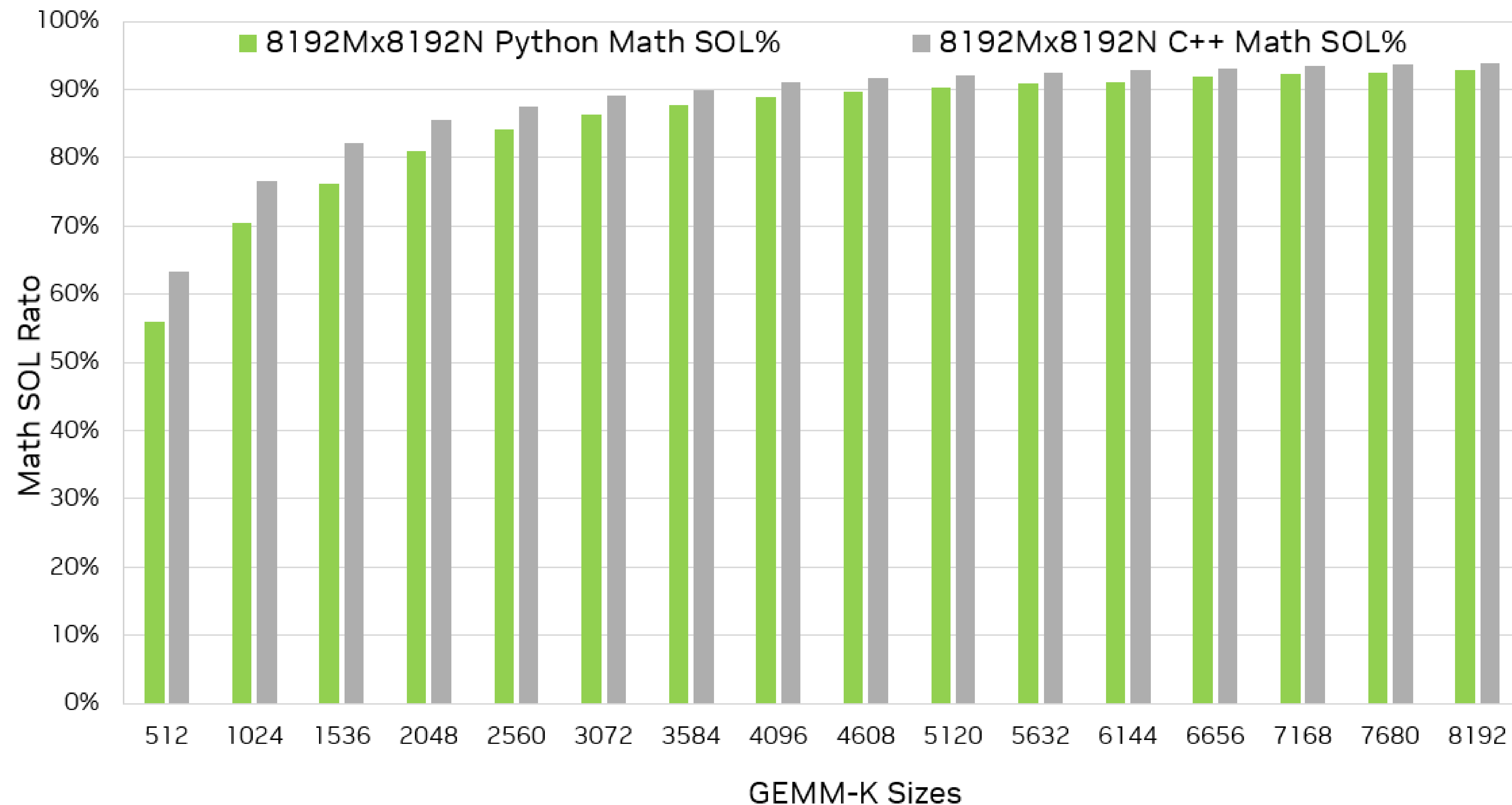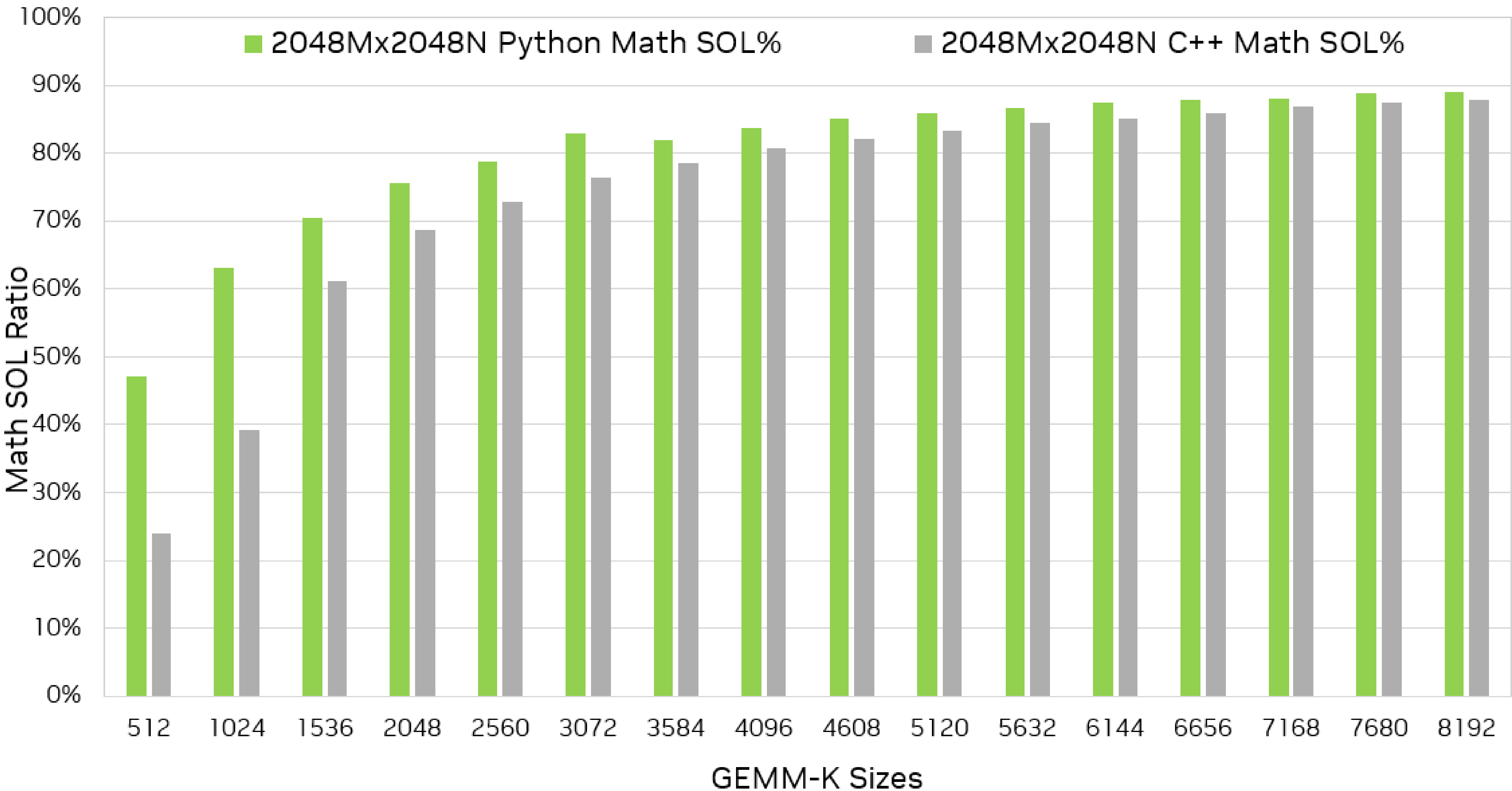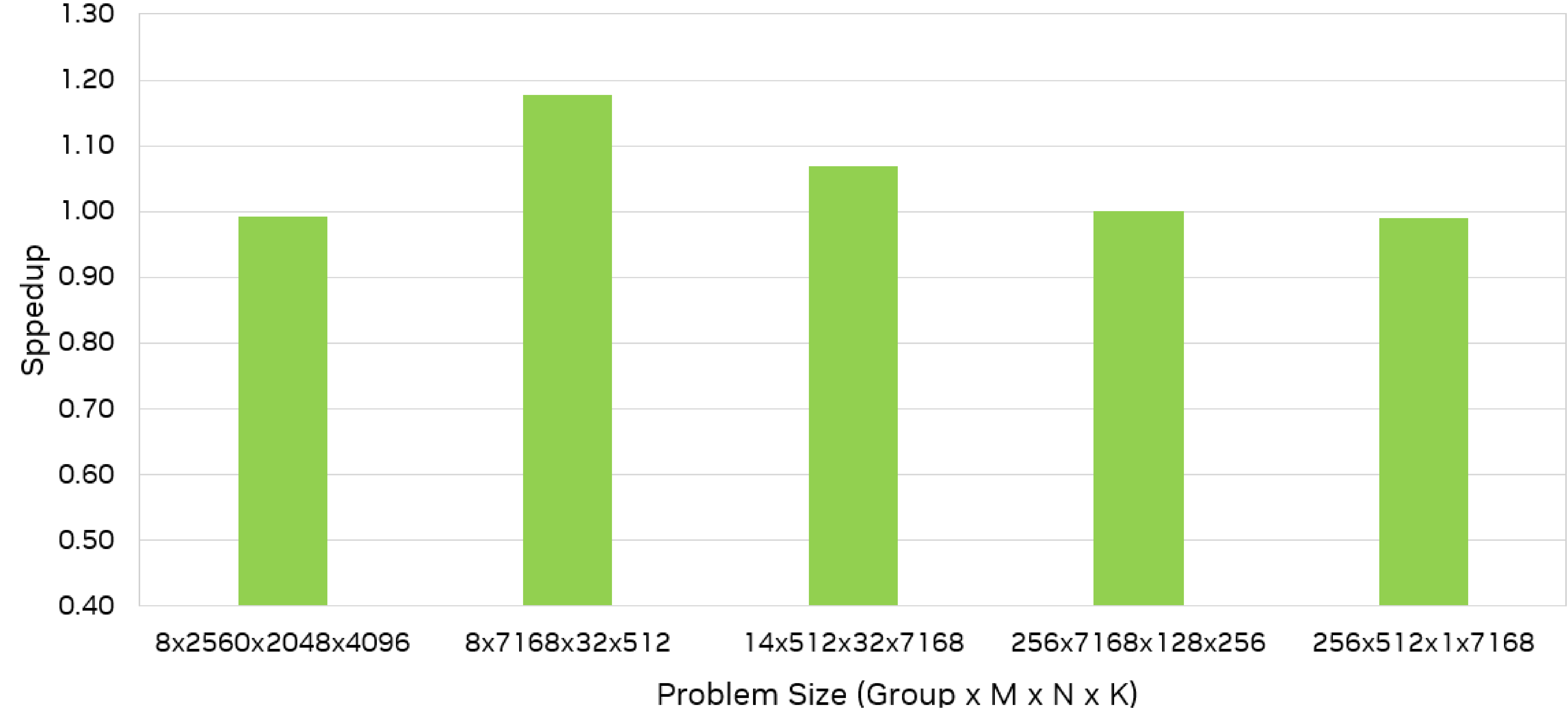