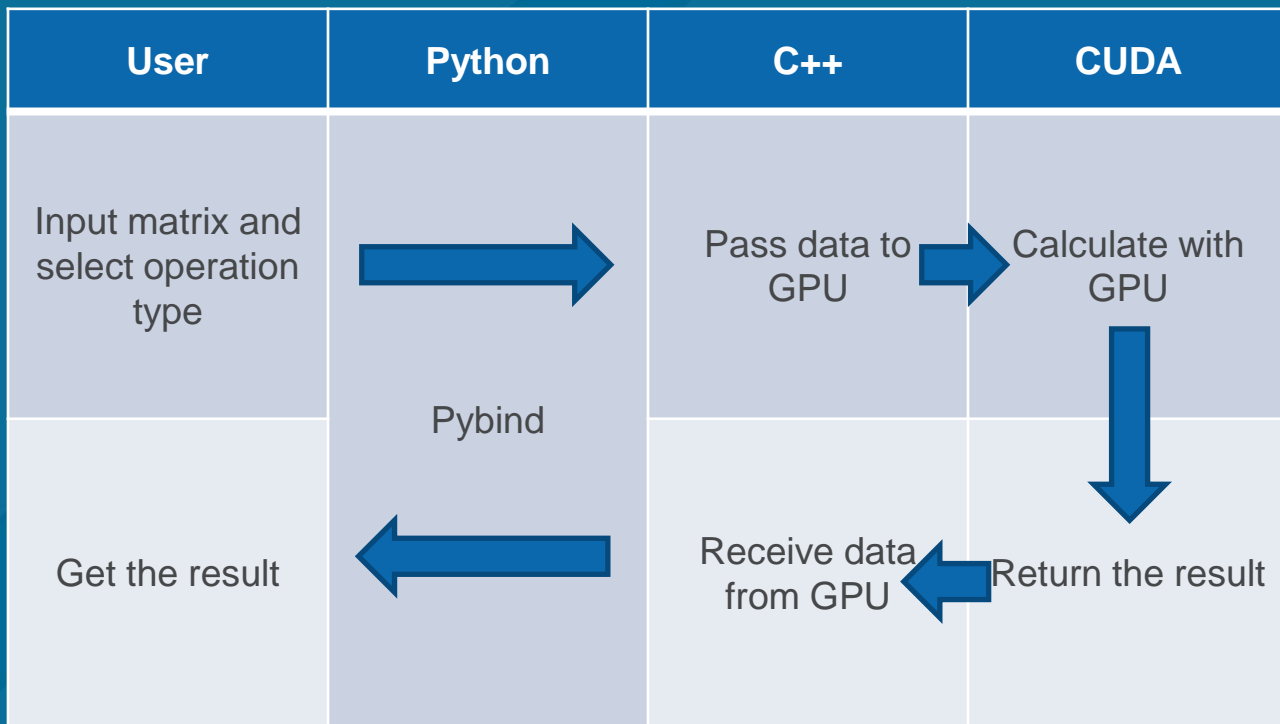# ECE 277 Project Presentation
## GPU Accelerated Common Matrix Operations with CUDA and Pybind

# GPU Accelerated Common Matrix Operations with CUDA and Pybind

- GPU is good at parallel computing such as matrix computation.
- This project creates a Python library and bind it with C++ using Pybind.
- Users can implement common matrix operations such as addition and multiplication on Python with GPU acceleration.
- The project only focuses on 2D matrices

# GPU Accelerated Common Matrix Operations with CUDA and Pybind



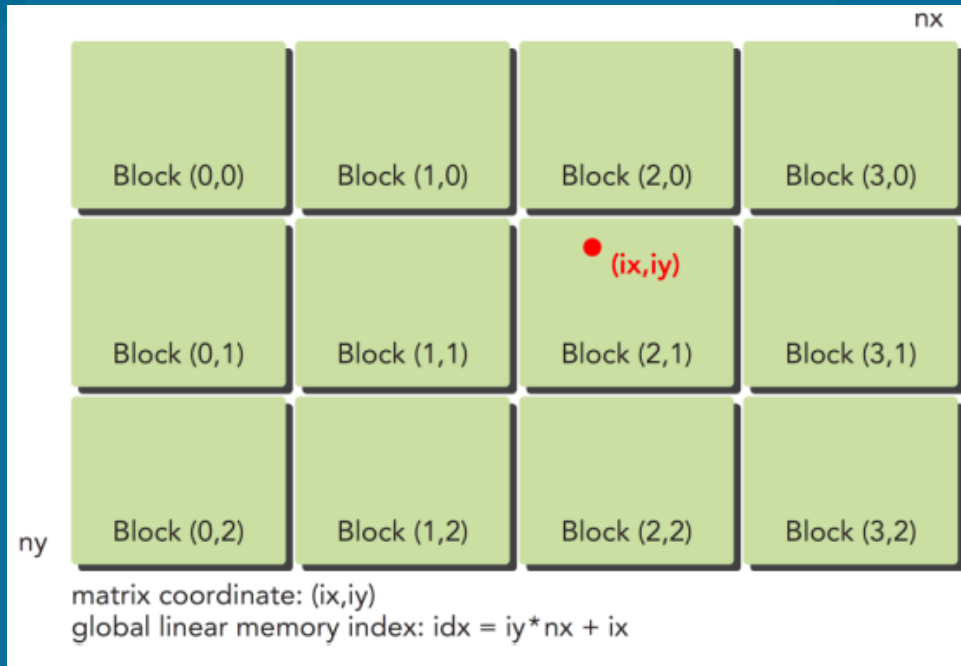| User | Python | C++ | CUDA |
|------|--------|-----|------|
| Input matrix and select operation type | Pybind | Pass data to GPU | Calculate with GPU |
| Get the result | | Receive data from GPU | Return the result |

# Operation Types

1. Addition
2. Subtraction
3. Multiplication
4. Transposition

*The project only focuses on 2D and integer matrices.

# Addition & Subtraction

## 2D grid and 2D block partition

- Map one thread to one matrix element

- One threadblock processes a sub-block of a matrix (multiple columns and multiple rows)



matrix coordinate: (ix,iy)
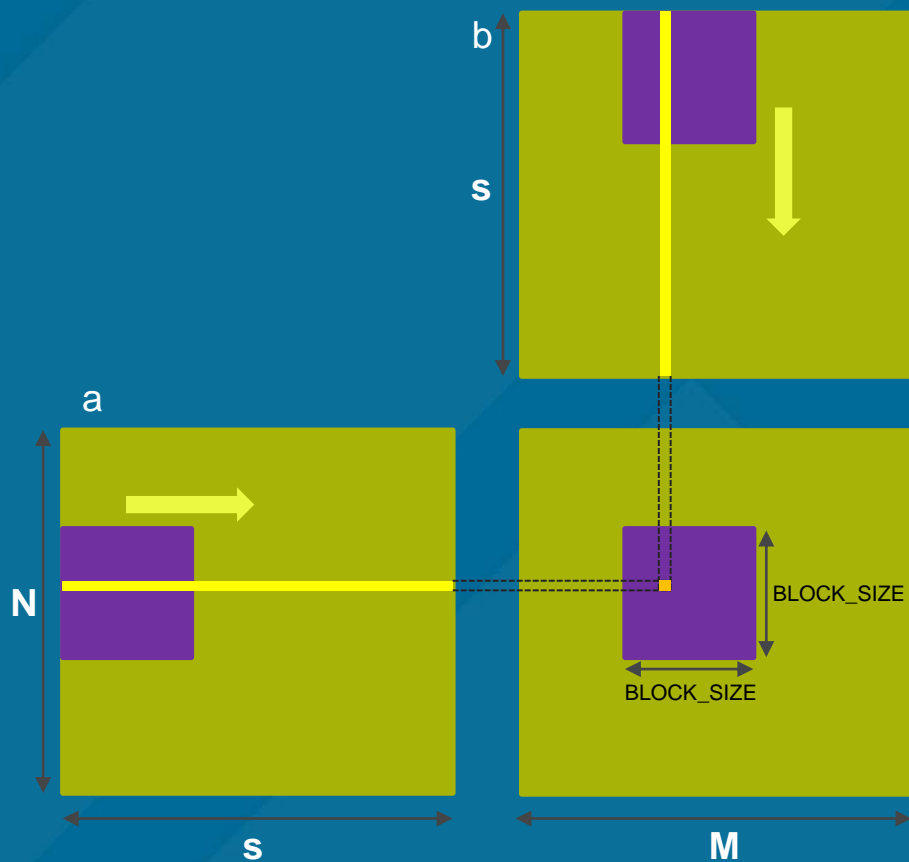global linear memory index: idx = iy*nx + ix

# Multiplication

## Use Shared Memory (SMEM) to accelerate the computation

- If we don't use SMEM, the program will access Global Memory (GMEM) multiple times during execution. Since the latency of GMEM is relatively high, the efficiency of code execution will be badly reduced.

- We introduced SMEM for optimization. We mainly use the Write Once/ Read Many (WO/RM) of SMEM to reduce the latency during data transferring.

```
40  __global__ void kernel_mmul(int* A, int* B, int* C, int M, int N, int s, dim3 blk)
41  {
42      __shared__ int smem_m[BLOCK_SIZE][BLOCK_SIZE];
43      __shared__ int smem_n[BLOCK_SIZE][BLOCK_SIZE];
44      unsigned int row = blockDim.y * blockIdx.y + threadIdx.y;
45      unsigned int col = blockDim.x * blockIdx.x + threadIdx.x;
46      unsigned int tmp = 0;
47
48      if (row + col == 0)
49          printf("cuda matrix ((%d, %d)(%d, %d)) multiplication\n", N, s, s, M);
50
51      for (int stride = 0; stride <= s / blk.y; stride++) {
52          int idm = stride * blk.y + row * s + threadIdx.x;
53          if (row < N && blk.y * stride + threadIdx.x < s) {
54              smem_m[threadIdx.y][threadIdx.x] = A[idm];
55          }
56          else {
57              smem_m[threadIdx.y][threadIdx.x] = 0;
58          }
59          int idn = stride * blk.y * M + col + threadIdx.y * M;
60          if (col < M && blk.y * stride + threadIdx.y < s) {
61              smem_n[threadIdx.y][threadIdx.x] = B[idn];
62          }
63          else {
64              smem_n[threadIdx.y][threadIdx.x] = 0;
65          }
66          __syncthreads();
67          for (int i = 0; i < blk.y; i++) {
68              tmp += smem_m[threadIdx.y][i] * smem_n[i][threadIdx.x];
69          }
70          __syncthreads();
71      }
72      if (row < N && col < M)
73      {
74          C[row * M + col] = tmp;
75      }
76  }
```
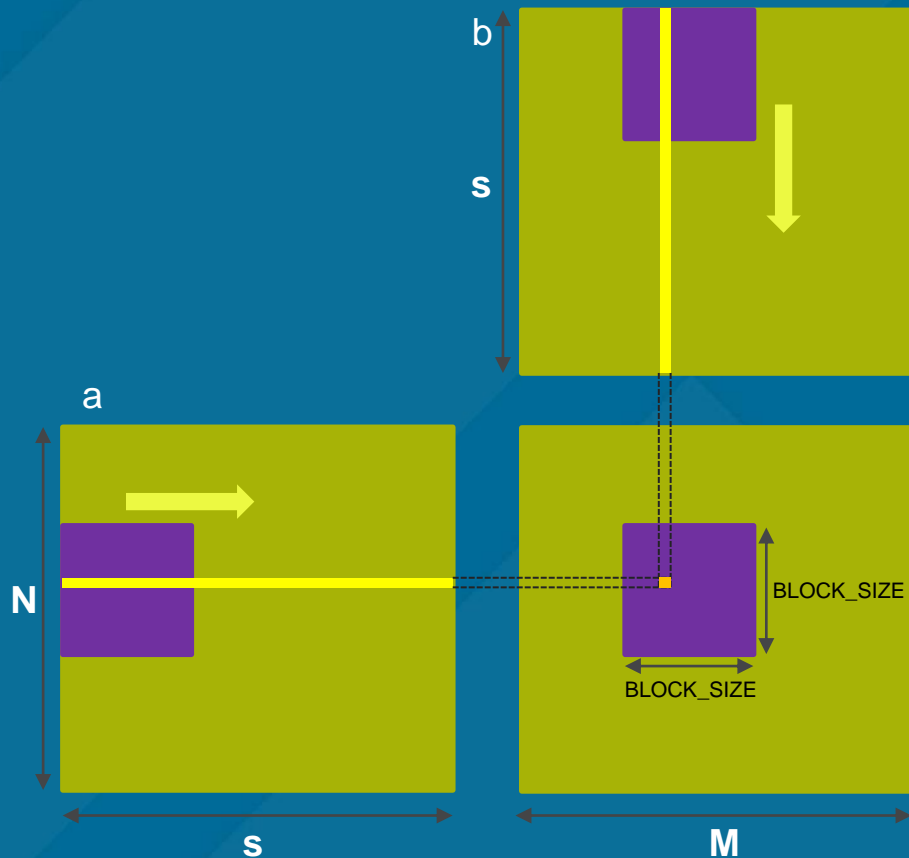
1. First, we define two tiles that exist in SMEM. The size of the square tile is the size of the current Block.
2. In the kernel function, each thread plays two roles: 1. Assign the data in the GMEM to the SMEM; 2. Calculate the value in the result matrix.
3. My codes use the method of moving tiles. While copying and calculating, tiles are moved towards certain direction. As shown in the figure, one tile will move to the positive direction of the x-axis of Matrix a, and the other tile will move to the positive direction of the y-axis of Matrix b, and the step size is the side length of the tile.
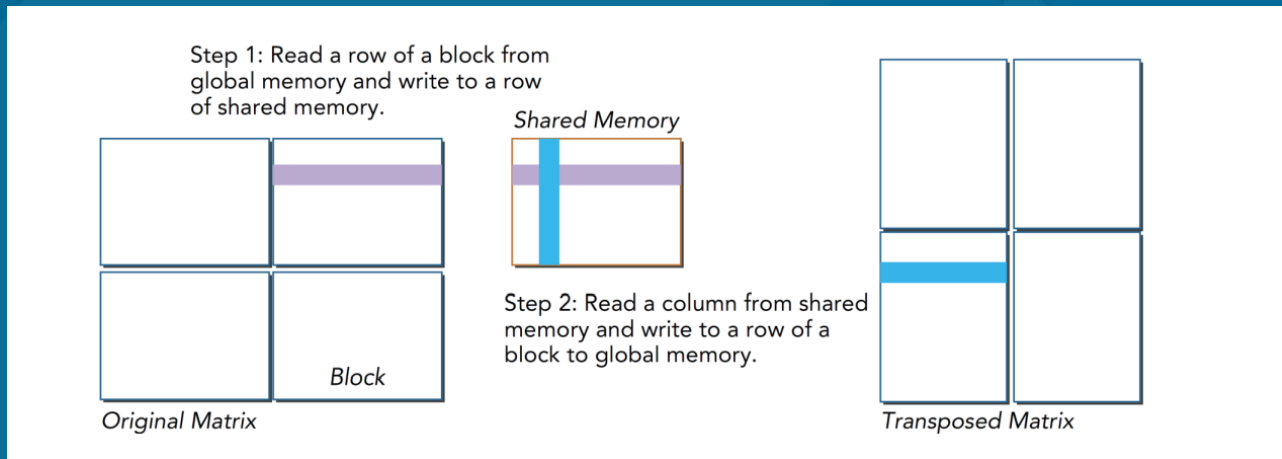
## Transpose with SMEM

1. **Row-wise GMEM read, Row-wise SMEM write**
2. **Column-wise SMEM read, Row-wise GMEM write**

**Both GMEM read and write can be coalesced access so we can make use of it to accelerate the operation.**



Step 1: Read a row of a block from global memory and write to a row of shared memory.

Shared Memory

Step 2: Read a column from shared memory and write to a row of a block to global memory.

Block

Original Matrix

Transposed Matrix

Test Result

# Addition

```
M=10240, N=10240
cuda matrix (10240, 10240) addition
Pass
Python Time: 97789.59120000001 ms
GPU Time: 513.2112000000006 ms
```

# Subtraction

```
M=10240, N=10240
cuda matrix (10240, 10240) subtraction
Pass
Python Time: 107122.3083 ms
GPU Time: 488.4299999999939 ms
```

# Multilication

```
M=256, N=256
cuda matrix ((256, 128)(128, 256)) multiplication
Pass
Python Time: 5290.227 ms
GPU Time: 185.84150000000045 ms
```

# Transposition

```
M=10240, N=5120
cuda matrix (5120, 10240) transposition
Pass
Python Time: 29242.1385 ms
GPU Time: 323.9198999999999 ms_
```