# HETEROGENEOUS COMPUTATION BASED ON GPU
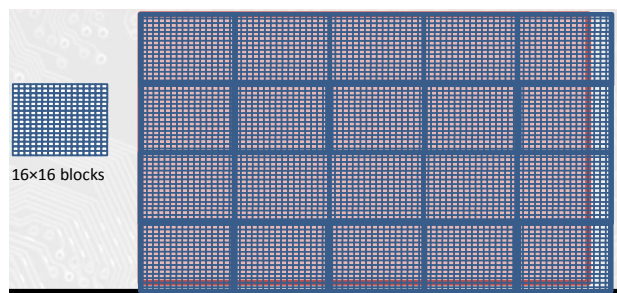
## INTRODUCTION

2021

---

# PROCESSING 2D MATRICES
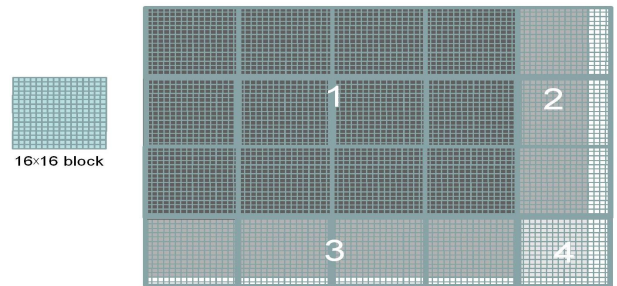
- Example: Define blocks of 16x16 blocks (256 parallel threads)

- Define N blocks, where N x Block size must fit in number of cores available in device.

- Each position will be addressed by r/o index provided by the execution context.
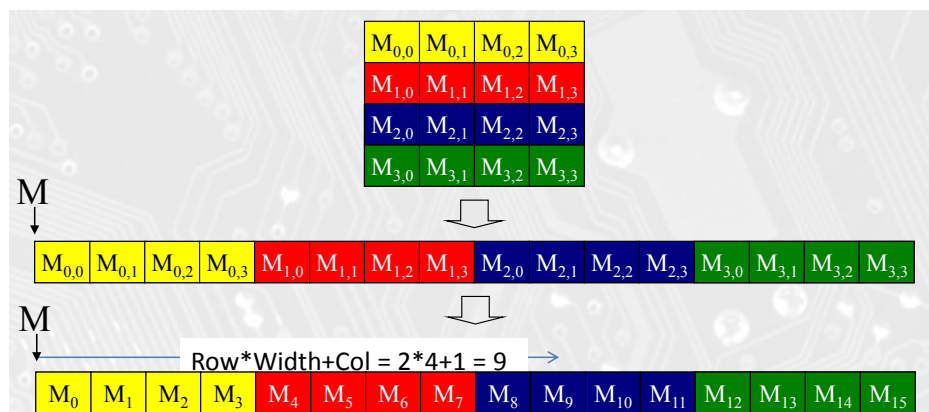
16×16 blocks

2021

# PROCESSING 2D MATRICES

- If block not fits in data:
  - Check if indexes index a defined memory position
    - If yes, execute program
    - If not, return null
  - BIG PROBLEM HERE
    - The GPU does not have branch prediction, then double the execution time.



16×16 block

---

# MEMORY INDEXING

All the data is a vector, then you have to construct the memory position using the indexes



Row*Width+Col = 2*4+1 = 9

# CALCULATE VECTOR POSITION

- GOAL: Calculate the position where the data will be stored in a row-wise representation of a 2 dimensional matrix
- How?
  - Based on the system variables blockDim, threadIdx and blockIdx, calculate the position in the matrix
  - Check if the position is a valid value (to be coded by the students)
  - If it is valid, write in the result vector, the value of his position
- Check in the notebook

---

# SQUARE VECTOR

- GOAL: Calculate the square value of input vector
- How
1. Calculate the position assigned using the blockDim, blockIdx and threadIdx,
2. Check if the memory position is valid
3. Fetch the value from the original vector a, calculate it square and save in the result vector b

# MATRIX MULTIPLICATION

- GOAL: Multiply 2 matrices of fixed size and recover the result

- How:

- First: each thread will calculate each cell in the result matrix c

1.  Calculate the assigned cell coordinates

2.  Check if the coordinates are valid

3.  If yes, use those coordinates to walk through the vector representation of the matrix a and b, and calculate the product and sum of each term

4.  Assign the result to the vector position in result vector c.

2021