



ECE408/CS483/CSE408 Fall 2022

Applied Parallel Programming

Lecture 6:
Part 1: Generalized Tiling
Part 2: DRAM Bandwidth

ECE408/CS483/CSE408 Fall 2022

Applied Parallel Programming

Lecture 6: Generalized Tiling // DRAM Bandwidth

Tiled Matrix Multiplication Kernel

```
__global__ void MatrixMulKernel(float* M, float* N, float* P, int Width)
{
1.  __shared__ float subTileM[TILE_WIDTH][TILE_WIDTH];
2.  __shared__ float subTileN[TILE_WIDTH][TILE_WIDTH];

3.  int bx = blockIdx.x;  int by = blockIdx.y;
4.  int tx = threadIdx.x; int ty = threadIdx.y;

    // Identify the row and column of the P element to work on
5.  int Row = by * TILE_WIDTH + ty;
6.  int Col = bx * TILE_WIDTH + tx;
7.  float Pvalue = 0;

    // Loop over the M and N tiles required to compute the P element
    // The code assumes that the Width is a multiple of TILE WIDTH!
8.  for (int q = 0; q < Width/TILE_WIDTH; ++q) {
        // Collaborative loading of M and N tiles into shared memory
9.      subTileM[ty][tx] = M[Row*Width + (q*TILE_WIDTH+tx)];
10.     subTileN[ty][tx] = N[(q*TILE_WIDTH+ty)*Width+Col];
11.     __syncthreads();
12.     for (int k = 0; k < TILE_WIDTH; ++k)
13.         Pvalue += subTileM[ty][k] * subTileN[k][tx];
14.     __syncthreads();
15. }
16. P[Row*Width+Col] = Pvalue;
}
```

Memory Bandwidth Consumption

- Using 16x16 tiling, we reduce the global memory by a factor of 16
 - Each global memory load is used by 16 floating-point operations
 - The 150GB/s bandwidth can now support $(150/4)*16 = 600$ GFLOPS!
- Using 32x32 tiling, we reduce the global memory accesses by a factor of 32
 - Each global memory load is used by 32 floating-point operations
 - The 150 GB/s bandwidth can now support $(150/4)*32 = 1,200$ GFLOPS!
 - The memory bandwidth is no longer a limiting factor for performance!

Shared Memory and Threading

- Each SM in Maxwell has 64KB shared memory (48KB max per block)
 - Shared memory size is implementation dependent!
 - For `TILE_WIDTH = 16`, each thread block uses $2 \times 256 \times 4B = 2KB$ of shared memory.
 - Shared memory can potentially support up to 32 active blocks
 - The threads per SM constraint (2048) will limit the number of blocks to 8
 - For `TILE_WIDTH = 32`, each thread block uses $2 \times 32 \times 32 \times 4B = 8KB$ of shared memory
 - Shared memory can potentially support up to 8 active blocks
 - The threads per SM constraint (2048) will limit the number of blocks to 2

Device Query

- Number of devices in the system

```
int dev_count;  
cudaGetDeviceCount( &dev_count);
```

- Capability of devices

```
cudaDeviceProp dev_prop;  
for (i = 0; i < dev_count; i++) {  
    cudaGetDeviceProperties( &dev_prop, i);  
}
```

- `cudaDeviceProp` is a built-in C structure type

- `dev_prop.maxThreadsPerBlock`
- `dev_prop.sharedMemoryPerBlock`

Handling Matrix of Arbitrary Size

The tiled matrix multiplication kernel in Lecture 5 can handle only the matrices whose dimensions are multiples of the tile dimensions

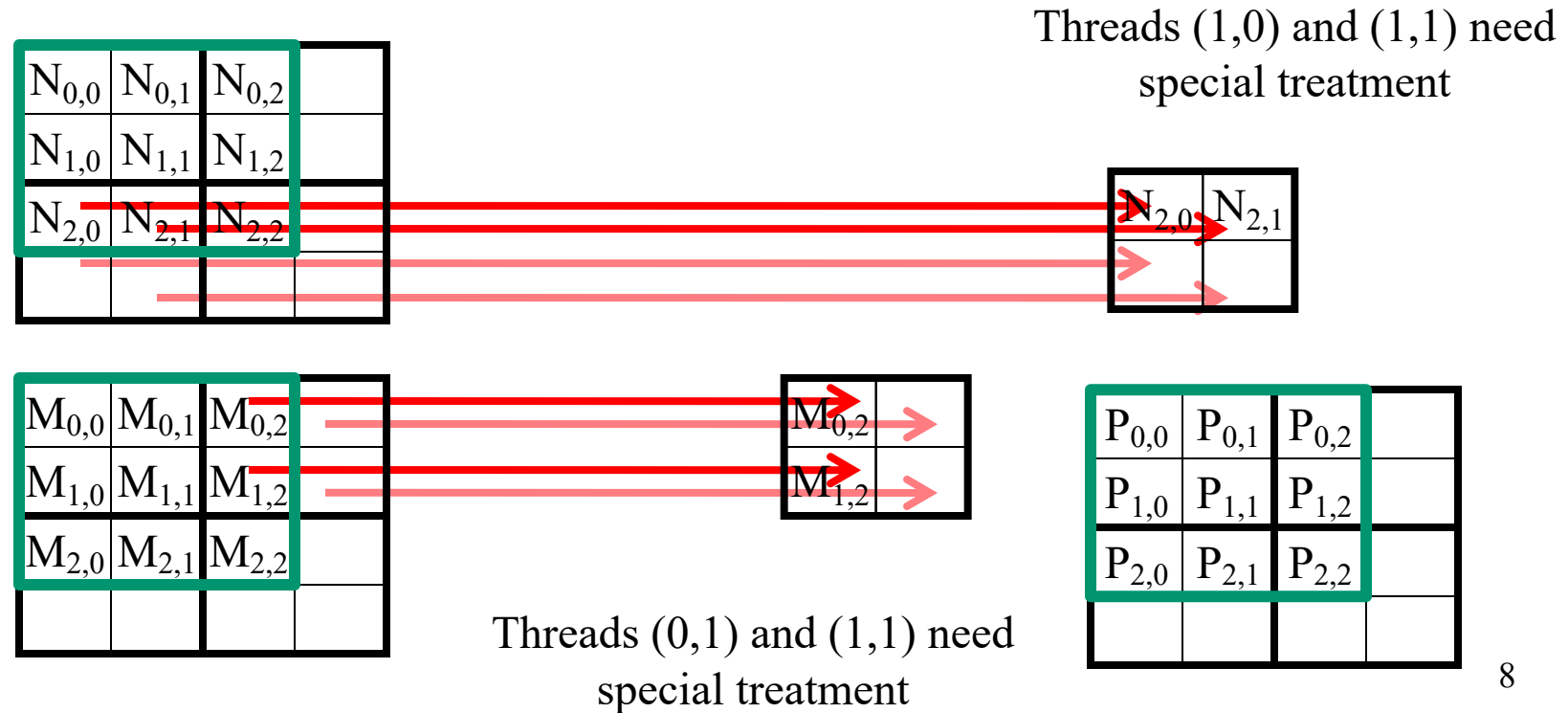
However, a real application needs to handle arbitrary sized matrices.

- We could pad (add elements to) the rows and columns into multiples of block size, but will have significant space and data transfer time overhead.
- We could add explicit checks in the code to handle boundaries

2x2 Tile on a 3x3 Multiply

`TILE_WIDTH = 2, width = 3`

Load of 2nd tile of Block (0,0)

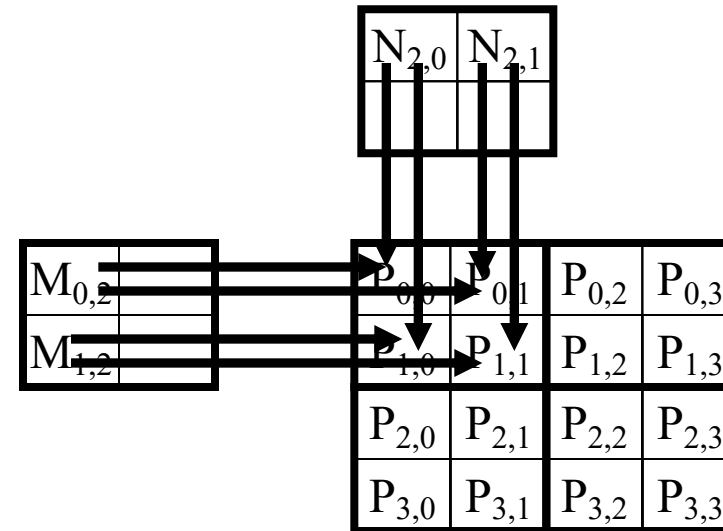


2x2 Tile on a 3x3 Multiply

Block (0,0), 2nd tile

$N_{0,0}$	$N_{0,1}$	$N_{0,2}$	
$N_{1,0}$	$N_{1,1}$	$N_{1,2}$	
$N_{2,0}$	$N_{2,1}$	$N_{2,2}$	

$M_{0,0}$	$M_{0,1}$	$M_{0,2}$	
$M_{1,0}$	$M_{1,1}$	$M_{1,2}$	
$M_{2,0}$	$M_{2,1}$	$M_{2,2}$	



2x2 Tile on a 3x3 Multiply

Block (0,0), 2nd tile

$N_{0,0}$	$N_{0,1}$	$N_{0,2}$	
$N_{1,0}$	$N_{1,1}$	$N_{1,2}$	
$N_{2,0}$	$N_{2,1}$	$N_{2,2}$	

$M_{0,0}$	$M_{0,1}$	$M_{0,2}$	
$M_{1,0}$	$M_{1,1}$	$M_{1,2}$	
$M_{2,0}$	$M_{2,1}$	$M_{2,2}$	

$M_{0,2}$	\equiv
$M_{1,2}$	\equiv

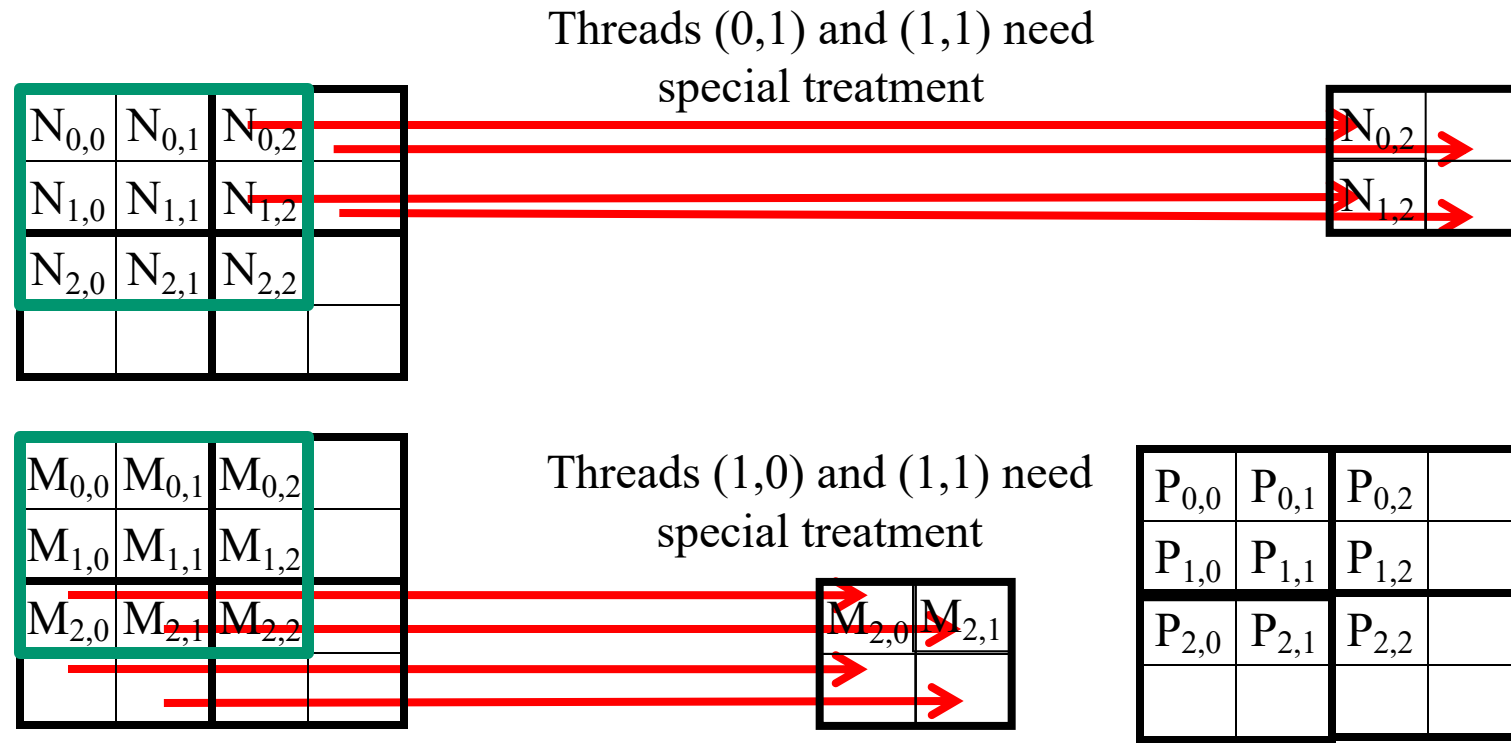
None of the threads should
take effect in this step

$N_{2,0}$	$N_{2,1}$

$P_{0,0}$	$P_{0,1}$	$P_{0,2}$	$P_{0,3}$
$P_{1,0}$	$P_{1,1}$	$P_{1,2}$	$P_{1,3}$
$P_{2,0}$	$P_{2,1}$	$P_{2,2}$	$P_{2,3}$
$P_{3,0}$	$P_{3,1}$	$P_{3,2}$	$P_{3,3}$

2x2 Tile on a 3x3 Multiply

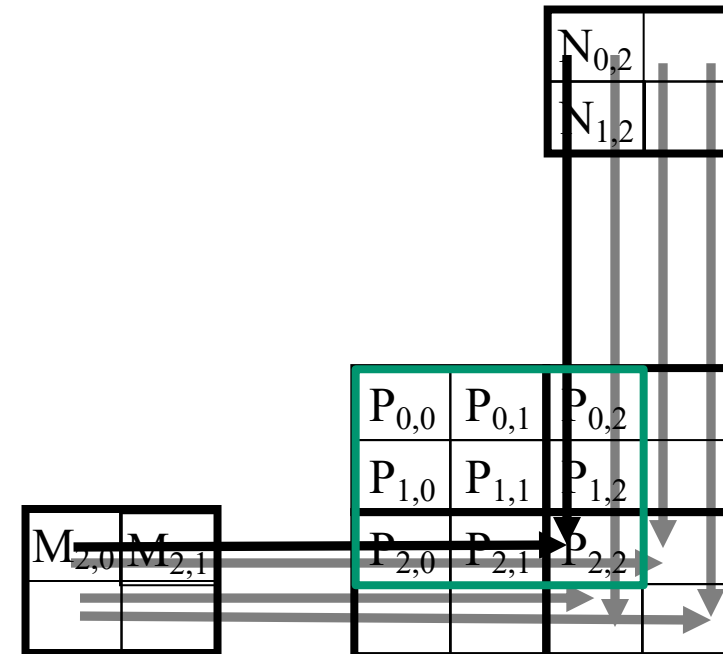
Load 1st Tile of Block (1,1)



1st Tile for Block (1,1)

$N_{0,0}$	$N_{0,1}$	$N_{0,2}$	
$N_{1,0}$	$N_{1,1}$	$N_{1,2}$	
$N_{2,0}$	$N_{2,1}$	$N_{2,2}$	

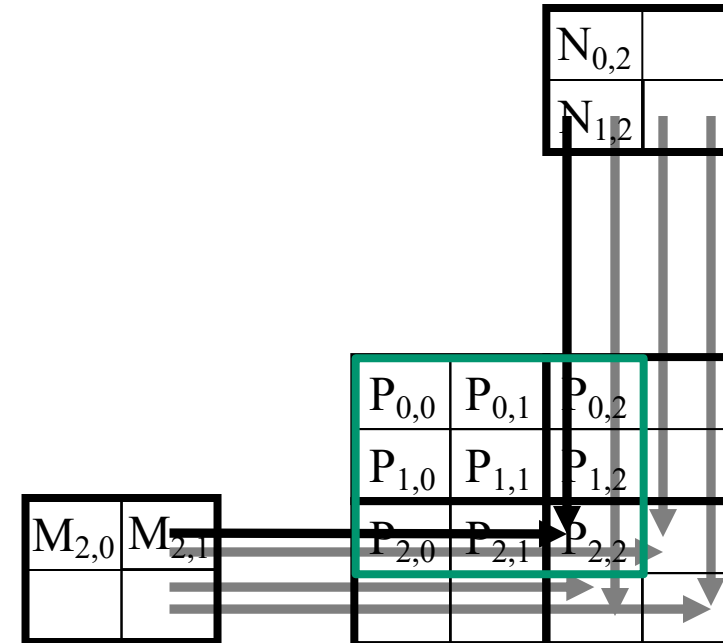
$M_{0,0}$	$M_{0,1}$	$M_{0,2}$	
$M_{1,0}$	$M_{1,1}$	$M_{1,2}$	
$M_{2,0}$	$M_{2,1}$	$M_{2,2}$	



1st Tile for Block (1,1)

$N_{0,0}$	$N_{0,1}$	$N_{0,2}$	
$N_{1,0}$	$N_{1,1}$	$N_{1,2}$	
$N_{2,0}$	$N_{2,1}$	$N_{2,2}$	

$M_{0,0}$	$M_{0,1}$	$M_{0,2}$	
$M_{1,0}$	$M_{1,1}$	$M_{1,2}$	
$M_{2,0}$	$M_{2,1}$	$M_{2,2}$	



Major Cases in 2x2 Example

- Threads that calculate valid P elements, but use invalid input
 - 1st Tile of Block(0,0), 2nd step, all threads
- Threads that calculate invalid P elements
 - Block(1,1), Thread(1,0), non-existent row
 - Block(1,1), Thread(0,1), non-existing column
 - Block(1,1), Thread(1,1), non-existing row/column

A “Simple” Solution

- Invalid input element, “load” a 0
 - Rationale: a 0 value will ensure that the multiply-add step does not affect the final value of the output element
- Invalid output element, don’t update global memory
 - Can still perform pvalue calculation (partial dot product), but doesn’t write to the global memory at the end of the kernel

2nd Tile for Block (0,0)

$N_{0,0}$	$N_{0,1}$	$N_{0,2}$	
$N_{1,0}$	$N_{1,1}$	$N_{1,2}$	
$N_{2,0}$	$N_{2,1}$	$N_{2,2}$	

$M_{0,0}$	$M_{0,1}$	$M_{0,2}$	
$M_{1,0}$	$M_{1,1}$	$M_{1,2}$	
$M_{2,0}$	$M_{2,1}$	$M_{2,2}$	

$N_{2,0}$	$N_{2,1}$
0	0

$M_{0,2}$	0
$M_{1,2}$	0

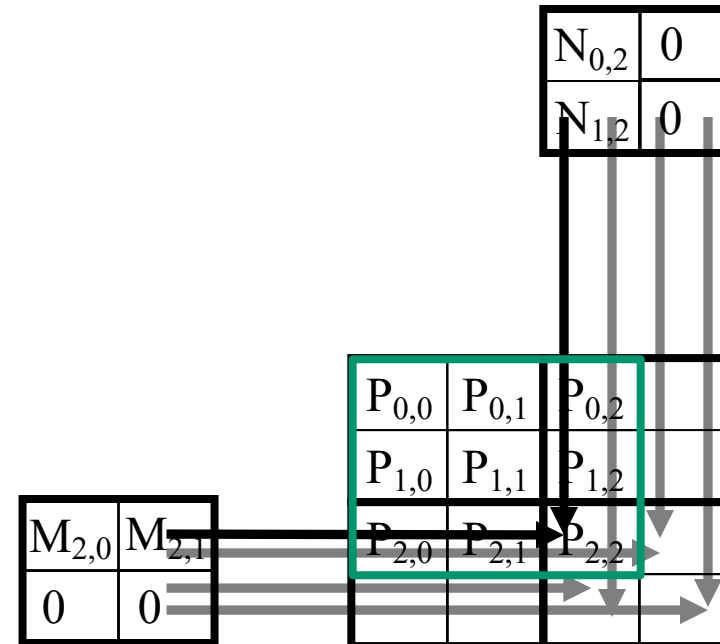
$P_{0,0}$	$P_{0,1}$	$P_{0,2}$	$P_{0,3}$
$P_{1,0}$	$P_{1,1}$	$P_{1,2}$	$P_{1,3}$
$P_{2,0}$	$P_{2,1}$	$P_{2,2}$	$P_{2,3}$
$P_{3,0}$	$P_{3,1}$	$P_{3,2}$	$P_{3,3}$

The multiply-add will not affect the output due to 0

1st Tile for Block (1,1)

$N_{0,0}$	$N_{0,1}$	$N_{0,2}$	
$N_{1,0}$	$N_{1,1}$	$N_{1,2}$	
$N_{2,0}$	$N_{2,1}$	$N_{2,2}$	

$M_{0,0}$	$M_{0,1}$	$M_{0,2}$	
$M_{1,0}$	$M_{1,1}$	$M_{1,2}$	
$M_{2,0}$	$M_{2,1}$	$M_{2,2}$	



Tiled Matrix Multiplication Kernel

```
// Loop over the M and N tiles required to compute the P element
// The code assumes that the Width is a multiple of TILE_WIDTH!
8.  for (int q = 0; q < (ceil((float)Width/TILE_WIDTH)); ++q) {
    // Collaborative loading of M and N tiles into shared memory
    if (Row < Width && (q*TILE_WIDTH+tx) < Width)
9.      subTileM[ty][tx] = M[Row*Width + q*TILE_WIDTH+tx];
10.  else
11.      subTileM[ty][tx] = 0;

12.  if (Col < Width && (q*TILE_WIDTH+ty) < Width)
13.      subTileN[ty][tx] = N[(q*TILE_WIDTH+ty)*Width+Col];
14.  else
15.      subTileN[ty][tx] = 0;
16.  __syncthreads();

17.  for (int k = 0; k < TILE_WIDTH; ++k)
18.      Pvalue += subTileM[ty][k] * subTileN[k][tx];
19.  __syncthreads();
20. }
21. if (Row < Width && Col < Width)
22.  P[Row*Width+Col] = Pvalue;
}
```

Some Important Points

- For each thread the conditions are different for
 - Loading M element
 - Loading N element
 - Calculation/storing output elements
- The effect of control divergence should be small for large matrices
- How about rectangular matrices?

Global Memory Bandwidth

Ideal

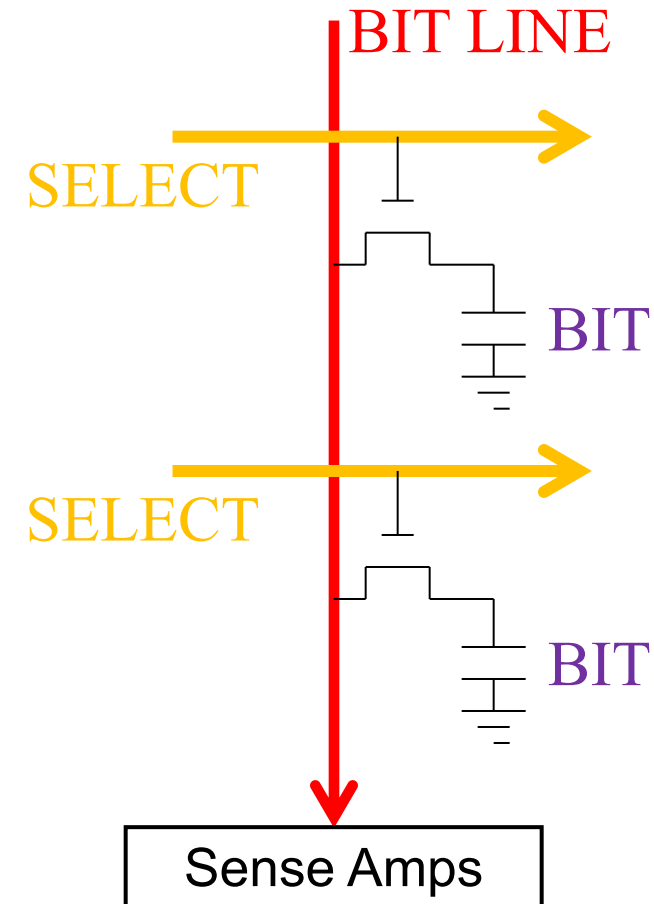


Reality

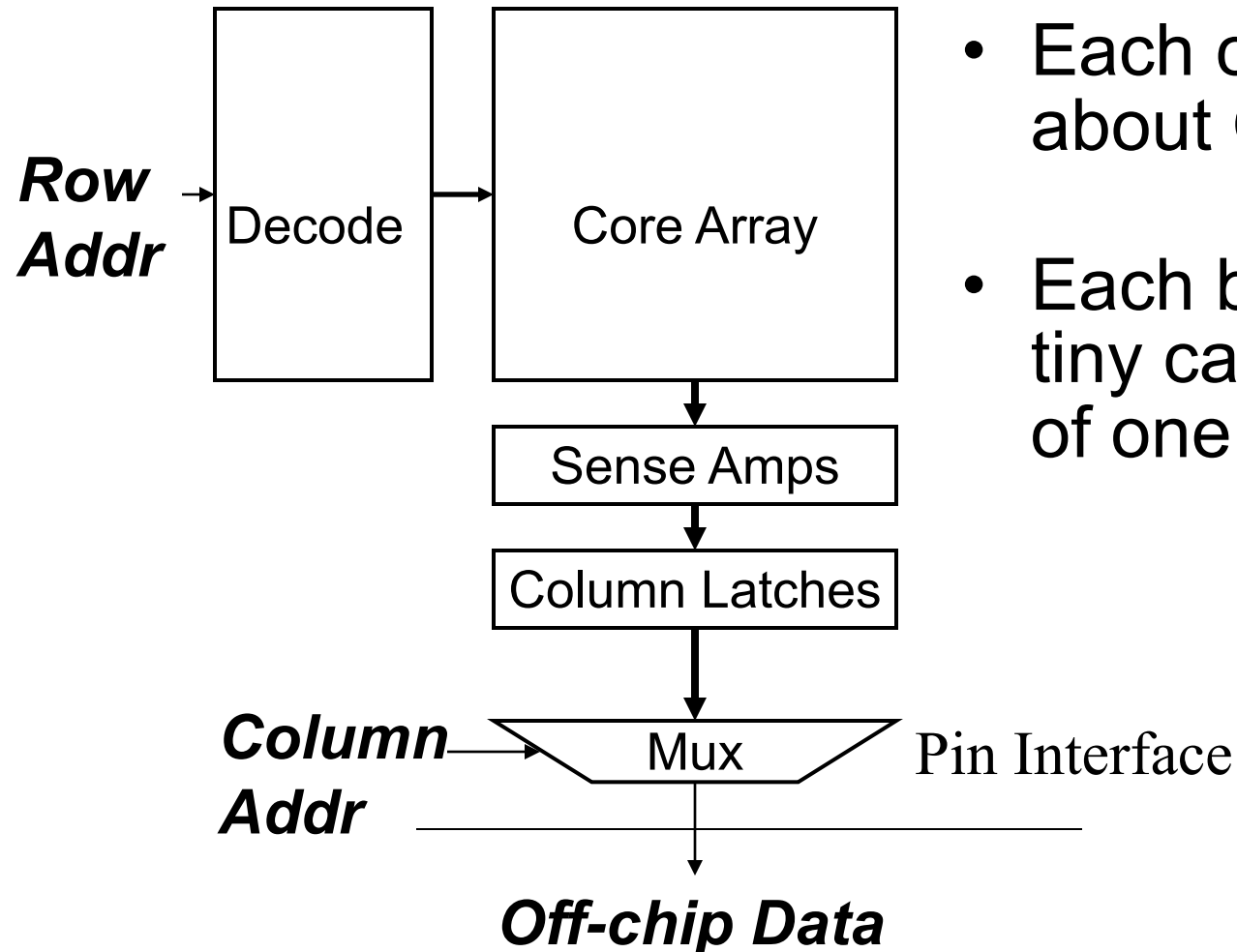


DRAM is Slow But Dense

- Capacitance...
 - tiny for the **BIT**, but
 - huge for the **BIT LINE**
- Use an amplifier for higher speed!
- Still **slow**...
- But only need **1 transistor per bit.**

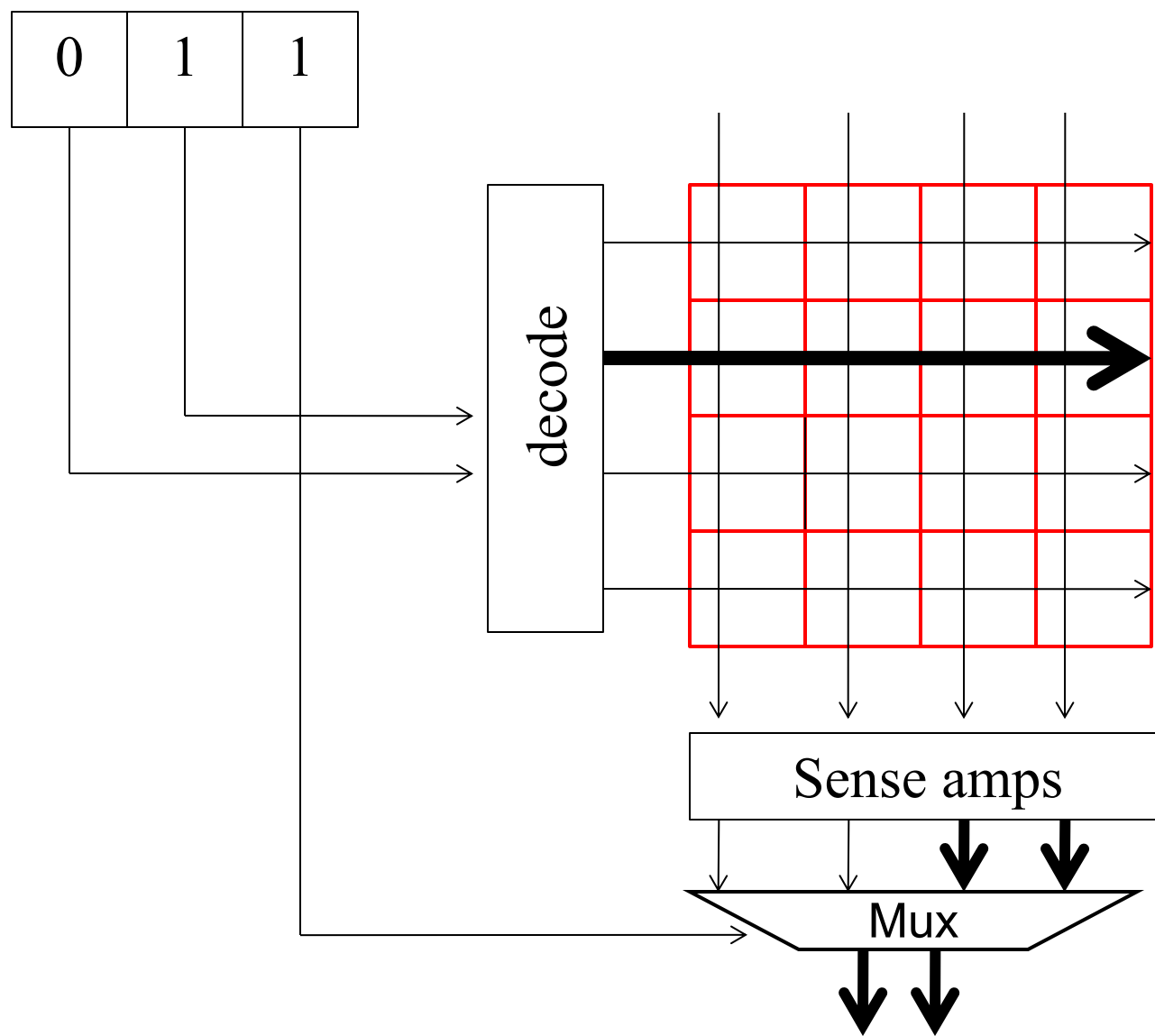


DRAM Bank Organization



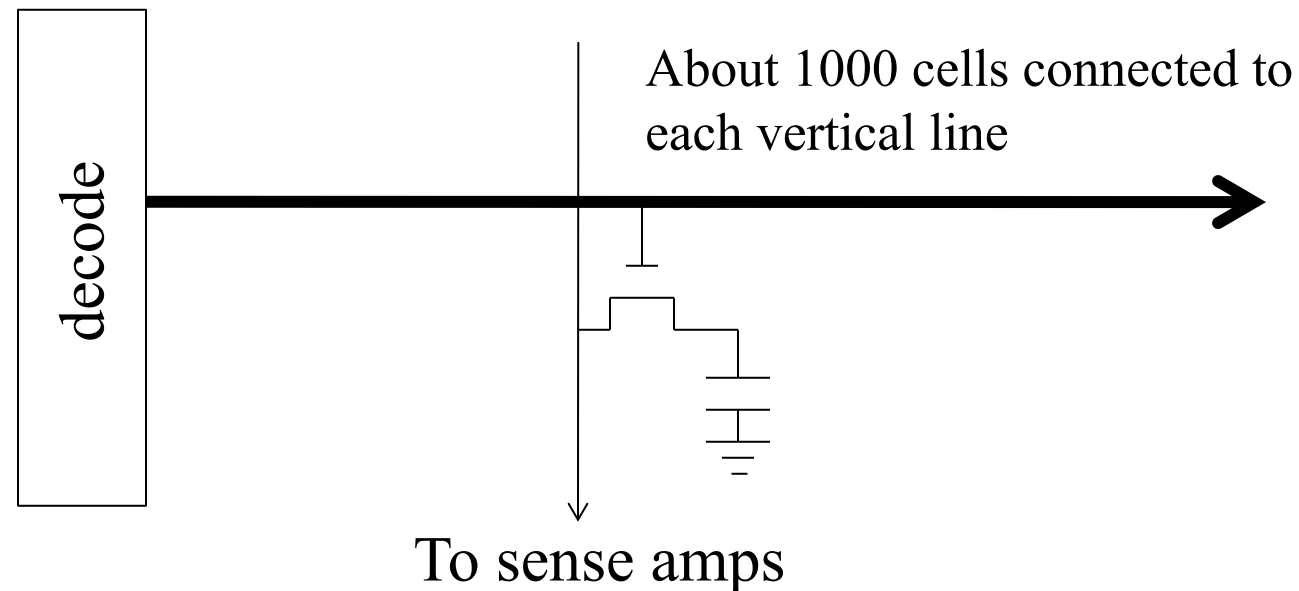
- Each core array has about $O(1M)$ bits
- Each bit is stored in a tiny capacitor, made of one transistor

A very small (8x2 bit) DRAM Bank

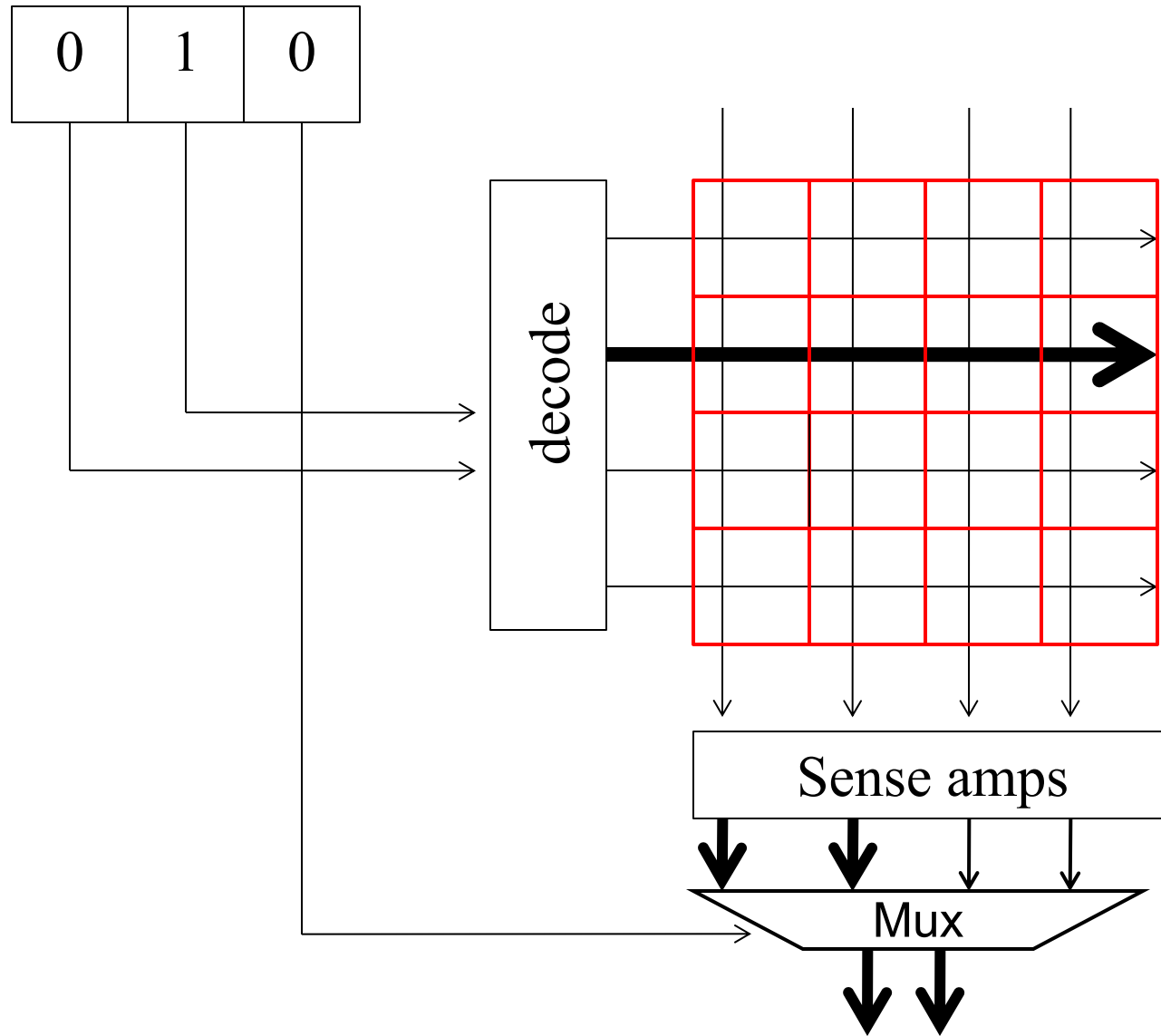


DRAM core arrays are slow.

- Reading from a cell in the core array is a very slow process
 - Current GDDR: Core speed = $\frac{1}{8}$ interface speed
 - ... likely to be worse in the future

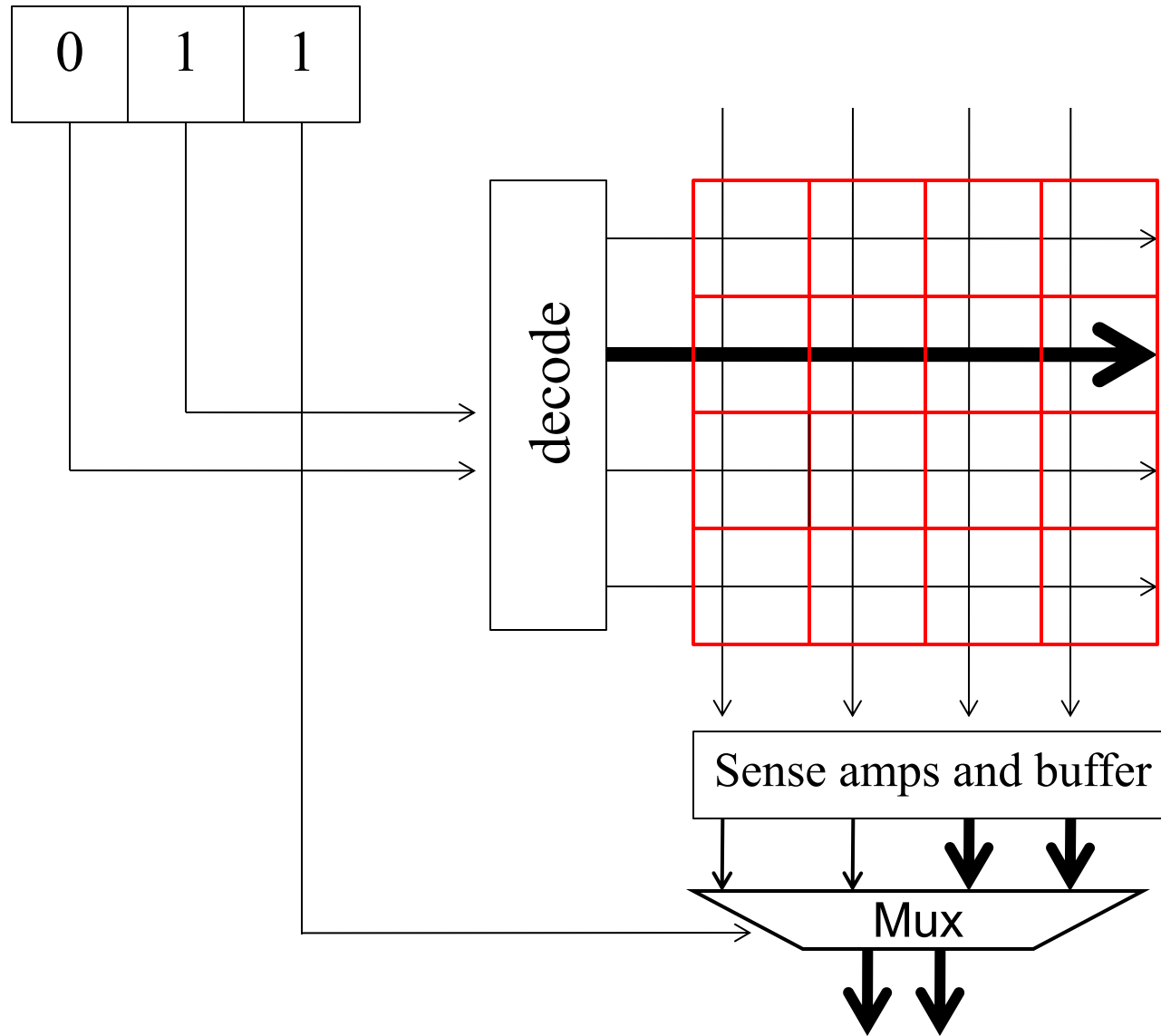


DRAM Bursting (burst size = 4 bits)

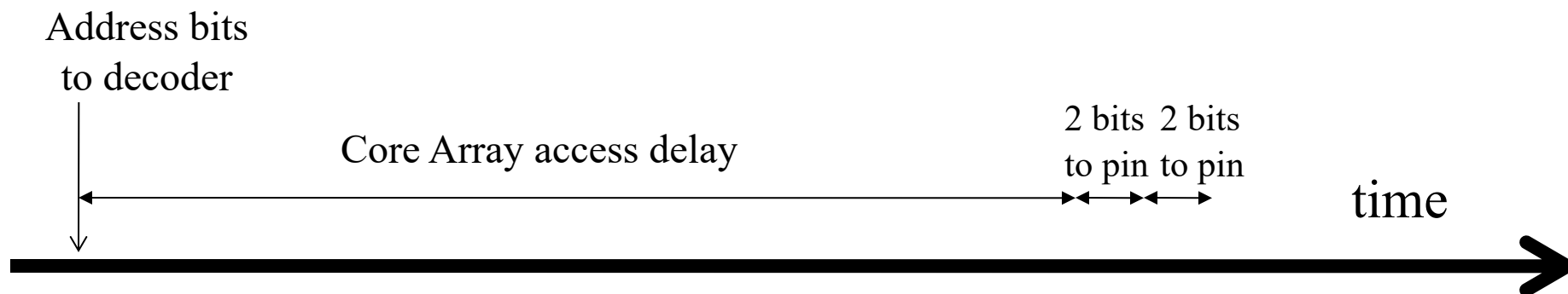


DRAM Bursting (cont.)

second part of the burst



DRAM Bursting for the 8x2 Bank



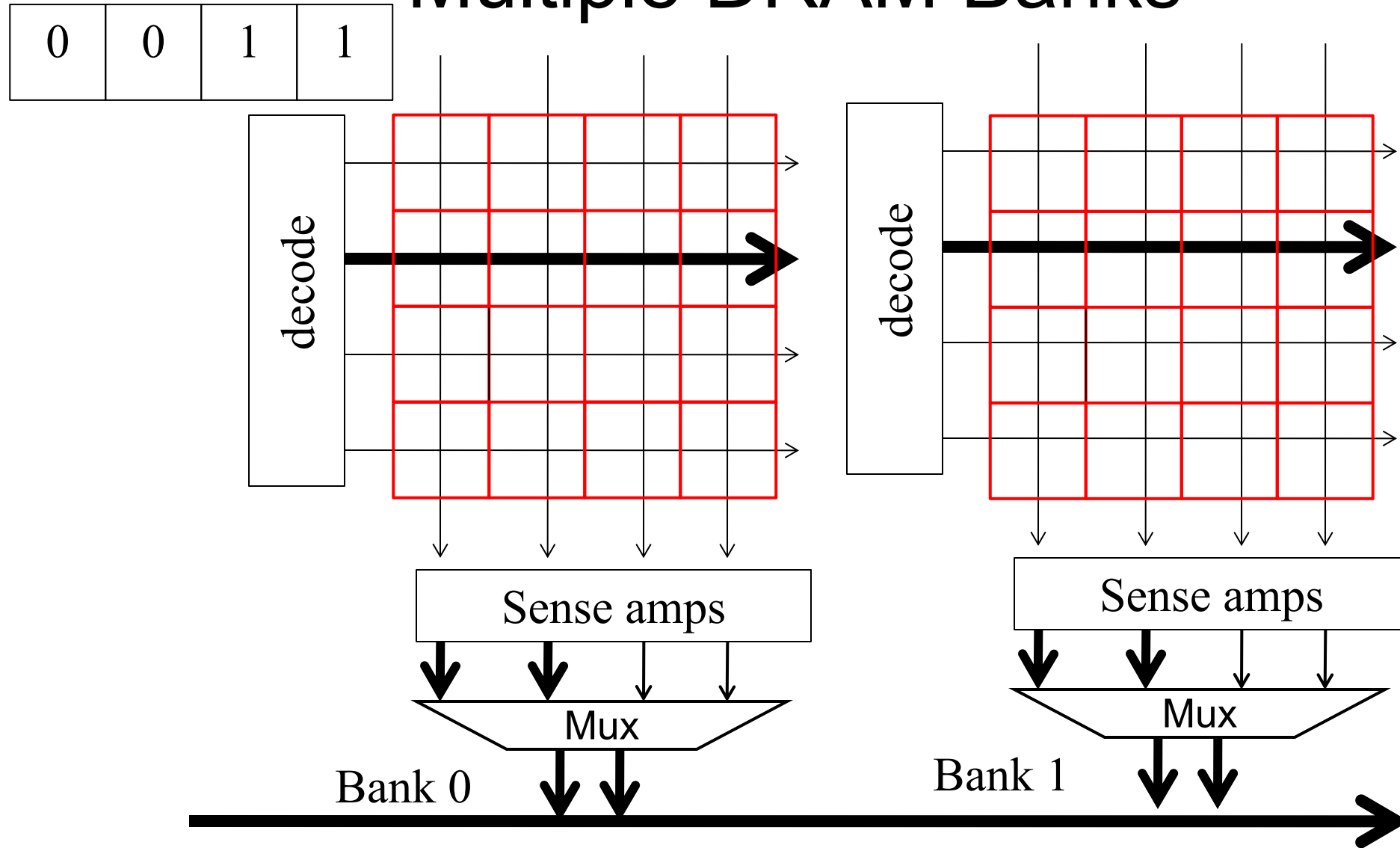
Non-burst timing



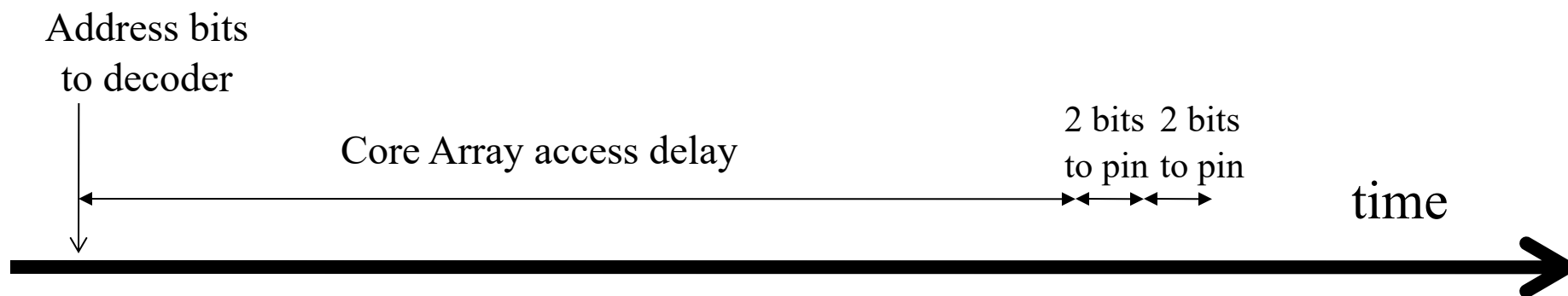
Burst timing

Modern DRAM systems are designed to be always accessed in burst mode. Burst bytes are transferred but discarded when accesses are not to sequential locations.

Multiple DRAM Banks



DRAM Bursting for the 8x2 Bank

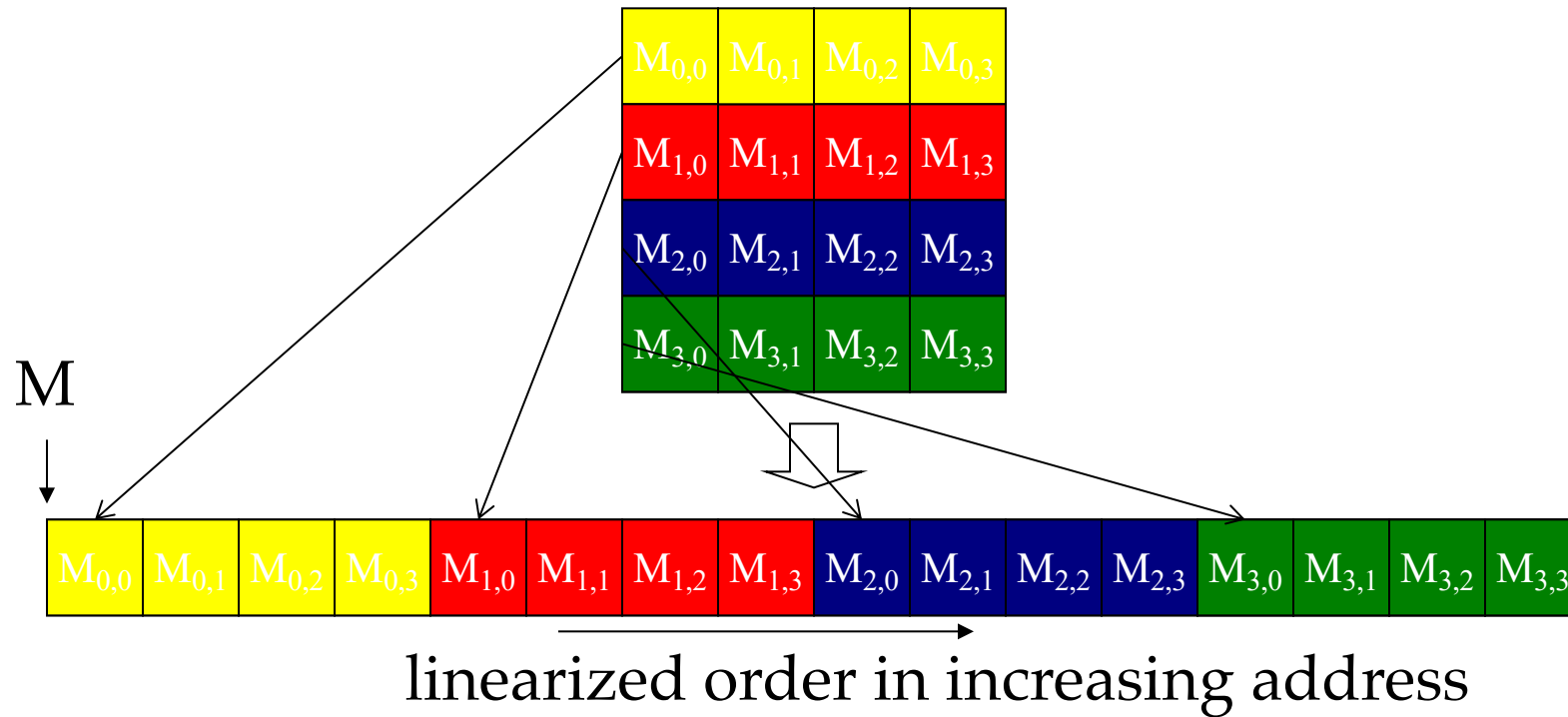


Single-Bank burst timing, dead time on interface



Multi-Bank burst timing, reduced dead time

Placing a 2D C array into linear memory space (review)



A Simple Matrix Multiplication Kernel

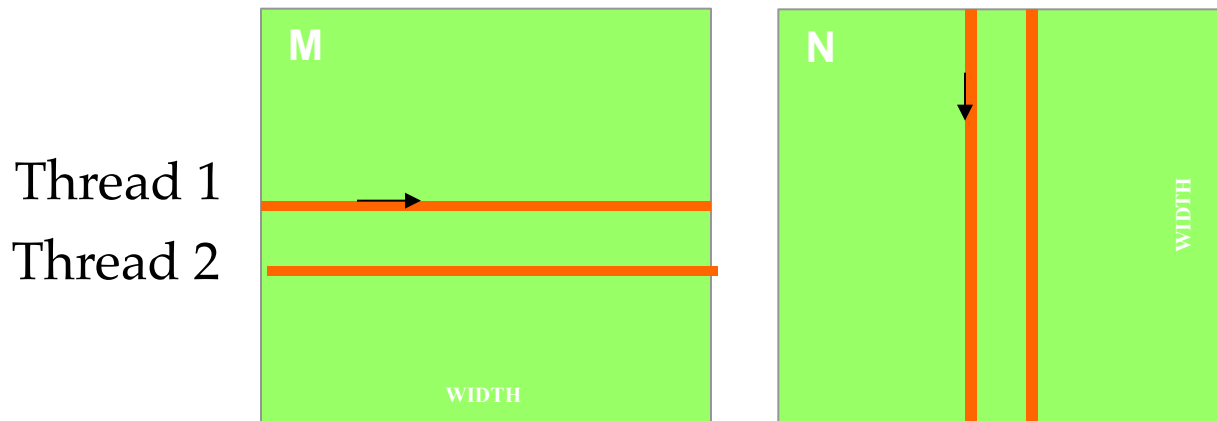
```
__global__ void MatrixMulKernel(float* M, float* N, float* P, int Width)
{
    // Calculate the row index of the P element and M
    int Row = blockIdx.y * blockDim.y + threadIdx.y;
    // Calculate the column index of P and N
    int Col = blockIdx.x * blockDim.x + threadIdx.x;

    if ((Row < Width) && (Col < Width)) {
        float Pvalue = 0;

        // each thread computes one element of the block sub-matrix
        for (int k = 0; k < Width; ++k)
            Pvalue += M[Row*Width+k] * N[k*Width+Col];

        P[Row*Width+Col] = Pvalue;
    }
}
```

Two Access Patterns

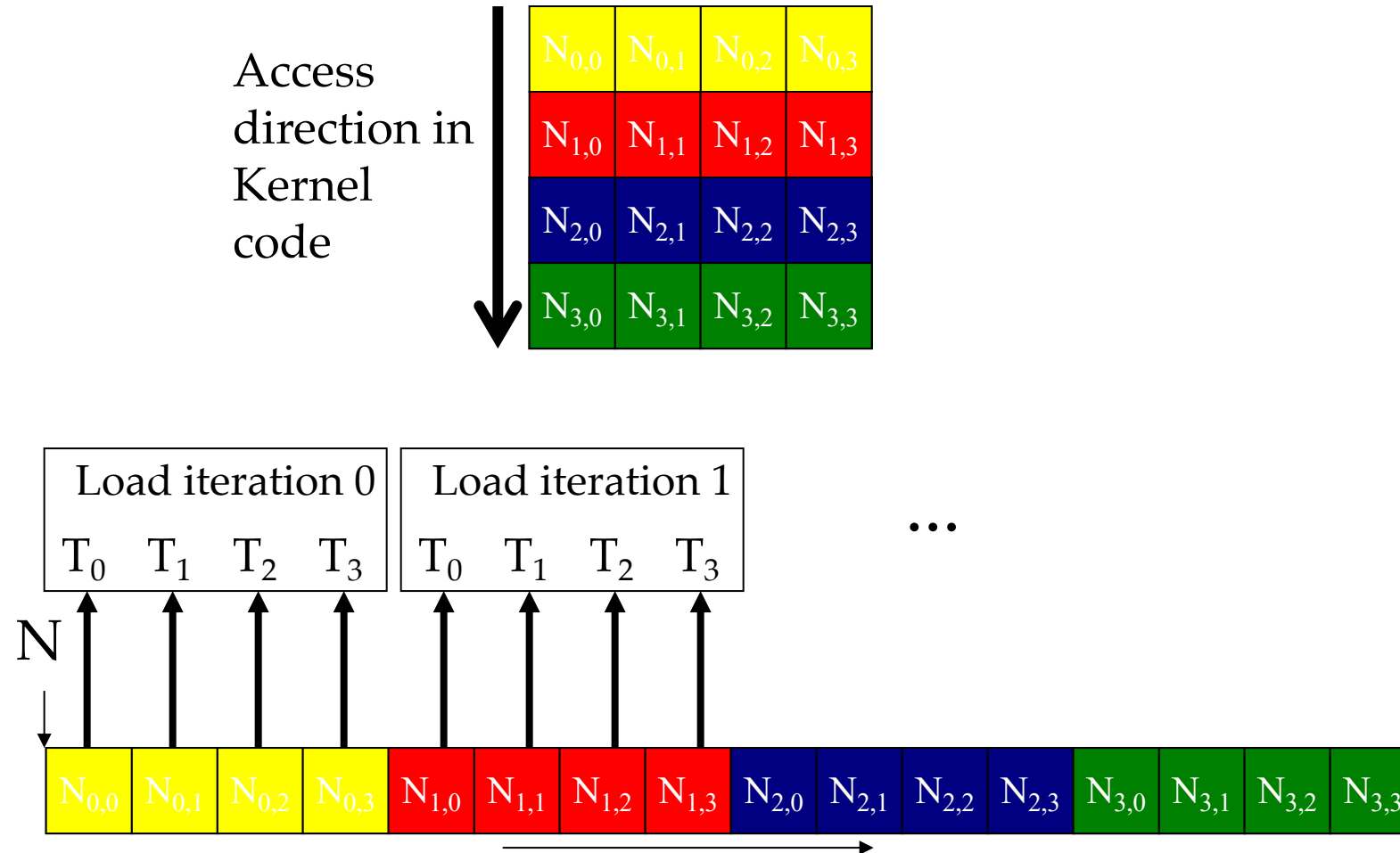


$M[\text{Row} * \text{Width} + k]$

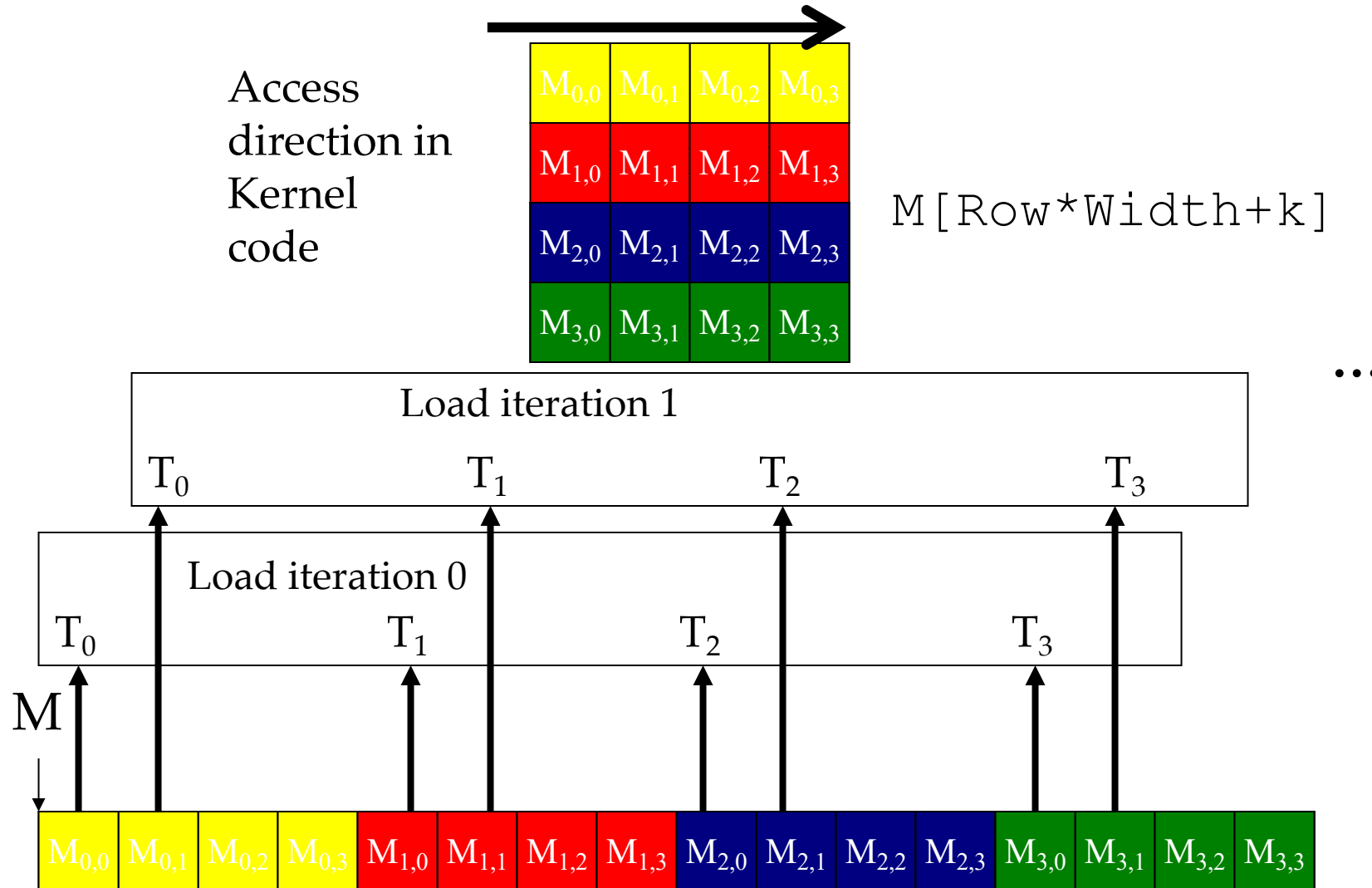
$N[k * \text{Width} + \text{Col}]$

k is loop counter in the inner product loop of the kernel code

N accesses are coalesced.

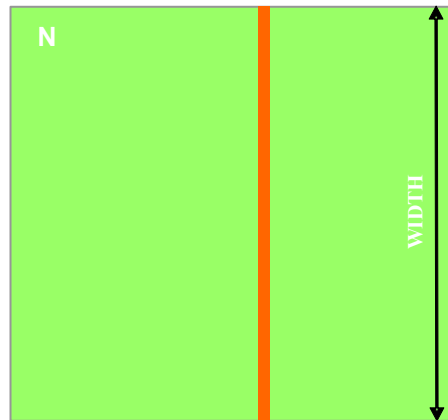
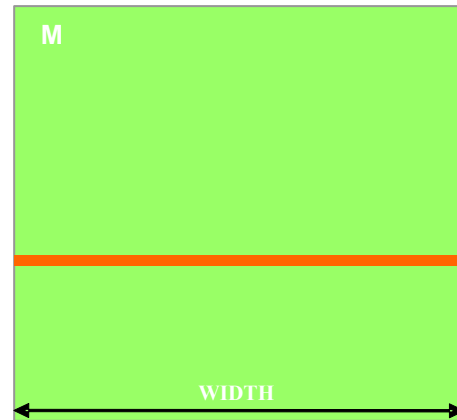


M accesses are not coalesced.

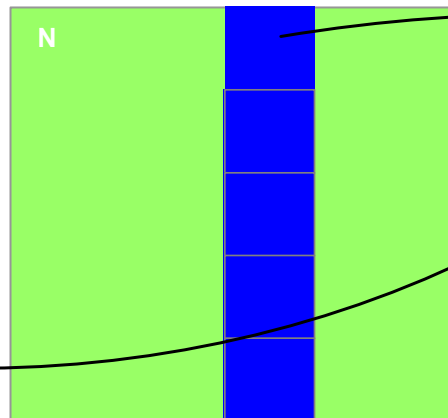
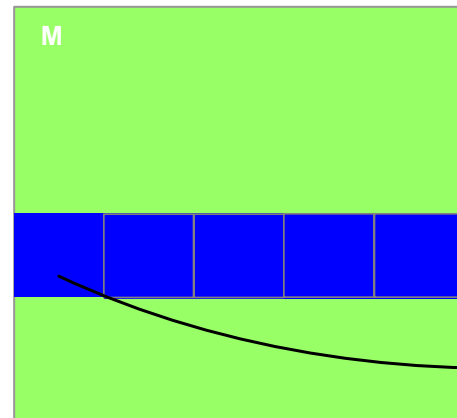


Use shared memory to enable coalescing in tiled matrix multiplication

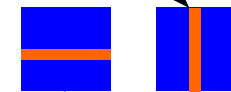
Original
Access
Pattern



Tiled
Access
Pattern



Copy into
Shared memory



Perform
multiplication
with shared memory
values

Two vertical lines, one blue and one orange, are positioned on the left side of the slide.

**ANY MORE QUESTIONS?
READ CHAPTER 5**