# Memory Coalescing

Instructions are issued in parallel at the warp level of 32 threads

Warp

Warp

Data

Warp



| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |

Data

For these slides we will treat 4 data elements as one of these fixed-length lines of contiguous memory

Warp

Data

The memory subsystem will attempt to minimize the number of lines required to fulfill the read/write requirements of the warp

Warp

| | | | |
|---|---|---|---|

Data

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|

All data in the line will be used

Warp

```
idx = threadIdx.x +
blockIdx.x * blockDim.x

a[idx] += 1
```

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |

Data

And the transfer will happen in as few lines as possible

Warp

```
idx = threadIdx.x +
blockIdx.x * blockDim.x

a[idx] += 1
```

Data

Warp

```
idx = threadIdx.x +
blockIdx.x * blockDim.x

a[idx] += 1
```



| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |

Data

Warp

More lines will have to be transferred to fulfil the needs of the warp

```
idx = blockIdx.x +
blockDim.x * threadIdx.x

a[idx] += 1
```

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |

Data

And more of the data being transferred will go unused

Warp

```
idx = blockIdx.x +
blockDim.x * threadIdx.x

a[idx] += 1
```

Data

The memory throughput is degraded, and additional time is required: a performance loss

Warp

```
idx = blockIdx.x +
blockDim.x * threadIdx.x

a[idx] += 1
```

Data

# Row and Column Sum Comparison

Consider a kernel that stores the sum of each row of a matrix (which here is 4 contiguous data elements) in a result vector

## Warp

## Result

| ? | ? | ? | ? |

## Data

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

A single thread could iterate over a row, summing it, and then write the result in the solution vector

Warp

Result

Data

A single thread could iterate over a row, summing it, and then write the result in the solution vector

Warp



Data

Result

| ? | ? | ? | ? |

Sum = 0

A single thread could iterate over a row, summing it, and then write the result in the solution vector

Warp

Result

? ? ? ?

Data

Sum = 6

A single thread could iterate over a row, summing it, and then write the result in the solution vector

Warp

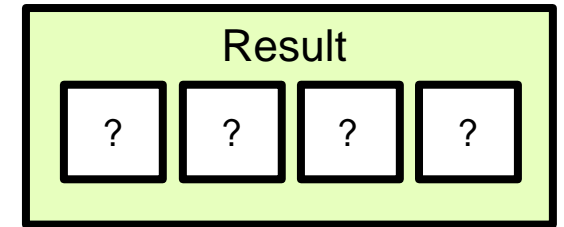Result

6 | ? | ? | ?

Sum = 6

Data

0 | 1 | 2 | 3
4 | 5 | 6 | 7
8 | 9 | 10 | 11
12 | 13 | 14 | 15

This seems natural, but look at what happens when we consider the parallel execution within the warp

Warp

Data

Which means (in our example) 4 lines of data will need to be loaded, and 75% of the data loaded will be unused

Warp

Data

Memory Line Size

Unfortunately, as each thread iterates over its row, the same uncoalesced pattern continues

Warp

Data

Memory Line Size

In this example we transferred 16 memory lines, and used 25% of the data for each line transferred

Warp



Data

|  |  |  |  |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Let's compare a kernel that stores the sum of each **column** of a matrix in a result vector

Warp

Result

? ? ? ?

Data

0 1 2 3

4 5 6 7

8 9 10 11

12 13 14 15

A single thread could iterate over a column, summing it, and then write the result in the solution vector

Warp

Result

? ? ? ?

Data

A single thread could iterate over a column, summing it, and then write the result in the solution vector

Warp

Data

Result

? ? ? ?

Sum = 0

A single thread could iterate over a column, summing it, and then write the result in the solution vector
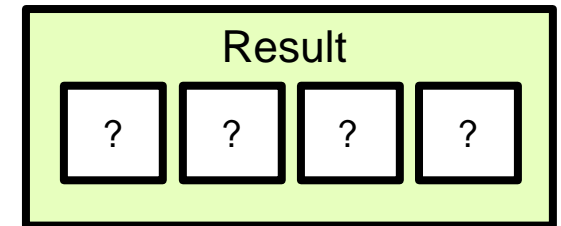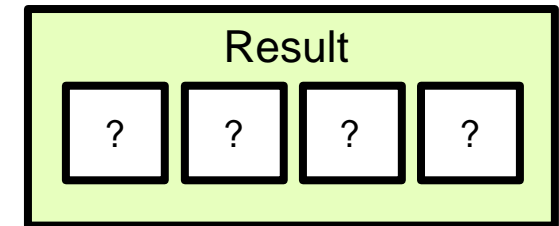
Warp



Result

? ? ? ?

Sum = 5

Data

A single thread could iterate over a column, summing it, and then write the result in the solution vector

Warp



Result

| ? | ? | ? | ? |

Sum = 24

Data

A single thread could iterate over a column, summing it, and then write the result in the solution vector

Warp

Result

24 | ? | ? | ?

Sum = 24

Data

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Here when we consider the parallel execution, we see that the warp's memory access is coalesced

Warp

Data

Memory Line Size

Here when we consider the parallel execution, we see that the warp's memory access is coalesced

Warp

Data

Memory Line Size

# Warp

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Data

Memory Line Size

In this example we transferred 4 memory lines (compared to 16), and used 100% of the data for each line transferred (compared to 25%)

Warp

Data

# Using Shared Memory to Support Coalesced Memory Access

We will examine a matrix transpose to demonstrate how shared memory can be used to promote coalesced data transfers to and from global memory

# Grid



Here we have a (2,2) grid, with each block containing (2,2) threads as well as (4,4) input and output matrices

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Input

Output

# Warp Size

# Memory Segment Size

# Grid

For these slides we will define a warp as 2 threads, and a memory segment as 2 data elements wide

| | | | |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Input

Output

Warp Size

Memory Segment Size

Grid

Our goal is to transpose the input by rotating all elements around the diagonal, writing the transposed elements to output

| Input | | | |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

| Output | | | |
|---|---|---|---|
| 0 | 4 | 8 | 12 |
| 1 | 5 | 9 | 13 |
| 2 | 6 | 10 | 14 |
| 3 | 7 | 11 | 15 |

Warp Size

Memory Segment Size

Grid

A naïve approach is to launch a grid with threads equal to input elements, and to have each thread read 1 element, then write it to output in the transposed location

```
x, y = cuda.grid(2)

out[x][y] = in[y][x]
```

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Input

Output

## Warp Size

## Memory Segment Size

## Grid

Observing the behavior of a single warp, is it the case that memory reads are coalesced? Let's dig into answering that question

```
x, y = cuda.grid(2)

out[x][y] = in[y][x]
```

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Input

Output

## Warp Size

## Memory Segment Size

## Grid

Rewriting the creation of the indexing variables, it is clearer that contiguous threads in the same warp are adjacent along the x axis

```
x = blockIdx.x * blockDim.x
        + threadIdx.x

y = blockIdx.y * blockDim.y
        + threadIdx.y

out[x][y] = in[y][x]
```

| | | | |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Input

Output

Warp Size

Memory Segment Size

Grid

Furthermore, these contiguous threads will read elements from the rows of input where data elements are contiguous

```
x = blockIdx.x * blockDim.x
        + threadIdx.x

y = blockIdx.y * blockDim.y
        + threadIdx.y

out[x][y] = in[y][x]
```

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Input

Output

Warp Size

Memory Segment
Size

Grid

Therefore, it makes sense that reads
from input are coalesced

```
x = blockIdx.x * blockDim.x
        + threadIdx.x

y = blockIdx.y * blockDim.y
        + threadIdx.y

out[x][y] = in[y][x]
```

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Input

Output

## Warp Size

## Memory Segment Size

## Grid

What about this warp's writes to output, will they be coalesced?

```
x = blockIdx.x * blockDim.x
        + threadIdx.x

y = blockIdx.y * blockDim.y
        + threadIdx.y

out[x][y] = in[y][x]
```

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Input

Output

## Warp Size

## Memory Segment Size

## Grid

Here we see that contiguous threads in the same warp will be writing along a column in output

```
x = blockIdx.x * blockDim.x
        + threadIdx.x

y = blockIdx.y * blockDim.y
        + threadIdx.y

out[x][y] = in[y][x]
```

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Input

Output

Warp Size

Memory Segment
Size

Grid

Input

Output

```
x = blockIdx.x * blockDim.x
        + threadIdx.x

y = blockIdx.y * blockDim.y
        + threadIdx.y

out[x][y] = in[y][x]
```

Warp Size

Memory Segment Size

Shared

Grid

We can use shared memory to make coalesced reads and writes. Here, each block will allocate a (2,2) shared memory tile

```
tile = cuda.shared.array(2,2)
```

Input

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Output

**Shared**

**Grid**

(It is worth reminding that in our slides, to preserve space, 2 threads is a warp length. A real warp is 32 threads)

```
tile = cuda.shared.array(2,2)
```

**Warp Size**

**Memory Segment Size**

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

**Input**

**Output**

## Warp Size

## Memory Segment Size

## Shared

| 0 | 1 |
|---|---|
|   |   |

## Grid

## Input

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

## Output

Now we can make coalesced reads from input, and write the values to the block's shared memory tile

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
```

**Warp Size**

**Memory Segment Size**

**Shared**

**Grid**

Because each shared memory tile is local to the block (not the grid) we index into it using thread indices, not grid indices

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
```

| 0 | 1 |
|---|---|
|   |   |

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

**Input**

**Output**

**Warp Size**

**Memory Segment Size**

**Shared**

| 0 | 1 |
|---|---|
| 4 | 5 |

**Grid**

After synchronizing on all threads in the block, the tile will contain all the data this block needs to begin the writes

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()
```

**Input**

| 0  | 1  | 2  | 3  |
|----|----|----|----|
| 4  | 5  | 6  | 7  |
| 8  | 9  | 10 | 11 |
| 12 | 13 | 14 | 15 |

**Output**

**Warp Size**

**Memory Segment Size**

**Shared**

| | |
|---|---|
| 0 | 1 |
| 4 | 5 |

**Grid**

So that the writes are coalesced, we want each warp to write to a row in output

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
```

**Input**

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

**Output**

Warp Size

Memory Segment Size

Shared

| 0 | 1 |
|---|---|
| 4 | 5 |

Grid

Notice that to write to output at the transposed locations we use blockIdx.y and blockDim.y to calculate the x axis index in output…

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
```

Input

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Output

NVIDIA DEEP LEARNING INSTITUTE

Warp Size

Memory Segment Size

Shared

Grid

Input

Output

…but to accomplish coalesced writes, we still map increments to threadIdx.x to be along the x output axis

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
```

Warp Size

Memory Segment Size

Shared

Grid

Because of this last detail, each warp will need to read from a column of the shared memory tile in order to perform the transpose

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```

Input

Output

Warp Size

Memory Segment Size

Shared

Grid

Input

Output

(There's more to come about efficient reads/writes to/from shared memory, but for now know that reading across the column in shared memory has very low impact compared to doing so with global memory)

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```

## Warp Size

## Memory Segment Size

## Shared

## Grid

## Shared

In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

## Input

## Output

## Warp Size

## Memory Segment Size

## Shared

## Grid

## Shared

In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory

```python
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

## Input

## Output

Warp Size

Memory Segment Size

Shared

Grid

Shared

| 0 | 1 |
|---|---|
| 4 | 5 |

In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```
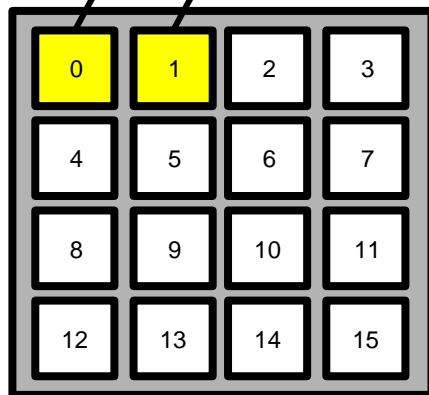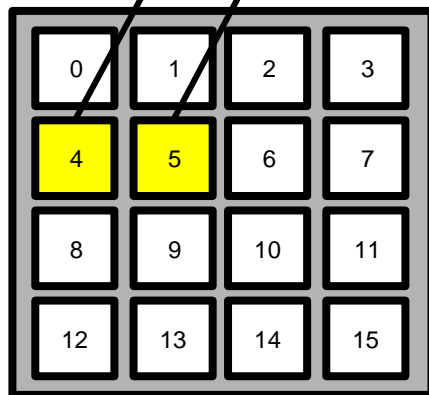
| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Input

Output

# Warp Size

# Memory Segment Size

## Shared

| | |
|---|---|
| 0 | 1 |
| 4 | 5 |
| 8 | 9 |
| | |

## Grid

## Shared

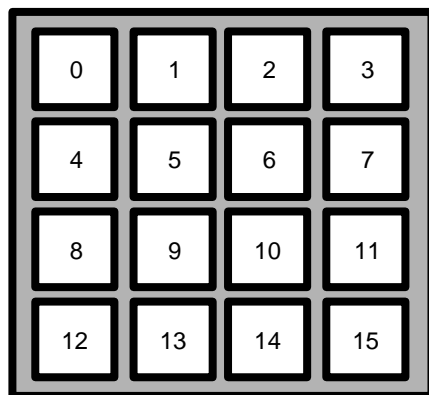In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory
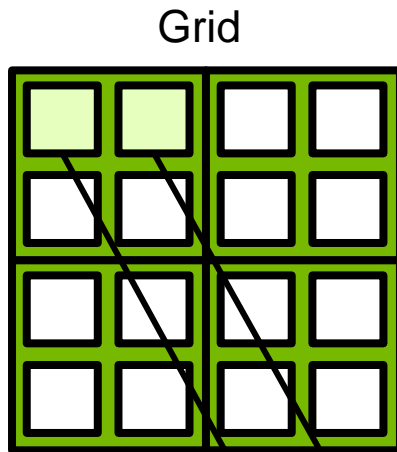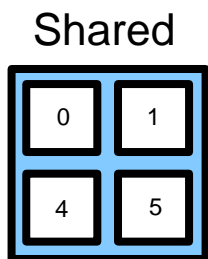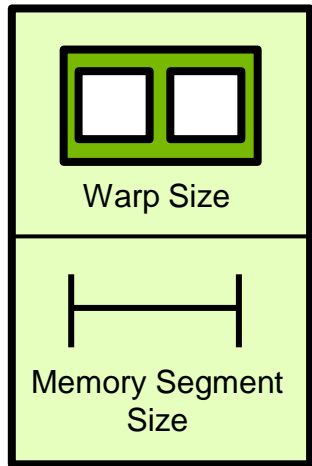
```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```

## Input

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

## Output

## Warp Size

## Memory Segment Size

## Shared

## Grid

## Shared

|    |    |
|----|----|
| 0  | 1  |
| 4  | 5  |
| 8  | 9  |
| 12 | 13 |

In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory
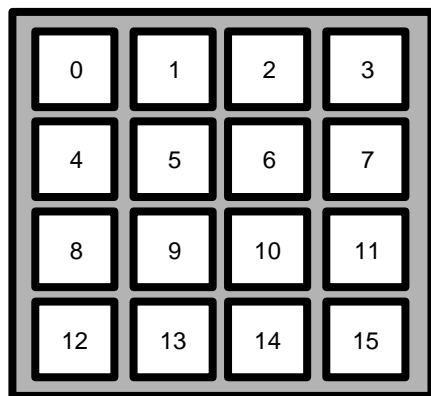
```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```
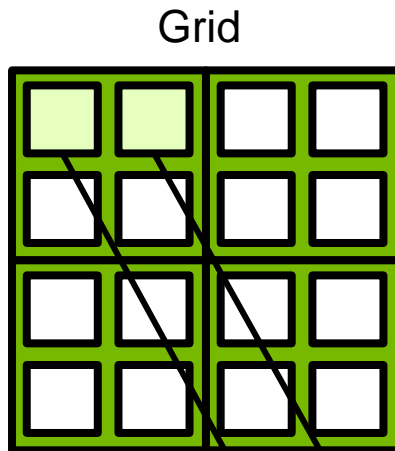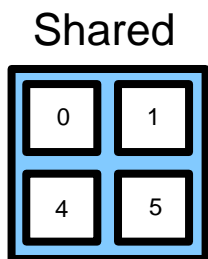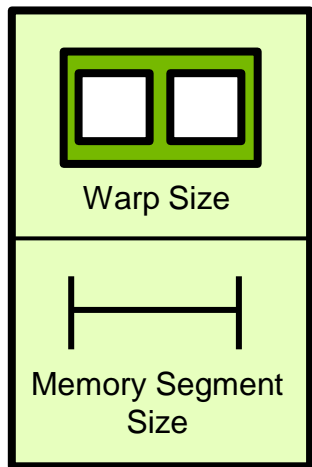
## Input

|    |    |    |    |
|----|----|----|----|
| 0  | 1  | 2  | 3  |
| 4  | 5  | 6  | 7  |
| 8  | 9  | 10 | 11 |
| 12 | 13 | 14 | 15 |

## Output

Warp Size

Memory Segment Size

Shared

Grid

Shared

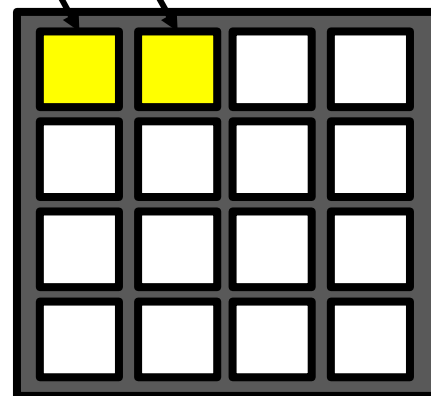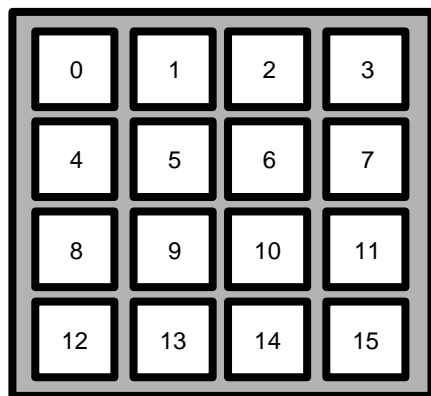In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory
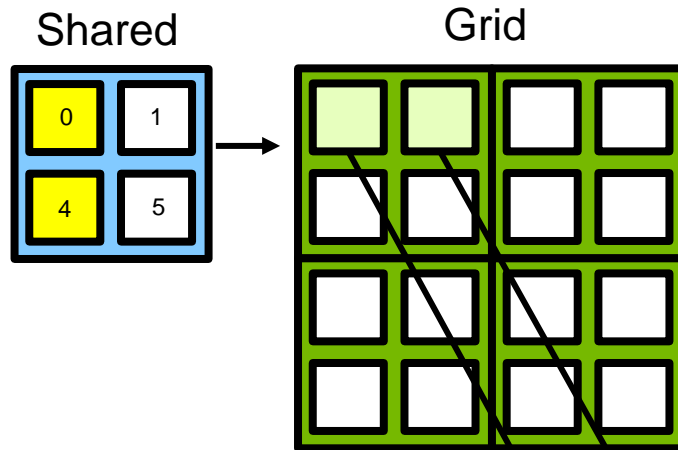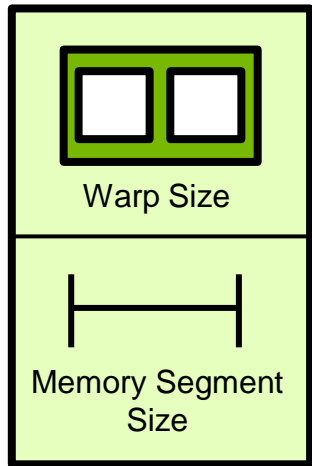
```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```

Input

Output

Warp Size

Memory Segment Size

Shared

Grid

Shared

| 0 | 1 |
| 4 | 5 |
| 8 | 9 |
| 12 | 13 |

| 2 | 3 |
| 6 | 7 |

In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory
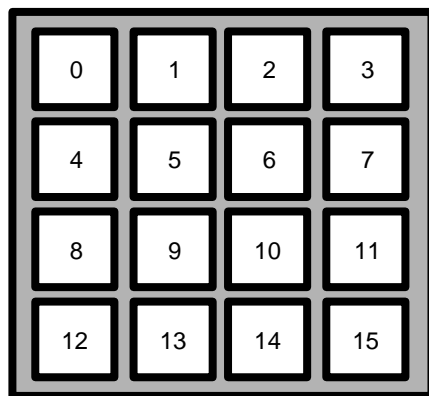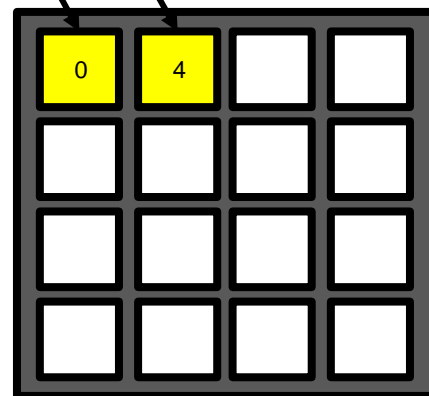
```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```
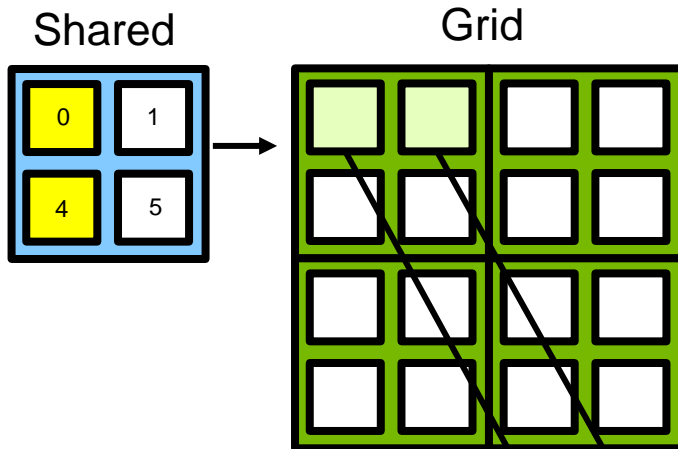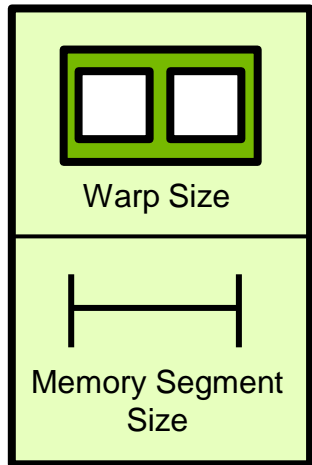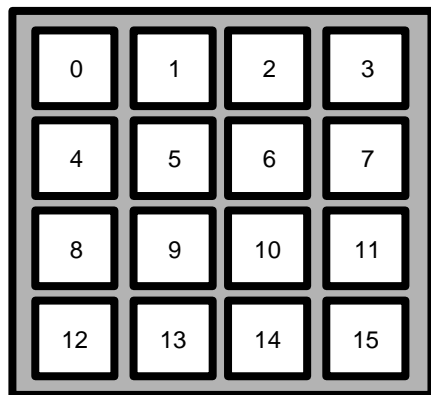
| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Input

Output

Warp Size

Memory Segment Size

Shared

| | |
|---|---|
| 0 | 1 |
| 4 | 5 |
| 8 | 9 |
| 12 | 13 |

Grid

Shared

| | |
|---|---|
| 2 | 3 |
| 6 | 7 |
| 10 | 11 |
| | |

In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```
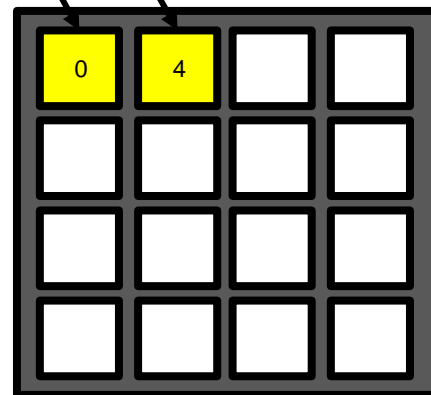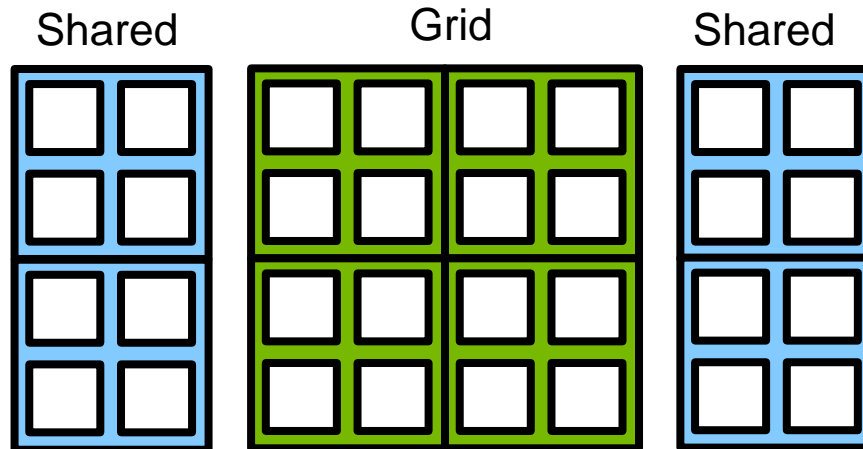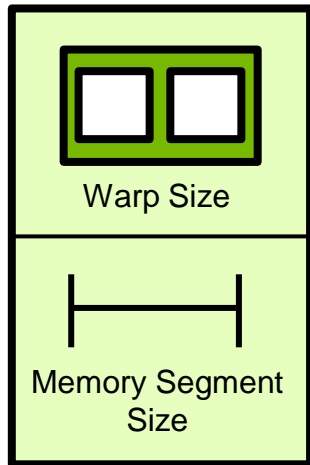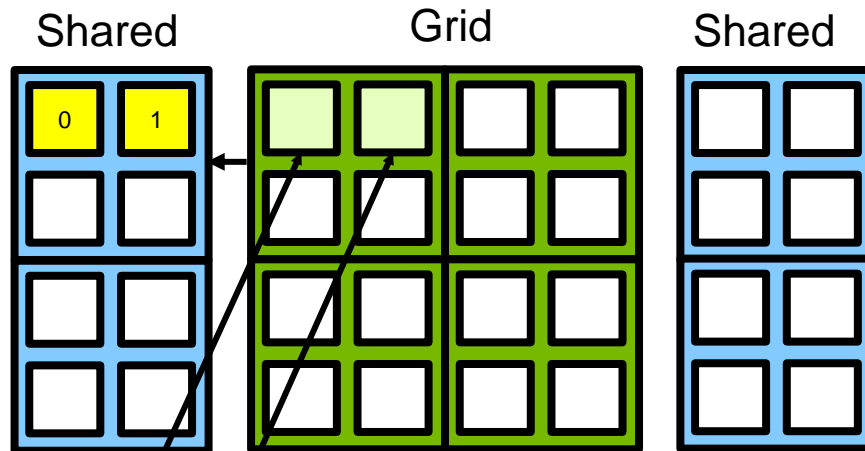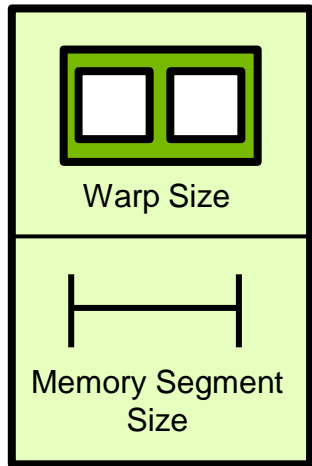
Input

| | | | |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Output

| | | | |
|---|---|---|---|
| | | | |
| | | | |
| | | | |
| | | | |

**Warp Size**

**Memory Segment Size**

**Shared**

| 0 | 1 |
| 4 | 5 |

| 8 | 9 |
| 12 | 13 |

**Grid**

**Shared**

| 2 | 3 |
| 6 | 7 |

| 10 | 11 |
| 14 | 15 |

In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory
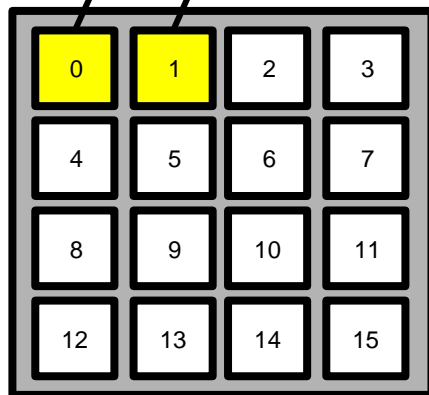
```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```
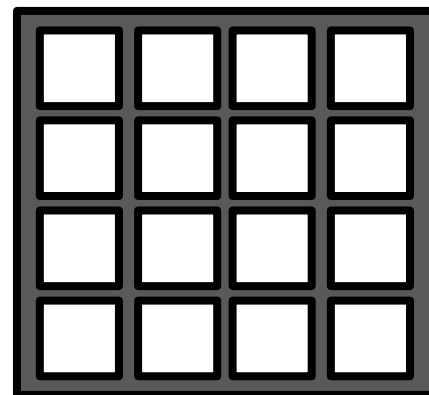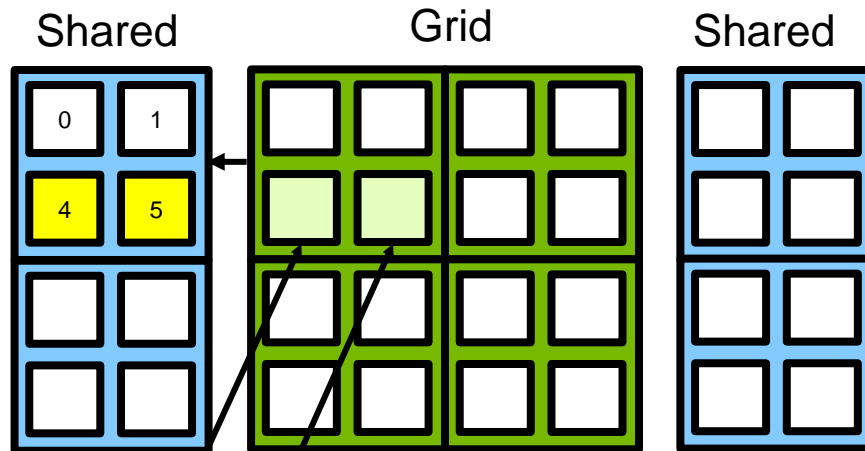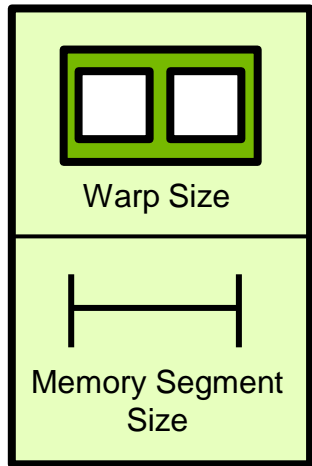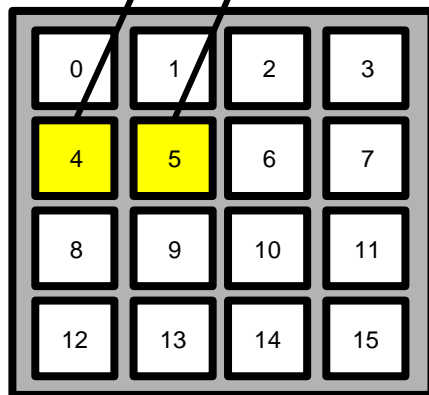
**Input**

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

**Output**

## Warp Size

## Memory Segment Size

## Shared

## Grid

## Shared

In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```
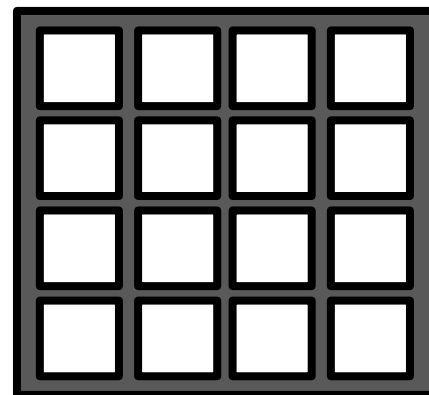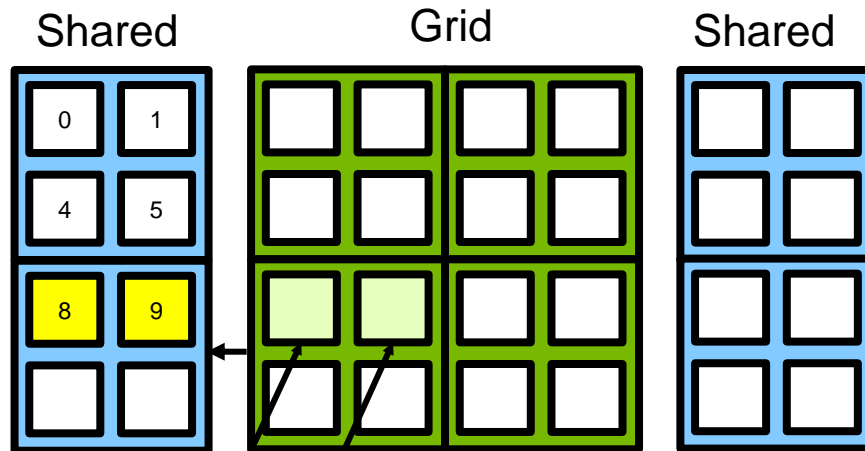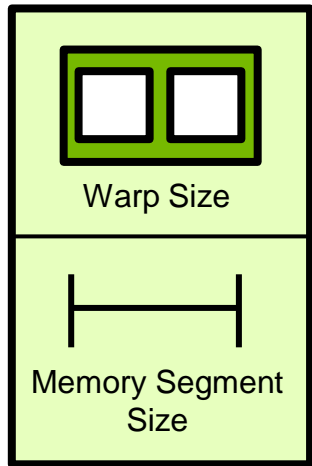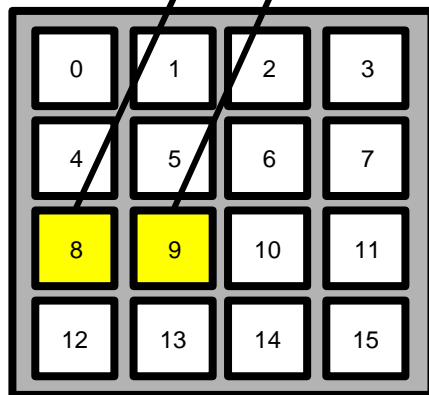
## Input

## Output

Warp Size

Memory Segment Size

Shared

Grid

Shared

In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory
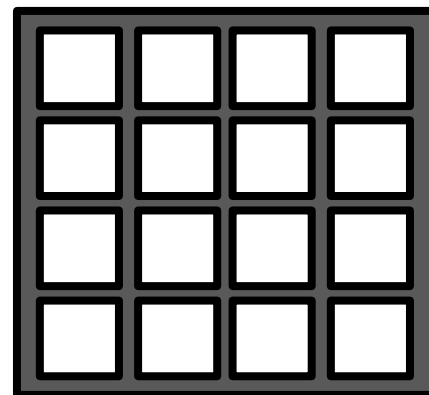
```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```
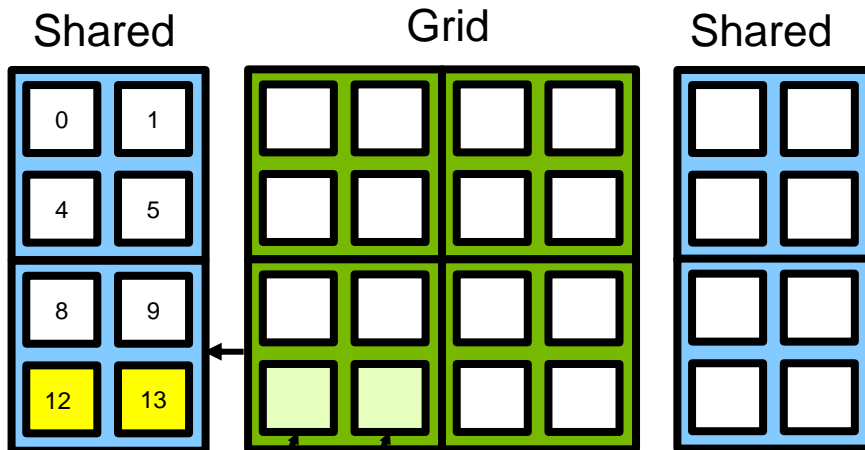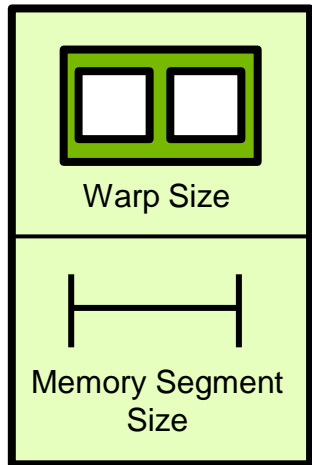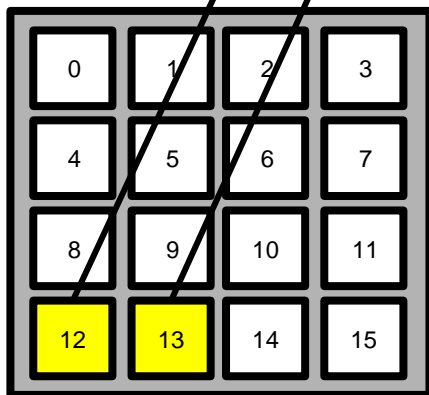
Input

Output

Warp Size

Memory Segment Size

Shared

| 0 | 1 |
| 4 | 5 |
| 8 | 9 |
| 12 | 13 |

Grid

Shared

| 2 | 3 |
| 6 | 7 |
| 10 | 11 |
| 14 | 15 |

In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```
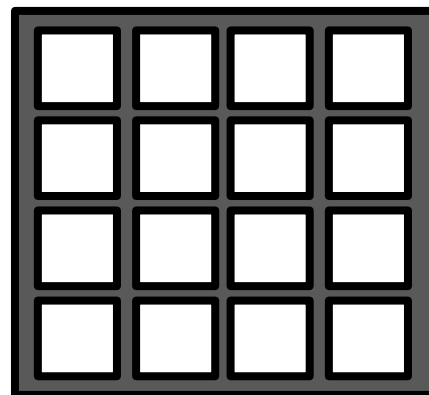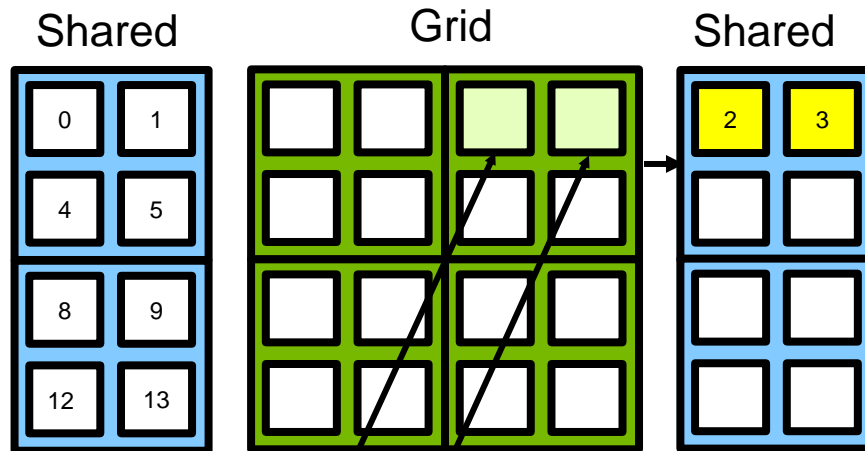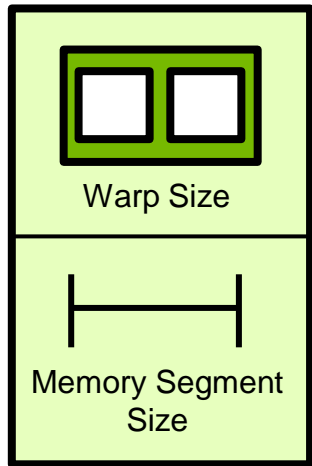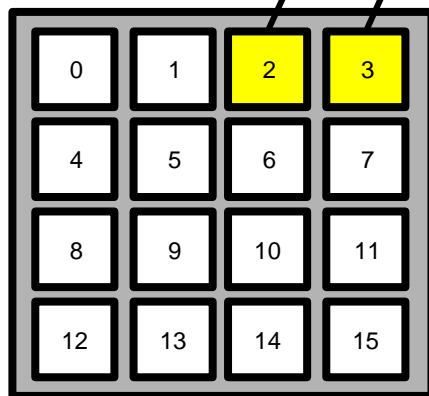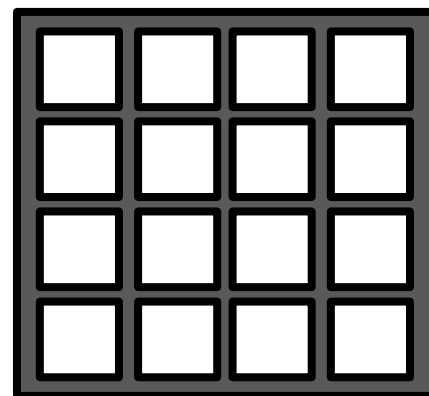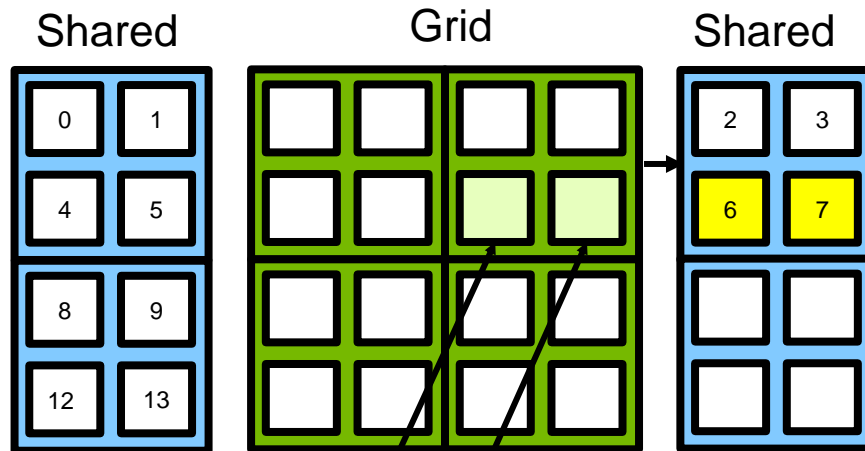
Input

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Output

| 0 | 4 | 8 | 12 |
| 1 | 5 | | |
| | | | |
| | | | |

## Warp Size

## Memory Segment Size

## Shared

| 0 | 1 |
| 4 | 5 |
| 8 | 9 |
| 12 | 13 |

## Grid

## Shared

| 2 | 3 |
| 6 | 7 |
| 10 | 11 |
| 14 | 15 |

In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```
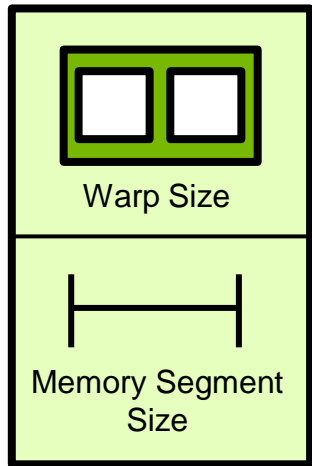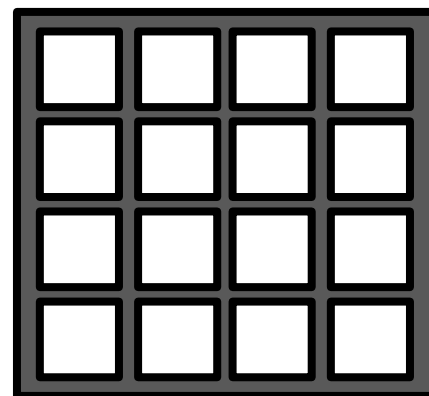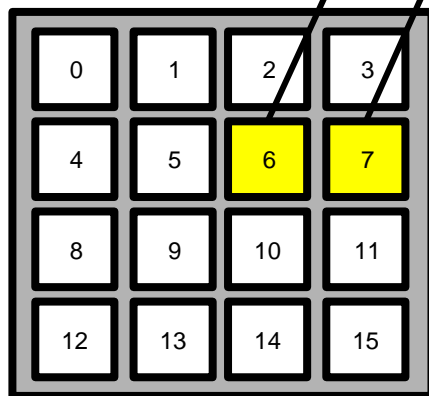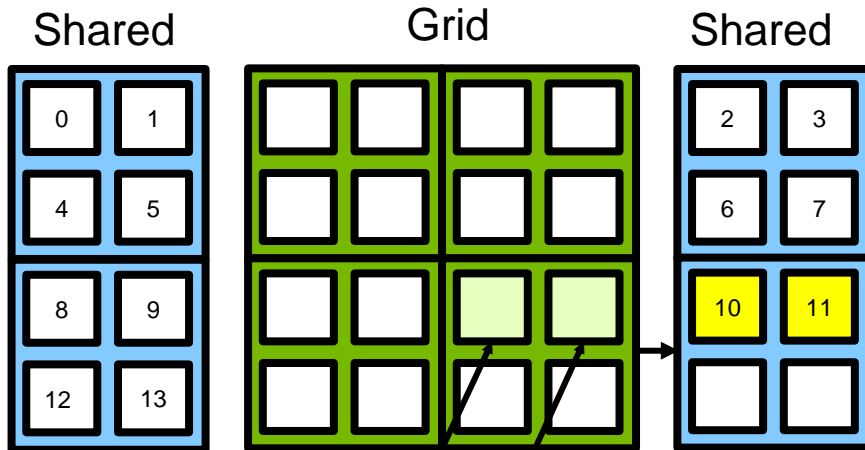
## Input

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

## Output

| 0 | 4 | 8 | 12 |
| 1 | 5 | 9 | 13 |
| | | | |
| | | | |

Warp Size

Memory Segment Size

Shared

Grid

Shared

In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```
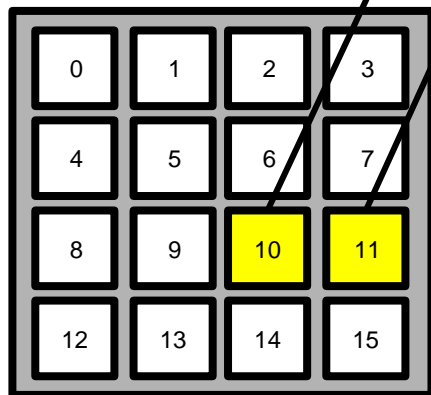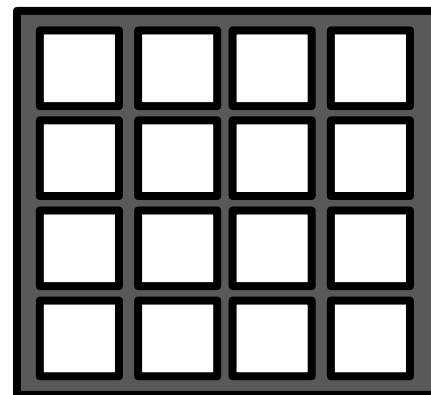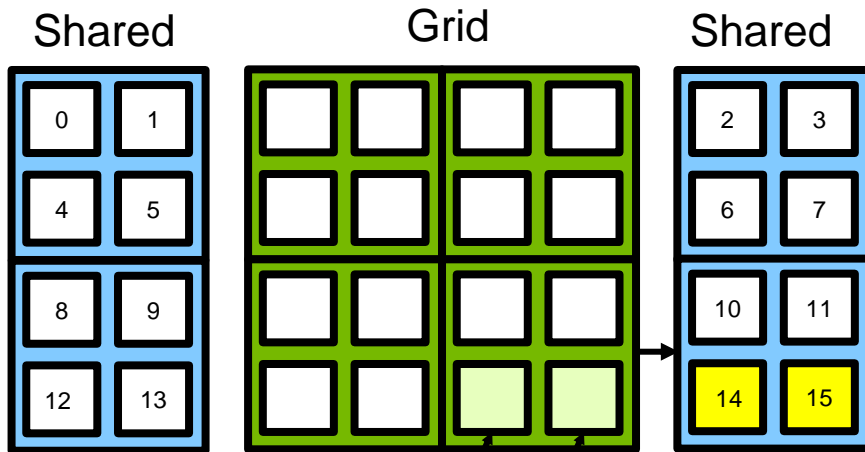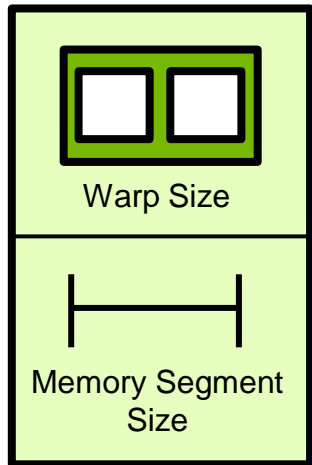
Input

Output

**Warp Size**

**Memory Segment Size**

**Shared**

| 0 | 1 |
|---|---|
| 4 | 5 |
| 8 | 9 |
| 12 | 13 |

**Grid**

**Shared**

| 2 | 3 |
|---|---|
| 6 | 7 |
| 10 | 11 |
| 14 | 15 |

In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory
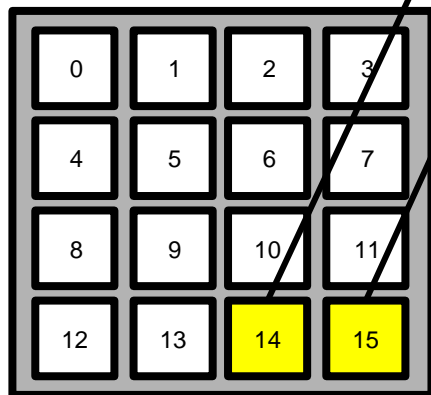
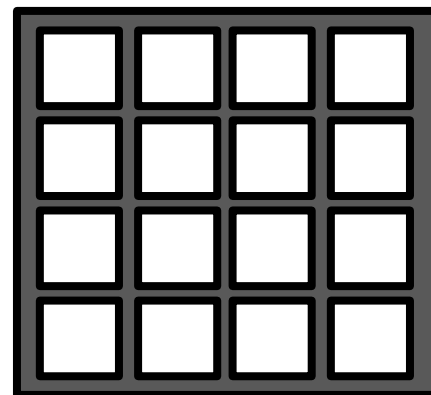```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```
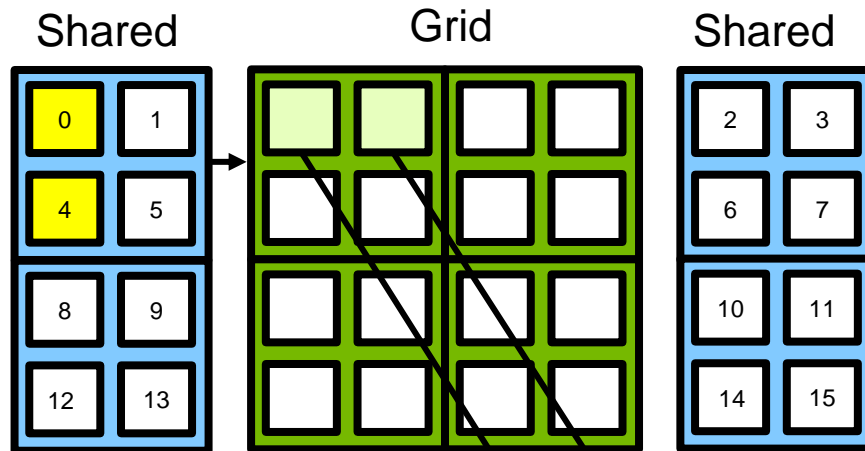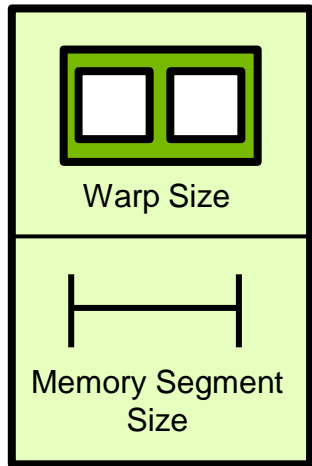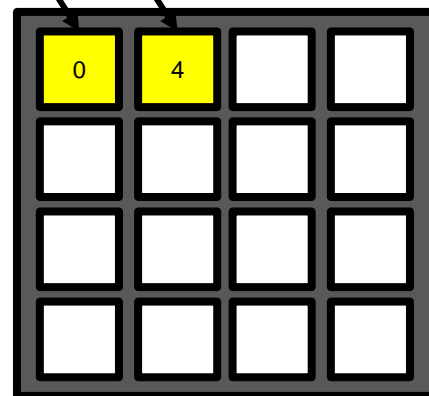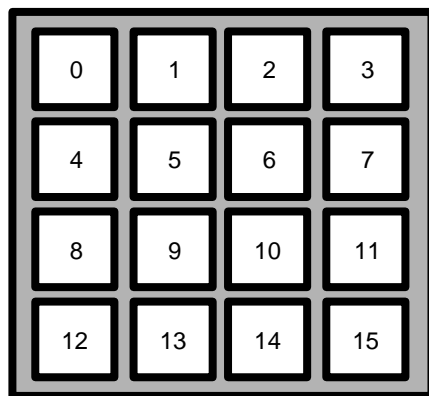
**Input**

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

**Output**

| 0 | 4 | 8 | 12 |
|---|---|---|---|
| 1 | 5 | 9 | 13 |
| 2 | 6 | | |
| 3 | 7 | | |

Warp Size

Memory Segment Size

Shared

| 0 | 1 |
| 4 | 5 |
| 8 | 9 |
| 12 | 13 |

Grid

Shared

| 2 | 3 |
| 6 | 7 |
| 10 | 11 |
| 14 | 15 |

In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```
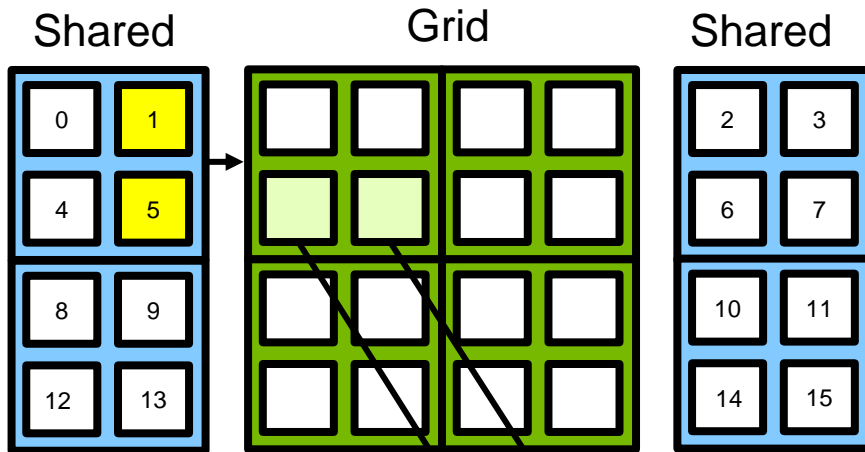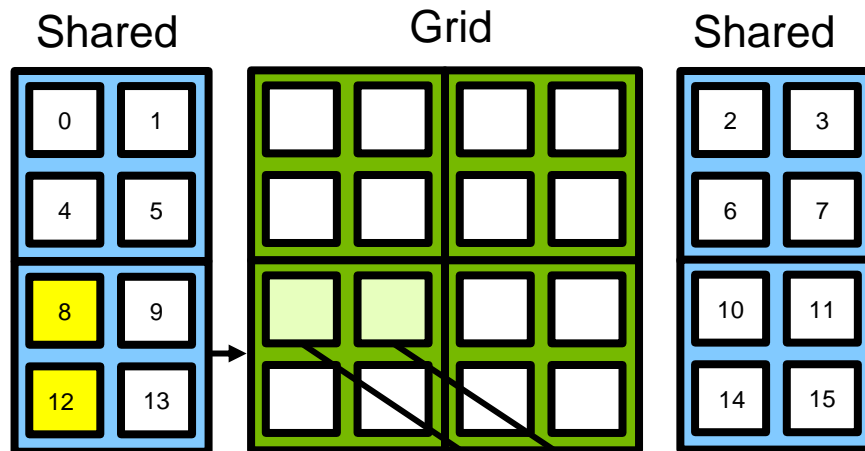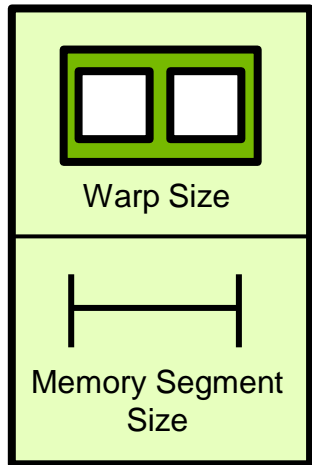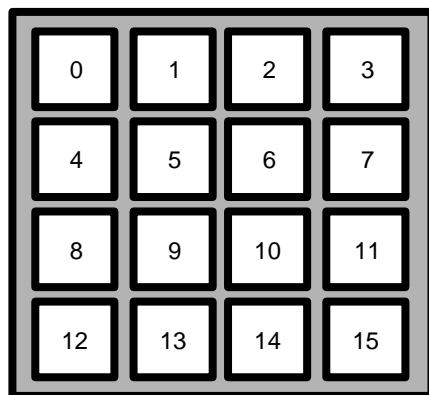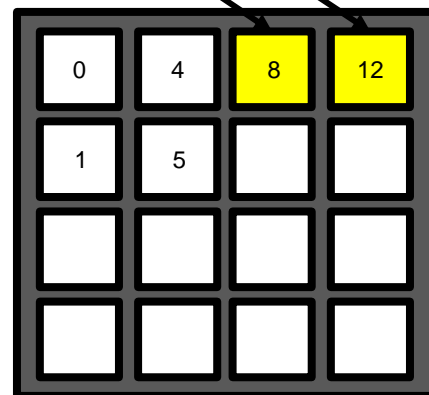
Input

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Output

| 0 | 4 | 8 | 12 |
| 1 | 5 | 9 | 13 |
| 2 | 6 | 10 | 14 |
| 3 | 7 | | |

Warp Size

Memory Segment Size

Shared

| | |
|---|---|
| 0 | 1 |
| 4 | 5 |
| 8 | 9 |
| 12 | 13 |

Grid

Shared

| | |
|---|---|
| 2 | 3 |
| 6 | 7 |
| 10 | 11 |
| 14 | 15 |

In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```
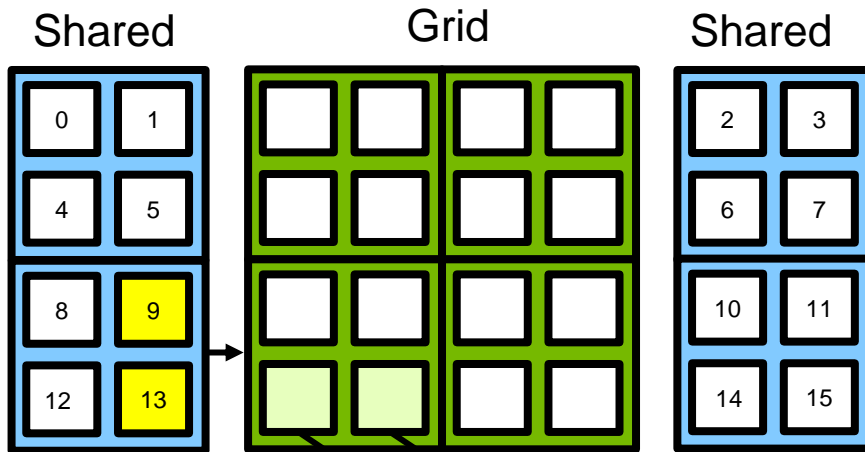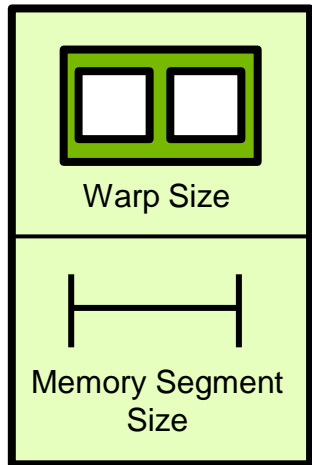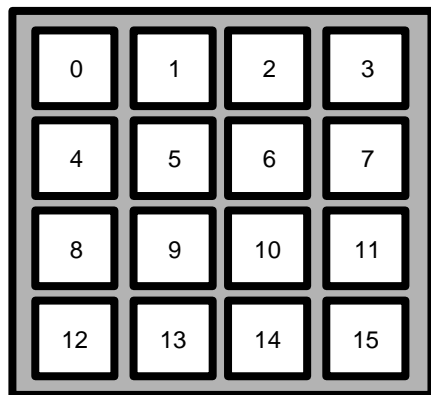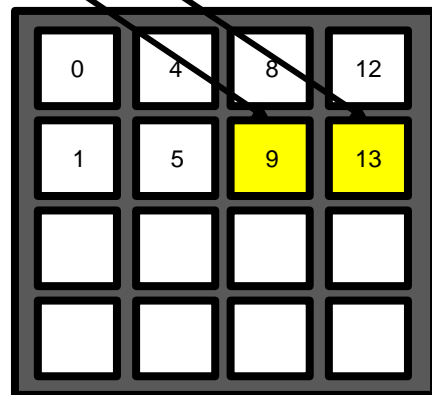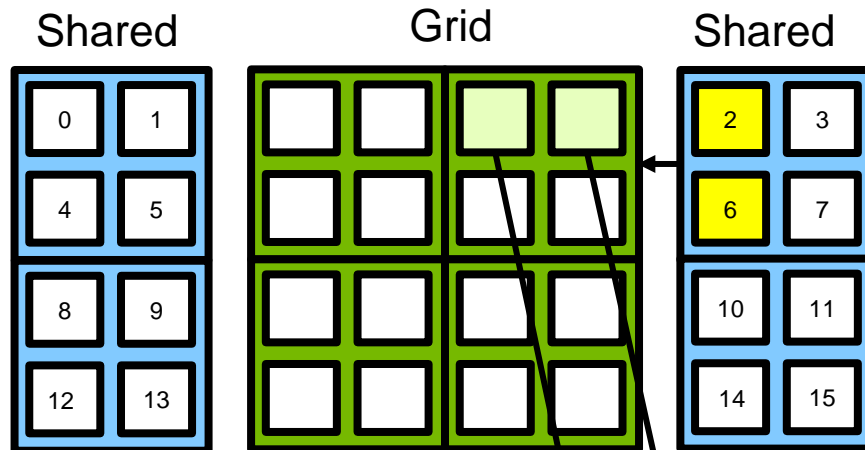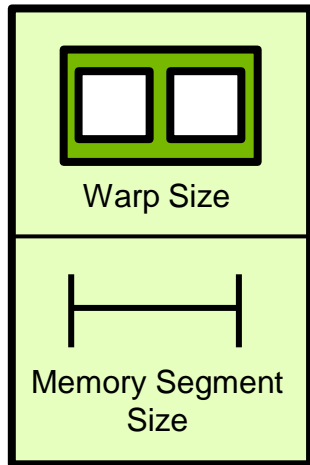
Input

| | | | |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Output

| | | | |
|---|---|---|---|
| 0 | 4 | 8 | 12 |
| 1 | 5 | 9 | 13 |
| 2 | 6 | 10 | 14 |
| 3 | 7 | 11 | 15 |

Warp Size

Memory Segment Size

Shared

| 0 | 1 |
| 4 | 5 |
| 8 | 9 |
| 12 | 13 |

Grid

Shared

| 2 | 3 |
| 6 | 7 |
| 10 | 11 |
| 14 | 15 |

In this way we can transpose the matrix while making fully coalesced reads from and writes to global memory
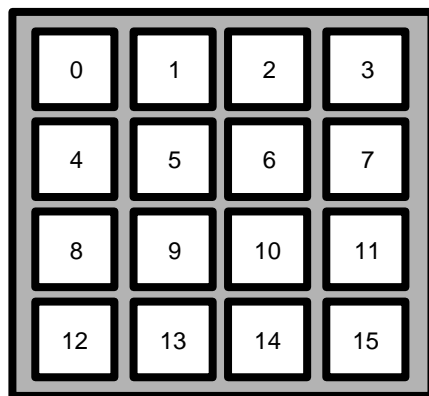
```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```
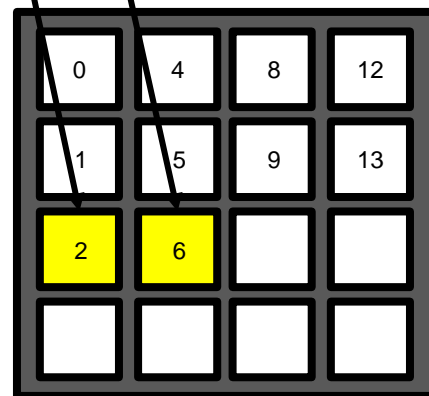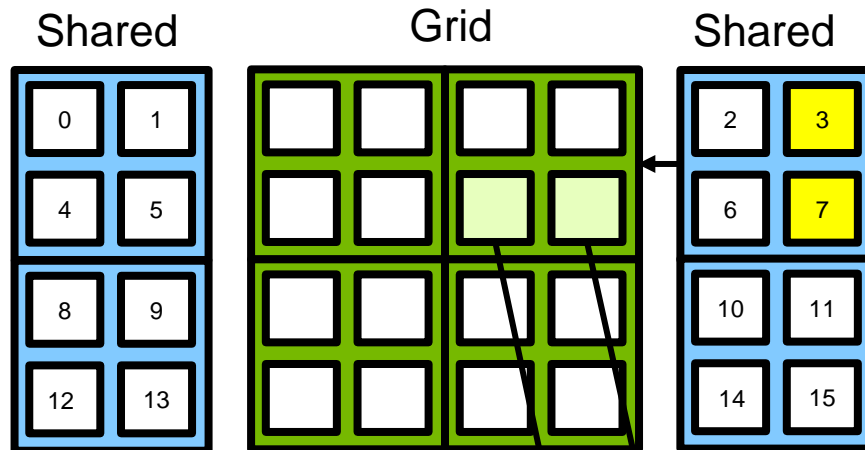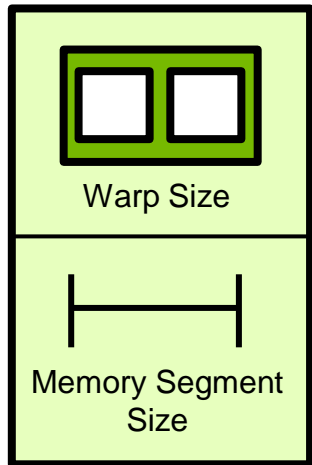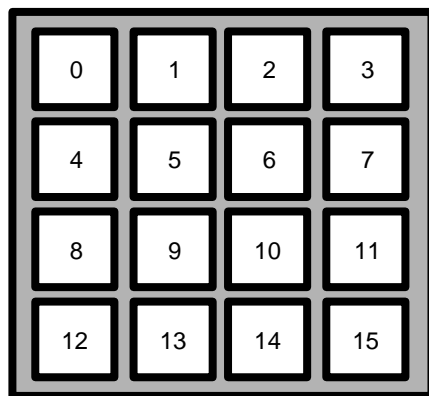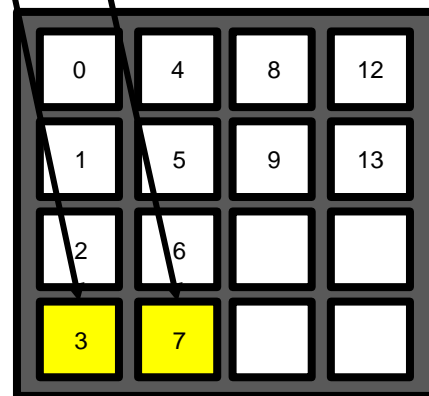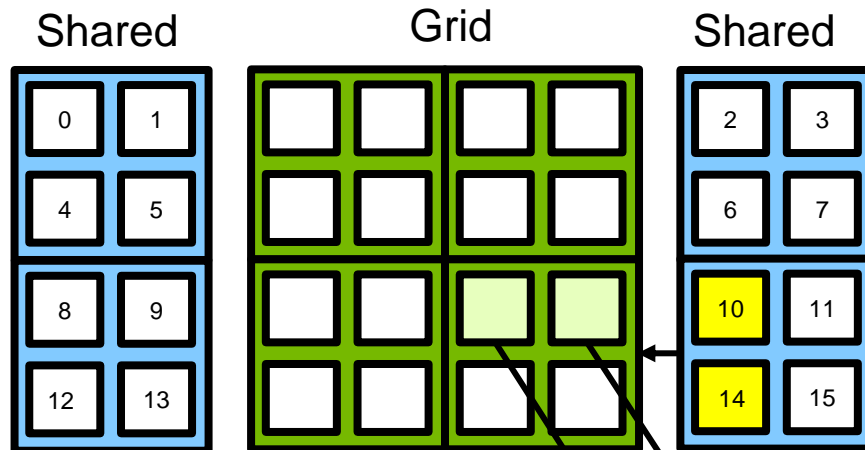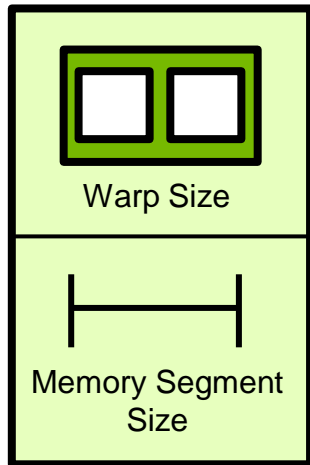
Input

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Output

| 0 | 4 | 8 | 12 |
| 1 | 5 | 9 | 13 |
| 2 | 6 | 10 | 14 |
| 3 | 7 | 11 | 15 |

# Shared Memory Bank Conflicts

Shared memory is physically stored in **banks**

Logical Shared Memory
`cuda.shared.array(4,4)`

Physical Shared Memory
in 4 banks

Warp

Successive 4-byte words (1 box in these slides) will belong to different banks

A  B  C  D

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Logical Shared Memory
`cuda.shared.array(4,4)`

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Physical Shared Memory
in 4 banks

Warp

A warp can access 4 bytes per bank, in parallel. This shared memory access would occur all at once

Logical Shared Memory
`cuda.shared.array(4,4)`

Physical Shared Memory
in 4 banks

Warp

So would this one, since each element is in a different bank

A    B    C    D

Logical Shared Memory
`cuda.shared.array(4,4)`

Physical Shared Memory
in 4 banks

# Warp



Memory accesses in the same bank result in the access operations being serialized. We call this a **bank conflict**.

(A) (B) (C) (D)

| Logical Shared Memory | | | |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Logical Shared Memory
`cuda.shared.array(4,4)`

| Physical Shared Memory | | | |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Physical Shared Memory
in 4 banks

Warp

In this scenario, we have a 2-way bank conflict that would require the memory access to be serialized over 2 cycles.

A    B    C    D

Logical Shared Memory
cuda.shared.array(4,4)

Physical Shared Memory
in 4 banks

Warp

Here have a 4-way bank conflict that would require the memory access to be serialized over 4 cycles.

A    B    C    D

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Logical Shared Memory
`cuda.shared.array(4,4)`

Physical Shared Memory
in 4 banks

Warp Size

Memory Segment Size

Shared

| 0 | 1 |
| 4 | 5 |
| 8 | 9 |
| 12 | 13 |

Grid

Shared

| 2 | 3 |
| 6 | 7 |
| 10 | 11 |
| 14 | 15 |

Recall from our earlier matrix transpose example that we were making this very kind of columnar read from shared memory, which means we had significant bank conflicts

```
tile = cuda.shared.array(2,2)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```

Input

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Output

| 0 | 4 | 8 | 12 |
| 1 | 5 | | |
| | | | |
| | | | |

Here is a technique we can use to avoid bank conflicts when we know we need to make columnar access to shared memory

# Warp



First, when we allocate our shared memory tile, we will pad it with an extra column

## Logical Shared Memory
`cuda.shared.array(4,5)`

# Warp



Next, when we write to the tile, we act as if the tile is (4,4) and only write to addresses in the range [0:4][0:4]

| | | | | |
|---|---|---|---|---|
| 0 | 1 | 2 | 3 | |
| 4 | 5 | 6 | 7 | |
| 8 | 9 | 10 | 11 | |
| 12 | 13 | 14 | 15 | |

Logical Shared Memory
`cuda.shared.array(4,5)`

Warp

The physical shared memory has a fixed size of 32 banks (4 banks in our slides to save space), so our padding of the shared memory array does not affect the number of memory banks

A   B   C   D

| 0 | 1 | 2 | 3 | |
| 4 | 5 | 6 | 7 | |
| 8 | 9 | 10 | 11 | |
| 12 | 13 | 14 | 15 | |

Logical Shared Memory
`cuda.shared.array(4,5)`

Physical Shared Memory
in 4 banks

Warp

A    B    C    D

| 0 | 1 | 2 | 3 |
| --- | --- | --- | --- |
|   |   |   |   |
|   |   |   |   |
|   |   |   |   |
|   |   |   |   |

| 0 | 1 | 2 | 3 |   |
| --- | --- | --- | --- | --- |
| 4 | 5 | 6 | 7 |   |
| 8 | 9 | 10 | 11 |   |
| 12 | 13 | 14 | 15 |   |

Logical Shared Memory
`cuda.shared.array(4,5)`

Physical Shared Memory
in 4 banks

# Warp

A  B  C  D

Physical Shared Memory

| A | B | C | D |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
|   | 4 | 5 | 6 |
| 7 |   |   |   |

Logical Shared Memory
`cuda.shared.array(4,5)`

| 0 | 1 | 2 | 3 | |
|---|---|---|---|---|
| 4 | 5 | 6 | 7 | |
| 8 | 9 | 10 | 11 | |
| 12 | 13 | 14 | 15 | |

Physical Shared Memory
in 4 banks

# Warp

So if we consider how the array is laid out within the memory banks, we see the following:

**A** **B** **C** **D**

| A | B | C | D |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
|   | 4 | 5 | 6 |
| 7 |   | 8 | 9 |
| 10 | 11 |   |   |
|   |   |   |   |

## Physical Shared Memory
in 4 banks

| 0 | 1 | 2 | 3 |   |
|---|---|---|---|---|
| 4 | 5 | 6 | 7 |   |
| 8 | 9 | 10 | 11 |   |
| 12 | 13 | 14 | 15 |   |

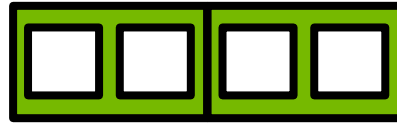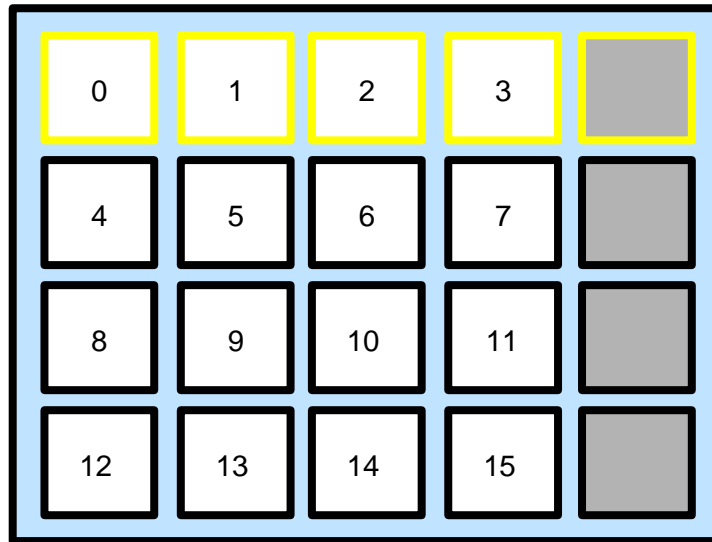## Logical Shared Memory
`cuda.shared.array(4,5)`

# Warp

So if we consider how the array is laid out within the memory banks, we see the following:

A  B  C  D

Physical Shared Memory
in 4 banks

| A | B | C | D |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
|   | 4 | 5 | 6 |
| 7 |   | 8 | 9 |
| 10 | 11 |   | 12 |
| 13 | 14 | 15 |   |

| 0 | 1 | 2 | 3 | |
|---|---|---|---|---|
| 4 | 5 | 6 | 7 | |
| 8 | 9 | 10 | 11 | |
| 12 | 13 | 14 | 15 | |

Logical Shared Memory
cuda.shared.array(4,5)

Warp

Now when we access a column of shared memory, each element resides in a different bank and there are no bank conflicts



Logical Shared Memory
`cuda.shared.array(4,5)`

Physical Shared Memory
in 4 banks

Warp

Now when we access a column of shared memory, each element resides in a different bank and there are no bank conflicts

A  B  C  D

| A | B | C | D |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
|   | 4 | 5 | 6 |
| 7 |   | 8 | 9 |
| 10 | 11 |  | 12 |
| 13 | 14 | 15 |  |

Physical Shared Memory
in 4 banks

| | | | | |
|---|---|---|---|---|
| 0 | 1 | 2 | 3 | |
| 4 | 5 | 6 | 7 | |
| 8 | 9 | 10 | 11 | |
| 12 | 13 | 14 | 15 | |

Logical Shared Memory
`cuda.shared.array(4,5)`

Warp

Now when we access a column of shared memory, each element resides in a different bank and there are no bank conflicts

Logical Shared Memory
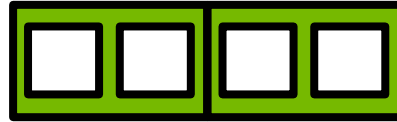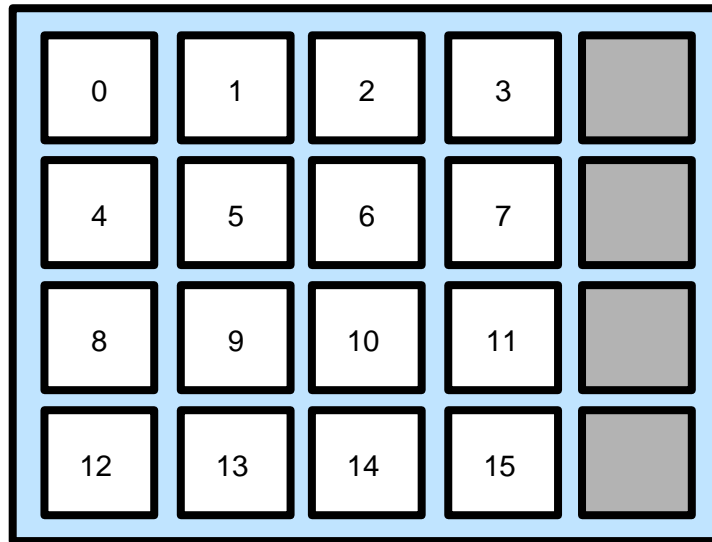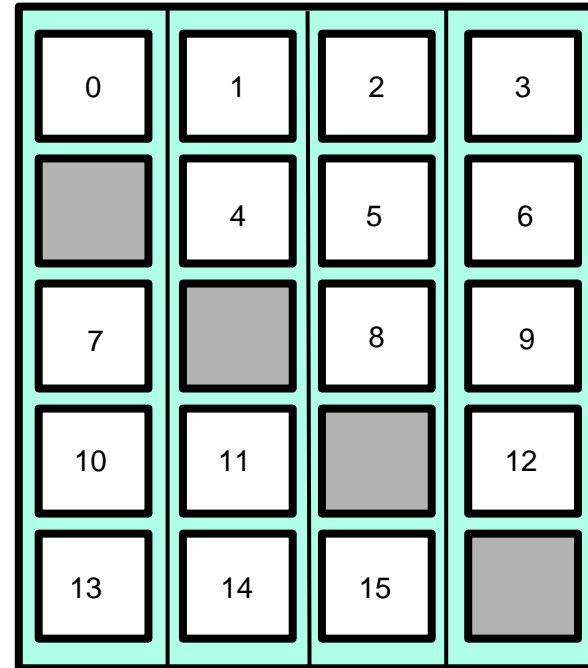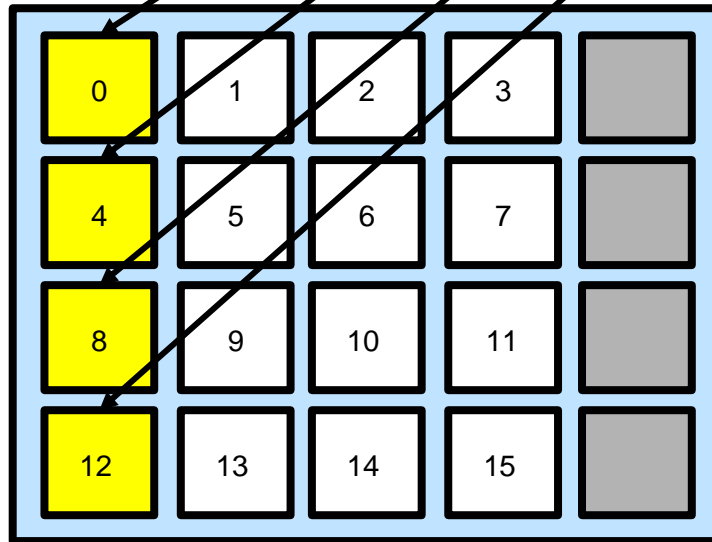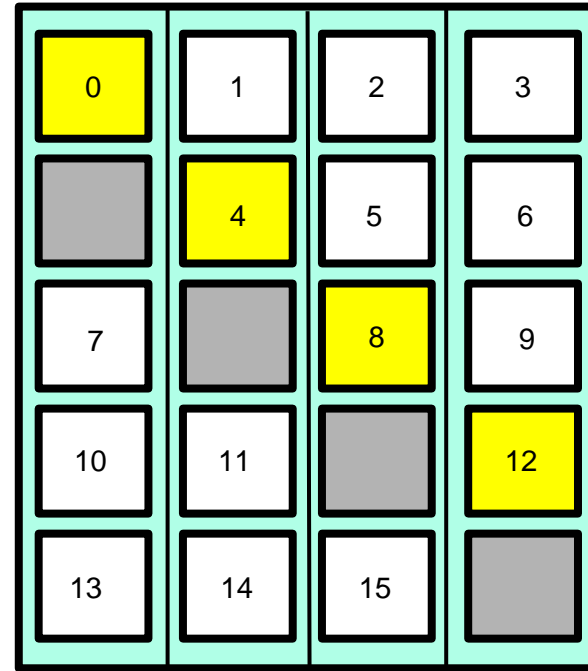`cuda.shared.array(4,5)`

Physical Shared Memory
in 4 banks

Warp

Now when we access a column of shared memory, each element resides in a different bank and there are no bank conflicts
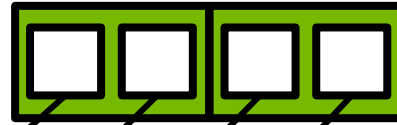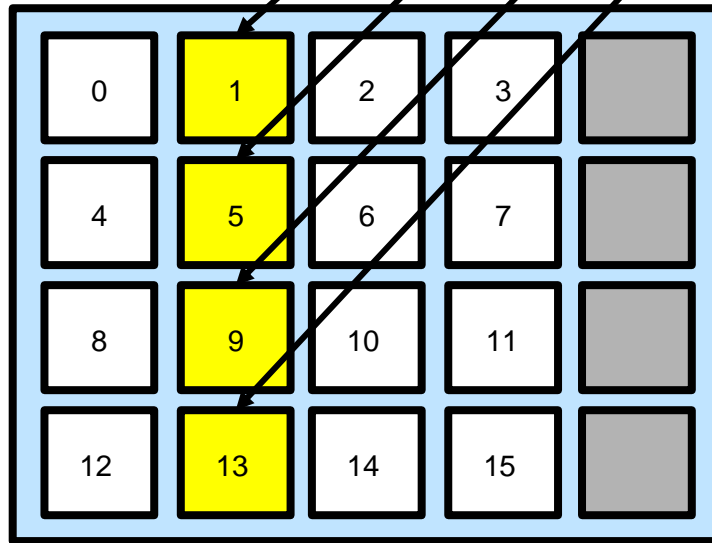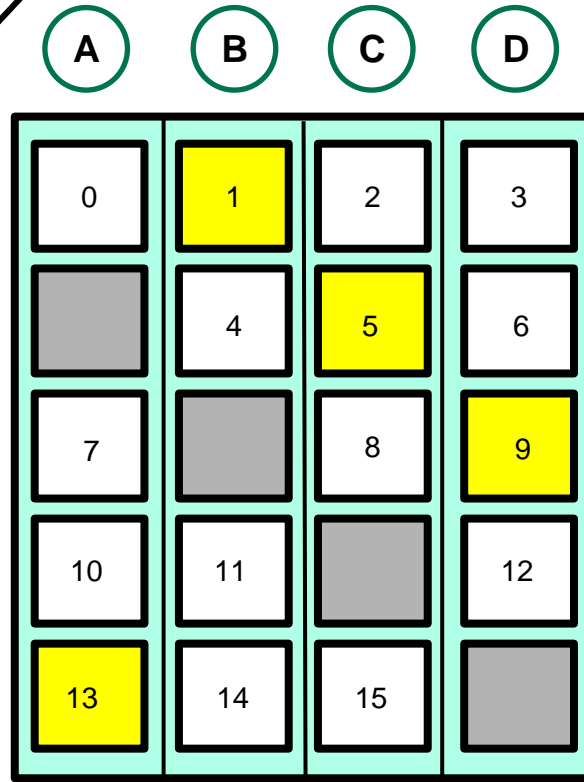
A  B  C  D

| | | | |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
| | 4 | 5 | 6 |
| 7 | | 8 | 9 |
| 10 | 11 | | 12 |
| 13 | 14 | 15 | |

Physical Shared Memory
in 4 banks

| 0 | 1 | 2 | 3 | |
|---|---|---|---|---|
| 4 | 5 | 6 | 7 | |
| 8 | 9 | 10 | 11 | |
| 12 | 13 | 14 | 15 | |

Logical Shared Memory
`cuda.shared.array(4,5)`

Warp

Now when we access a column of shared memory, each element resides in a different bank and there are no bank conflicts

A B C D

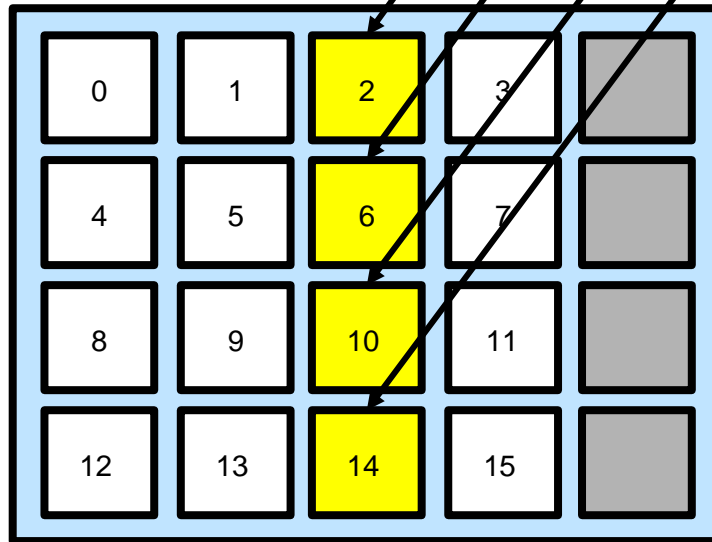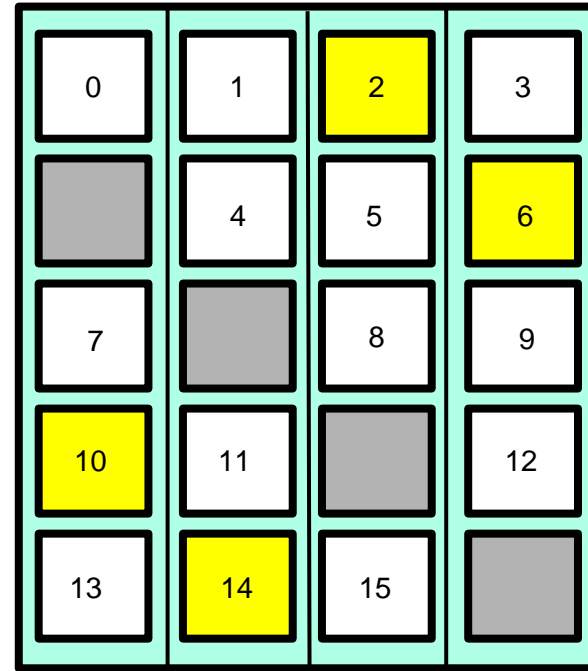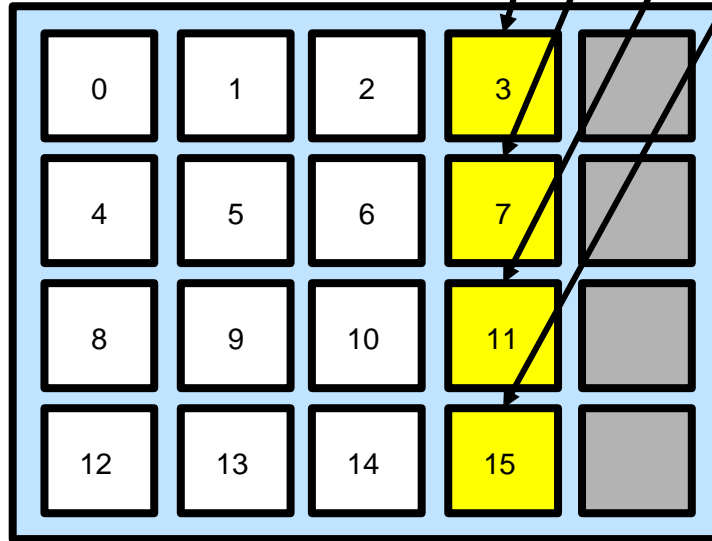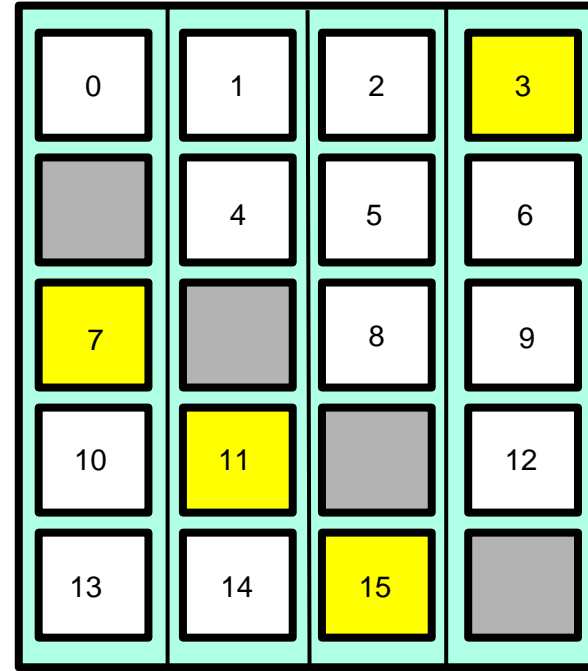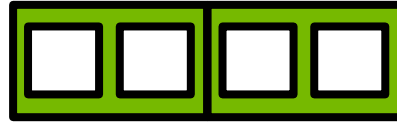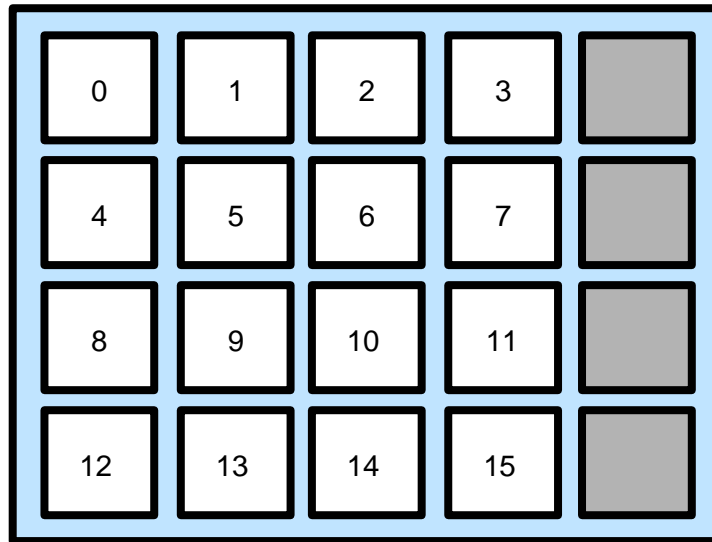Logical Shared Memory
`cuda.shared.array(4,5)`

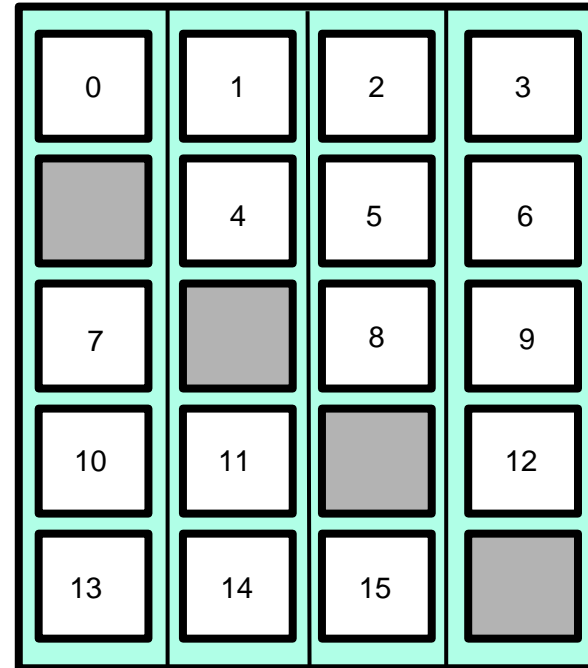Physical Shared Memory in 4 banks

# Warp



Worth mentioning that to use this technique for this example, the only change we had to make to our code was add one extra column to our shared memory allocation
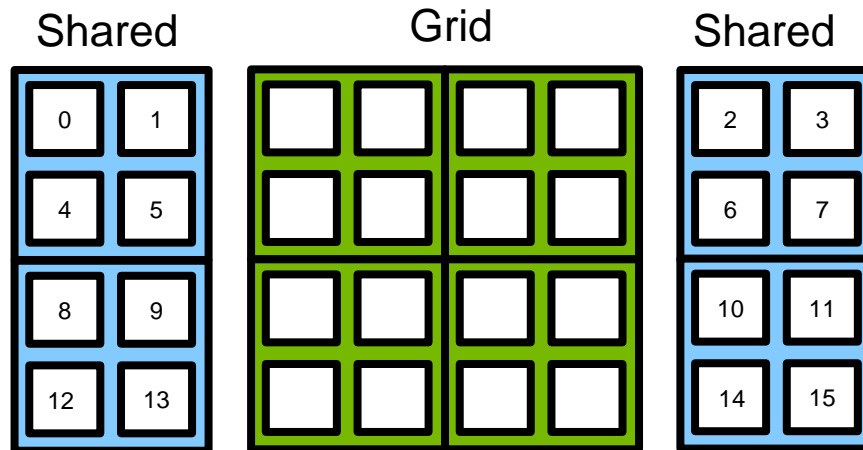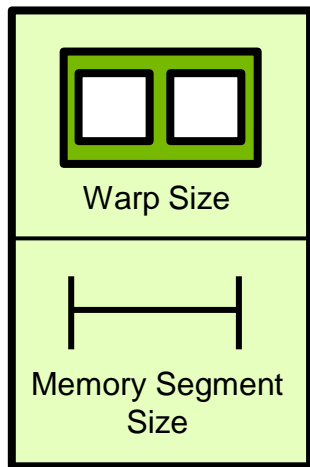
**A** **B** **C** **D**

## Logical Shared Memory
`cuda.shared.array(4,5)`

## Physical Shared Memory
in 4 banks

**Warp Size**

**Memory Segment Size**

**Shared**

| | |
|---|---|
| 0 | 1 |
| 4 | 5 |
| 8 | 9 |
| 12 | 13 |

**Grid**

**Shared**

| | |
|---|---|
| 2 | 3 |
| 6 | 7 |
| 10 | 11 |
| 14 | 15 |

From our earlier matrix transpose example, the single change in green below would suffice to avoid bank conflicts while retaining correctness

```
tile = cuda.shared.array(2,3)
x, y = cuda.grid(2)

tile[tIdx.y][tIdx.x] = in[y][x]
cuda.syncthreads()

o_x = bId.y*bDim.y + tId.x
o_y = bId.x*bDim.x + tId.y
o[o_y][o_x] = tile[tIdx.x][tIdx.y]
```

**Input**

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

**Output**

| 0 | 4 | 8 | 12 |
|---|---|---|---|
| 1 | 5 | 9 | 13 |
| 2 | 6 | 10 | 14 |
| 3 | 7 | 11 | 15 |