# Memory Coalescing

Recall that thread blocks are divided into **warps** of 32 threads
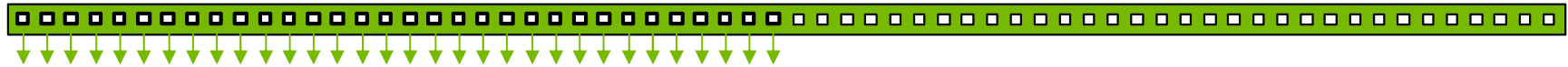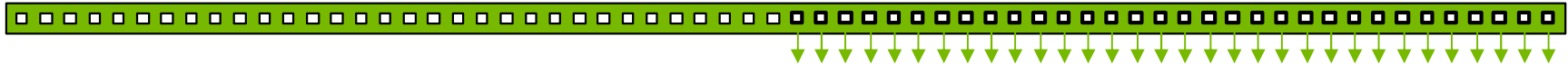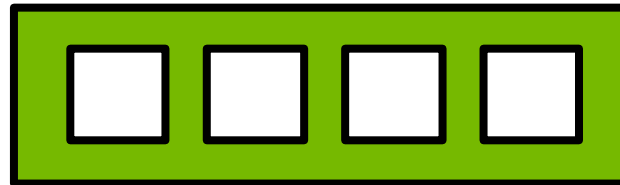
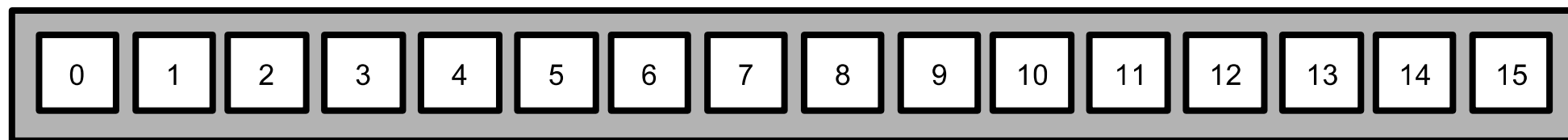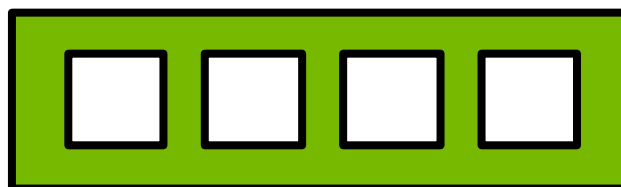Instructions are issued in parallel at the warp level of 32 threads

For space on these slides, we will treat just 4 threads as a warp

Warp

Data is transferred to and from global
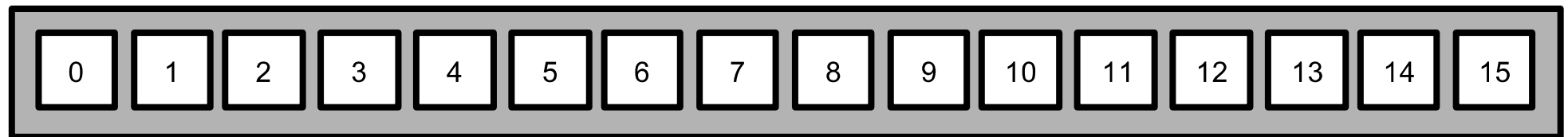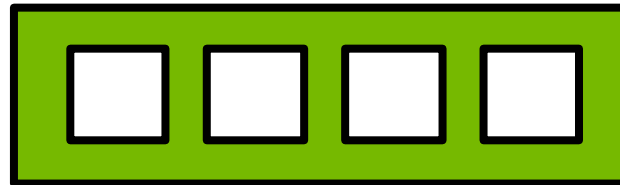device memory in 32-byte segments*

Warp

Data

(* If the data is in the L1 cache it will be transferred in 128-byte cache lines – see the notebook for details)
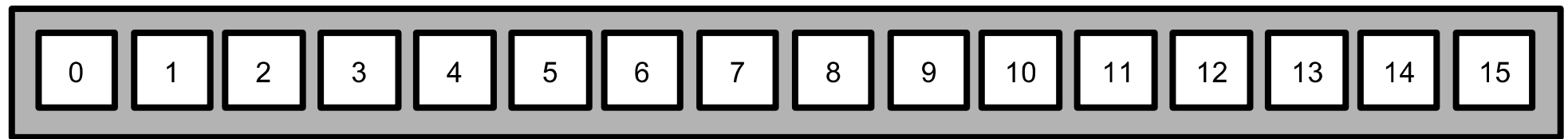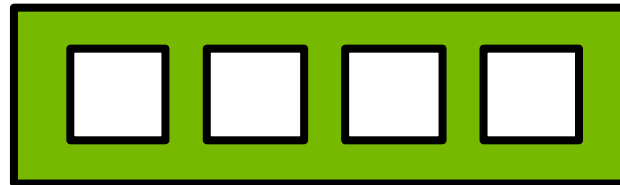
Warp

Data

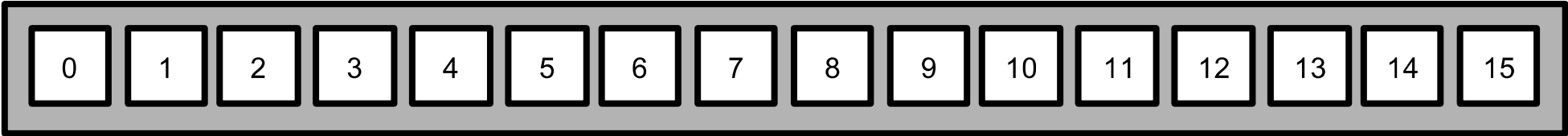For these slides we will treat 4 data elements as one of these fixed-length lines of contiguous memory

Warp

Data

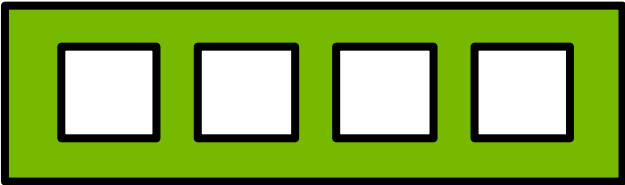| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |

Warp



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

Data
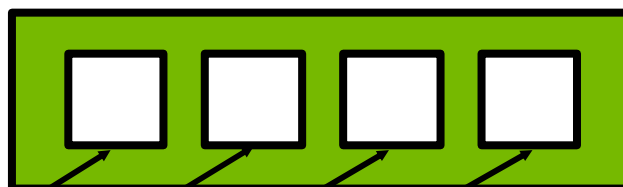
Warp



If the addresses requested are contiguous

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

a[idx] += 1
```

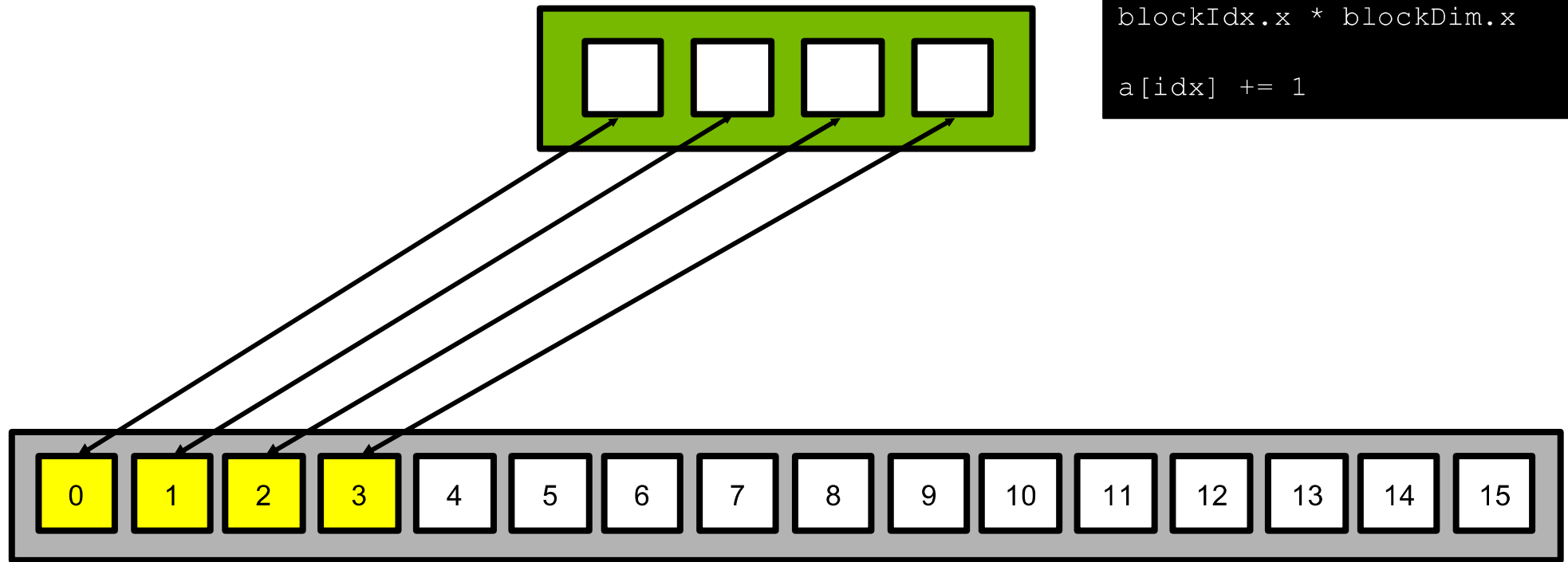| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |

Data

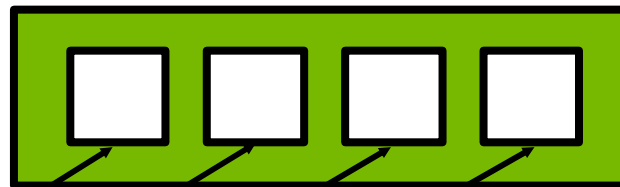All data in the line will be used

Warp

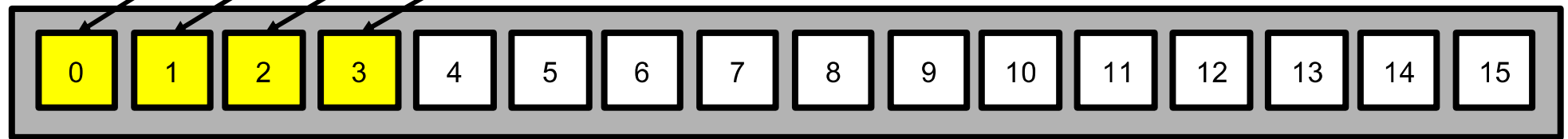```
idx = threadIdx.x +
blockIdx.x * blockDim.x

a[idx] += 1
```

Data

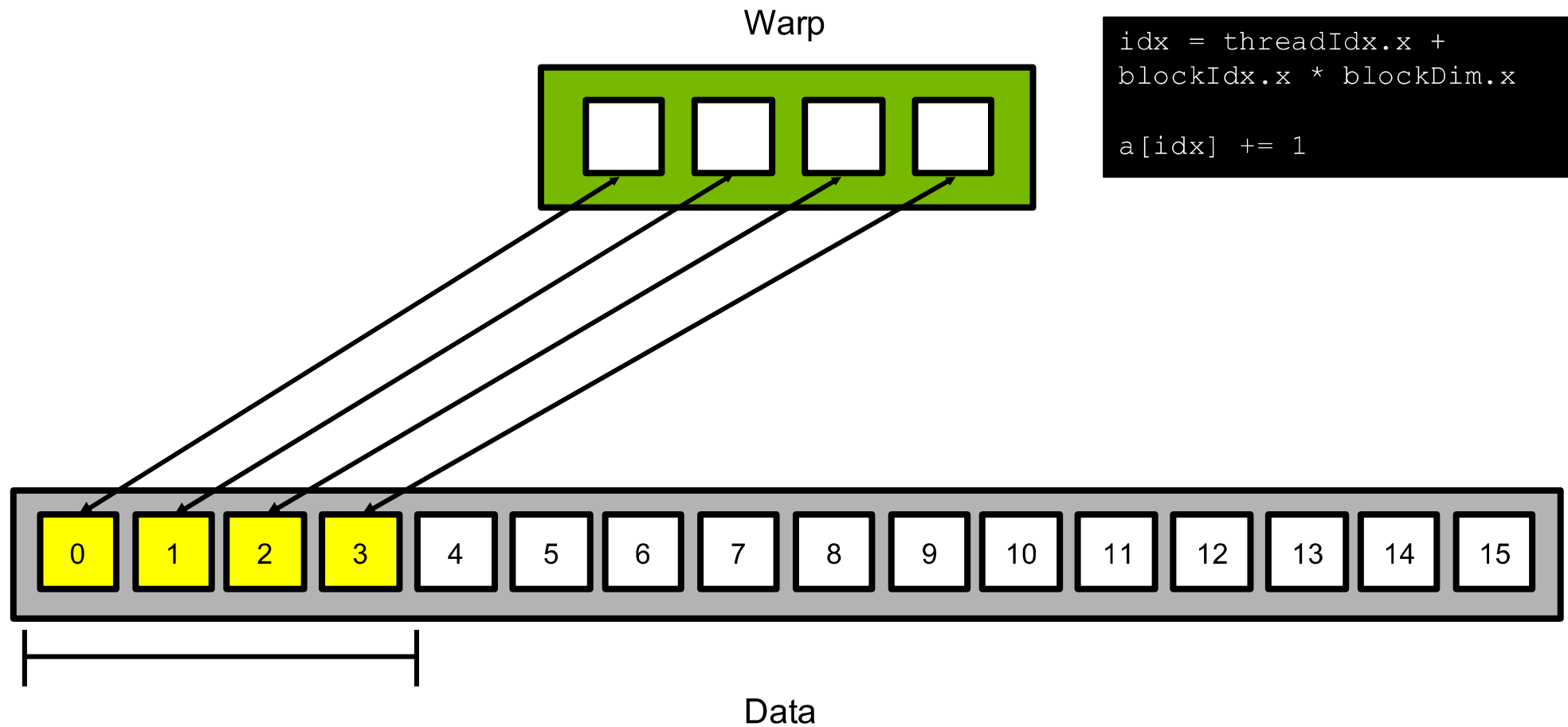And the transfer will happen in as few lines as possible

Warp

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

a[idx] += 1
```
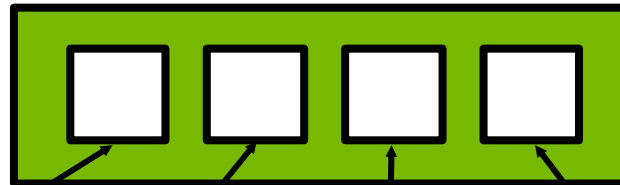
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |

Data

Warp

When this occurs, the memory access
is fully **coalesced**

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

a[idx] += 1
```

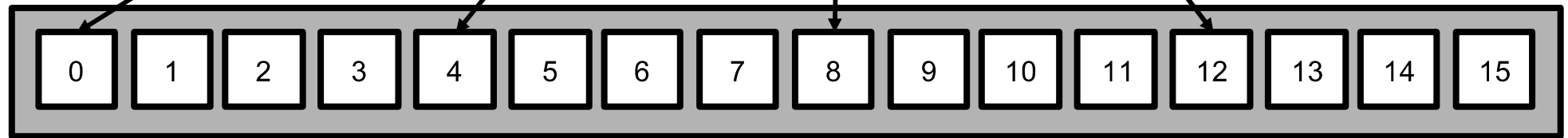| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |

Data

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

Warp

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

a[idx] += 1
```
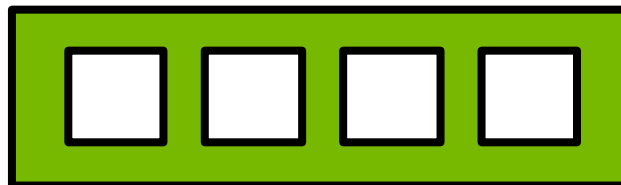
Data

# Row and Column Sum Comparison

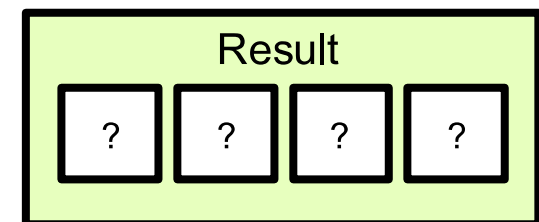Consider a kernel that stores the sum of each row of a matrix (which here is 4 contiguous data elements) in a result vector
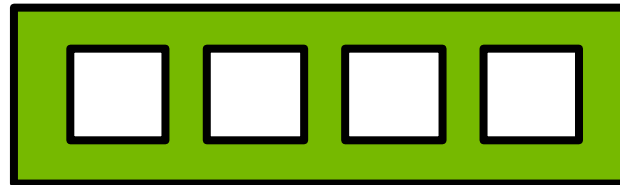
Warp

Result

? ? ? ?

Data

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

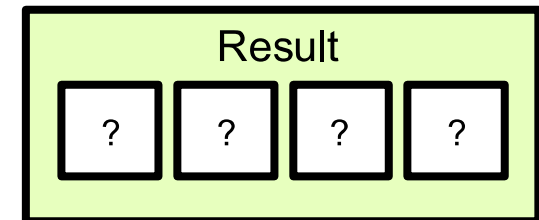A single thread could iterate over a row, summing it, and then write the result in the solution vector
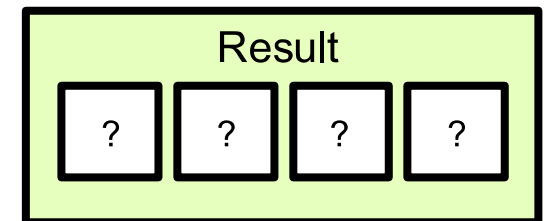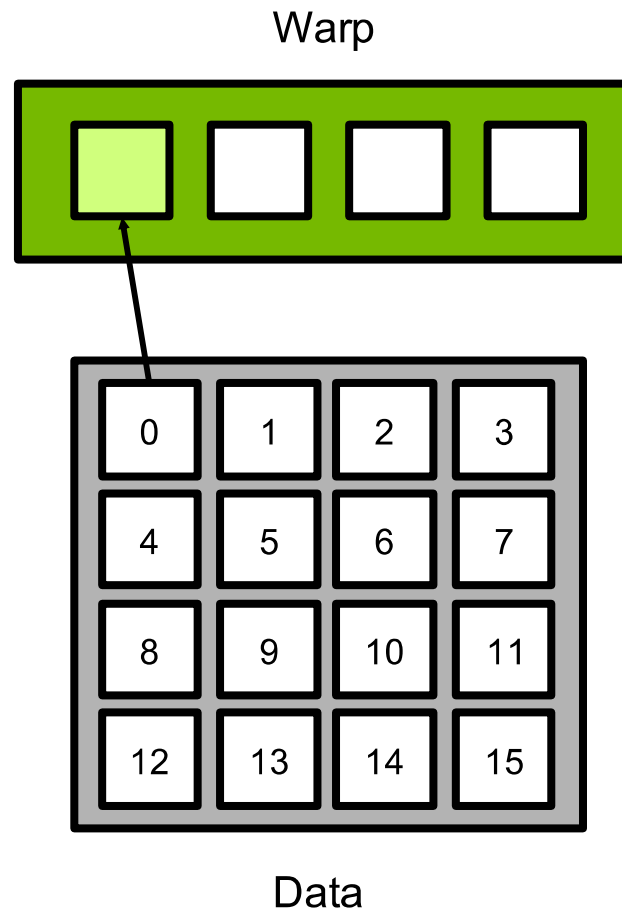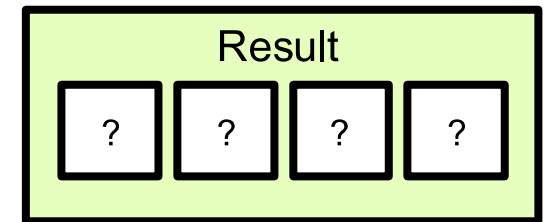
Warp

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

Result

| ? | ? | ? | ? |
|---|---|---|---|

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

Data

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

Warp

Result

? ? ? ?

Sum = 0

Data

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

**Warp**

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

**Data**

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

**Result**

| ? | ? | ? | ? |
|---|---|---|---|

Sum = 1

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

Warp

Result

? ? ? ?

Sum = 3

Data

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

Warp

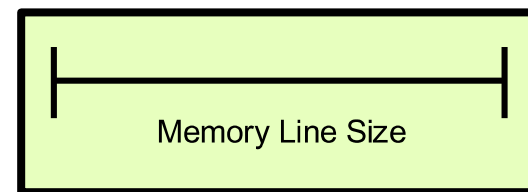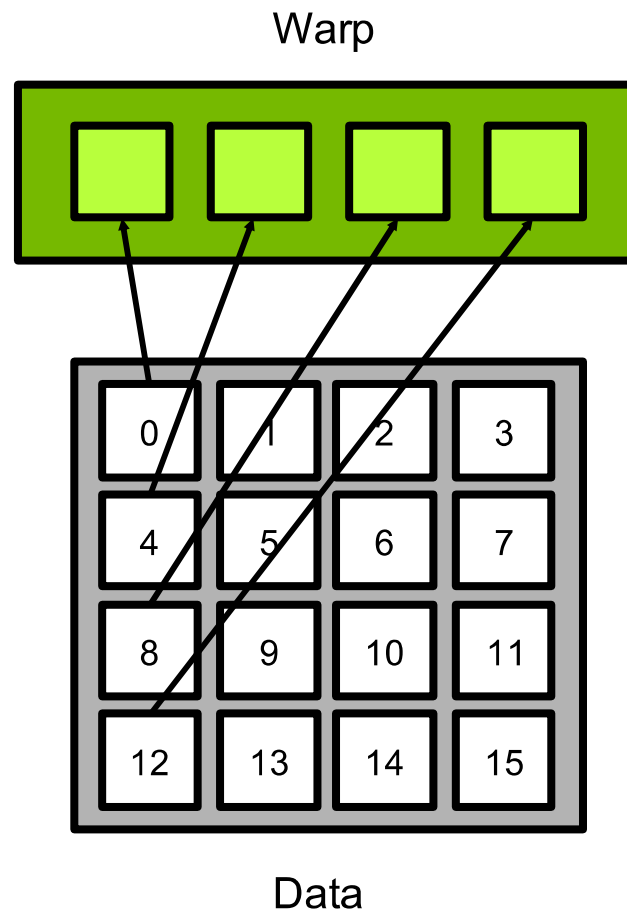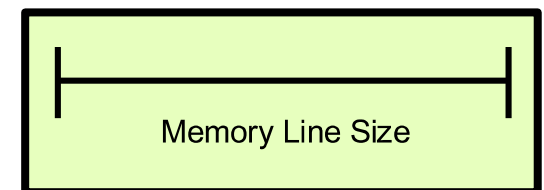Result

? ? ? ?

Sum = 6

Data

Warp



Data

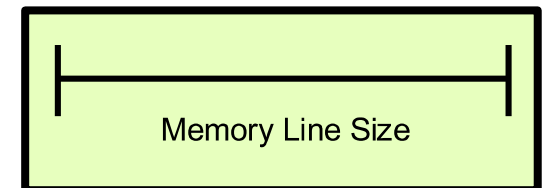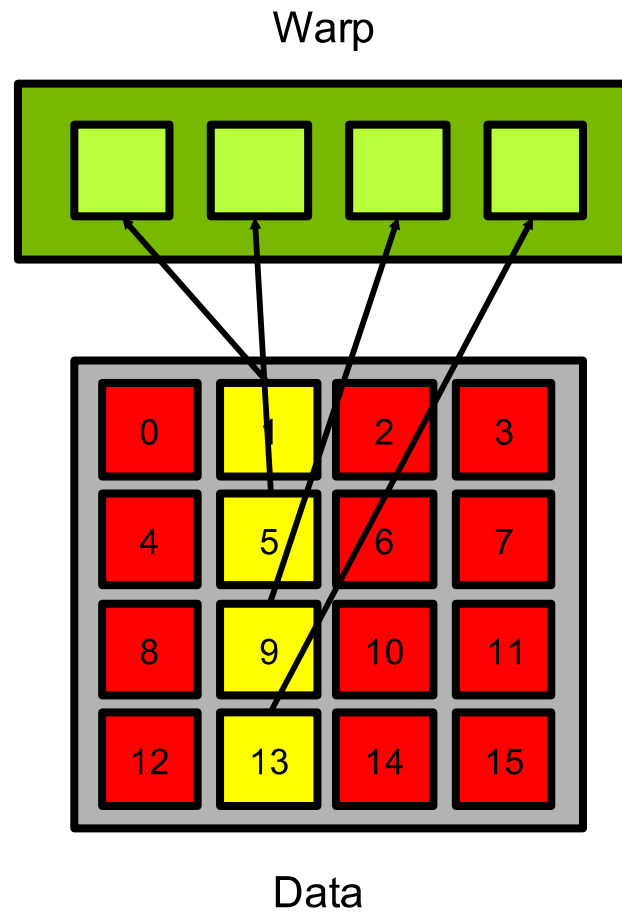Each thread in the warp is requesting data in a different line of memory

Warp

Data

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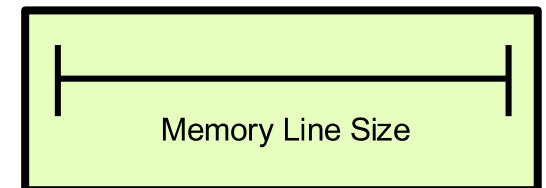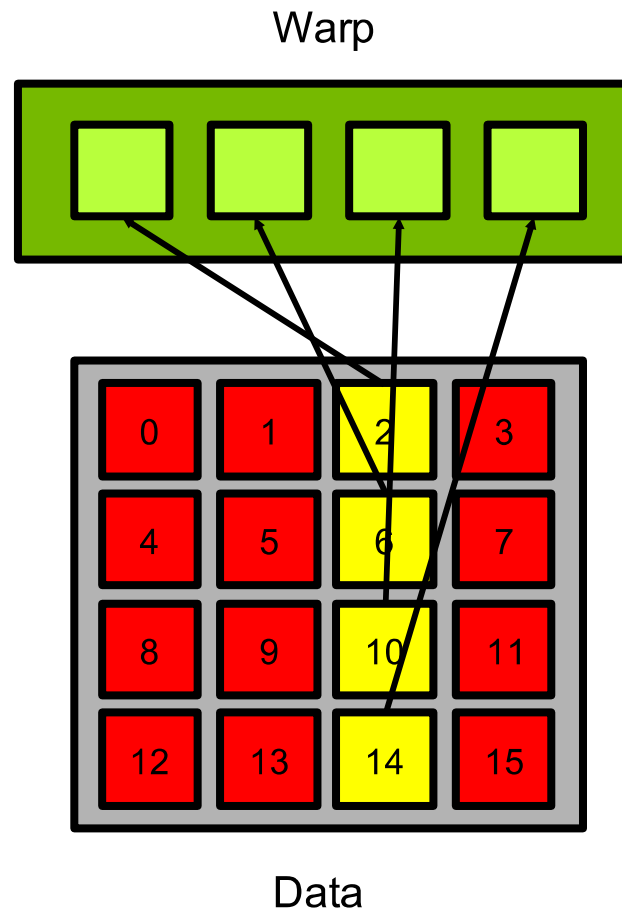12 13 14 15

Memory Line Size

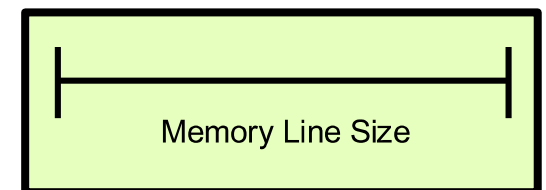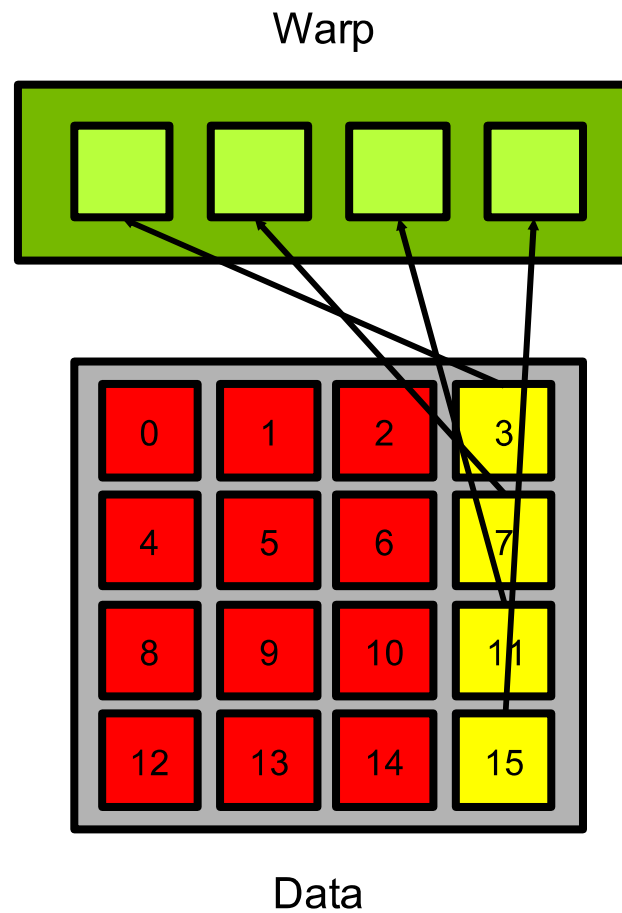Which means (in our example) 4 lines of data will need to be loaded, and 75% of the data loaded will be unused

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

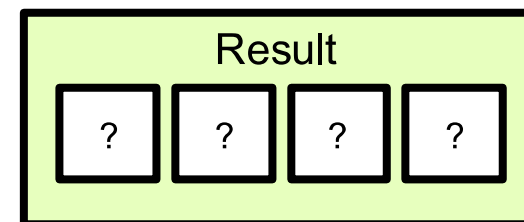Let's compare a kernel that stores the sum of each **column** of a matrix in a result vector
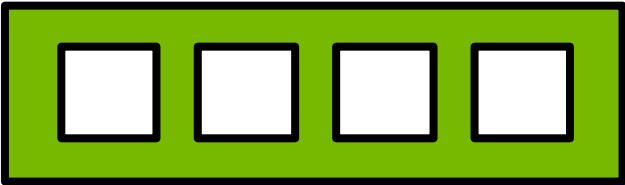
Warp

Result

| ? | ? | ? | ? |

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

Data

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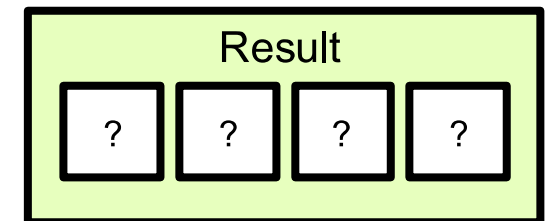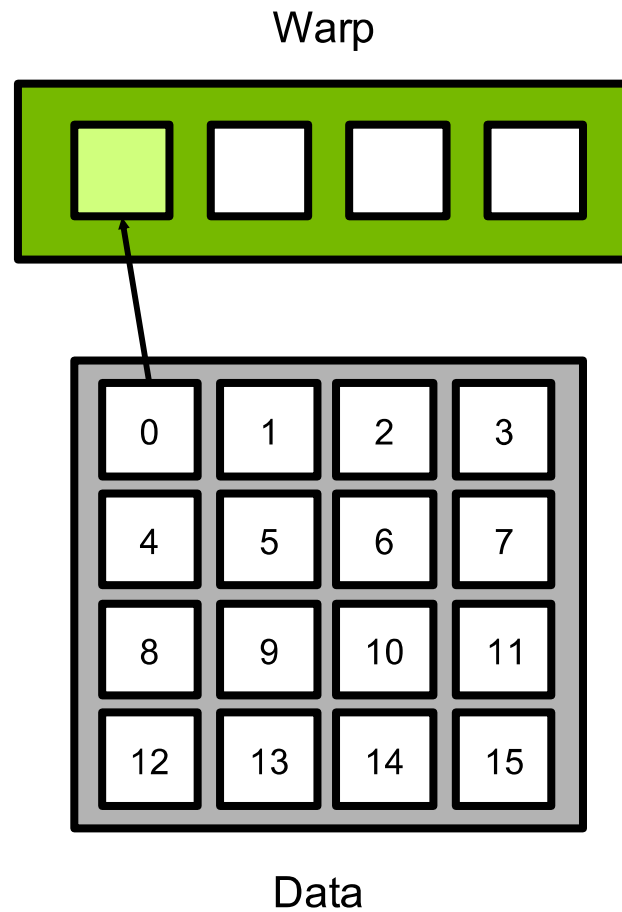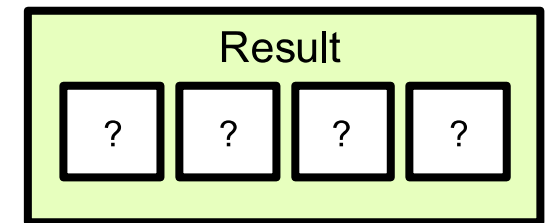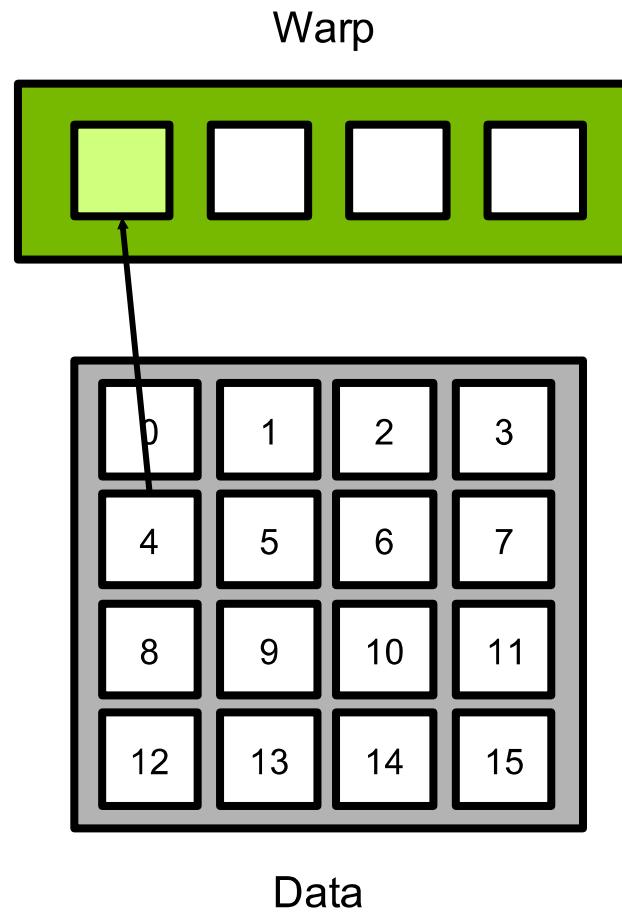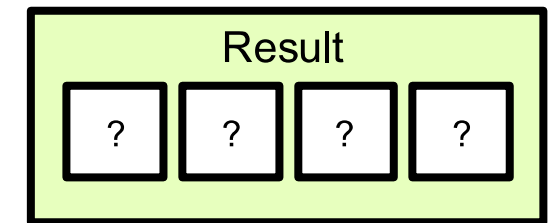
Warp



Result

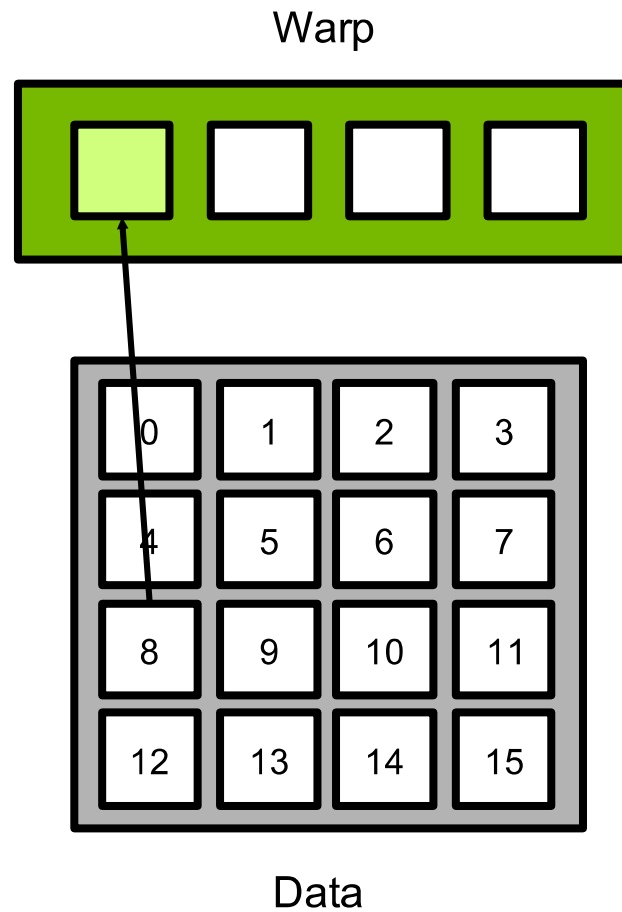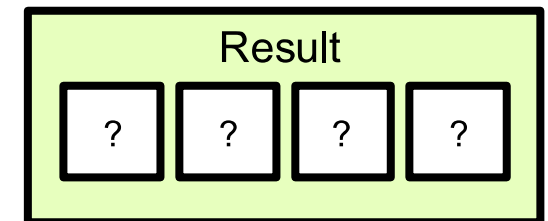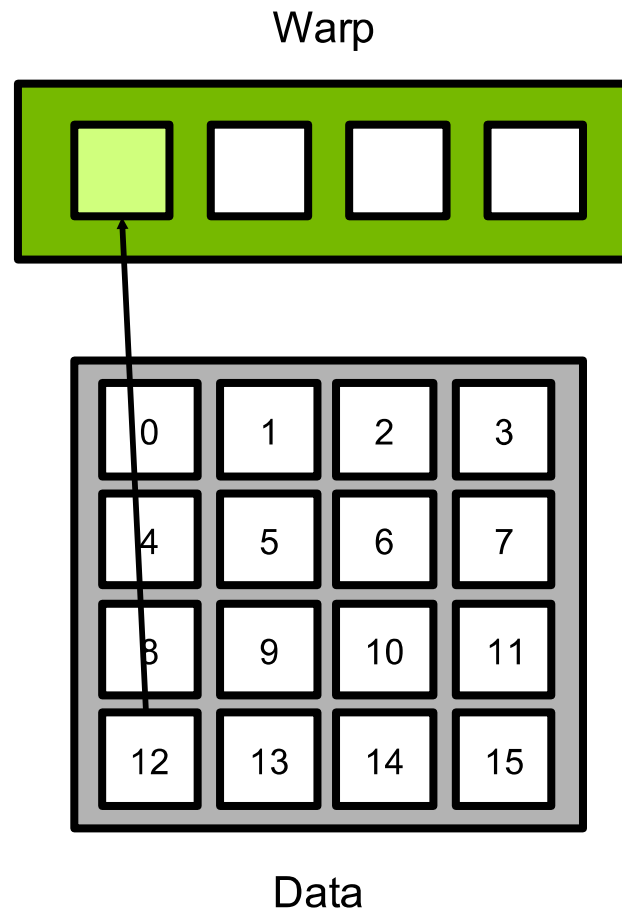| ? | ? | ? | ? |

Sum = 5

Data

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

Warp



Result

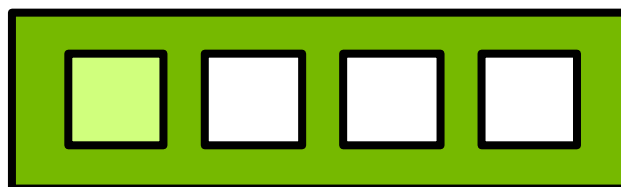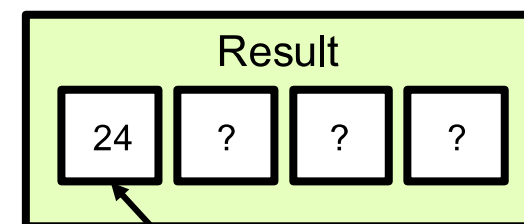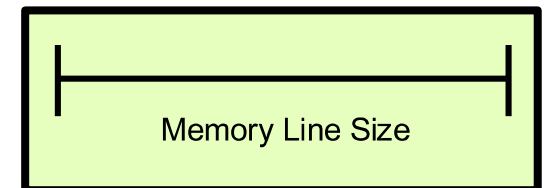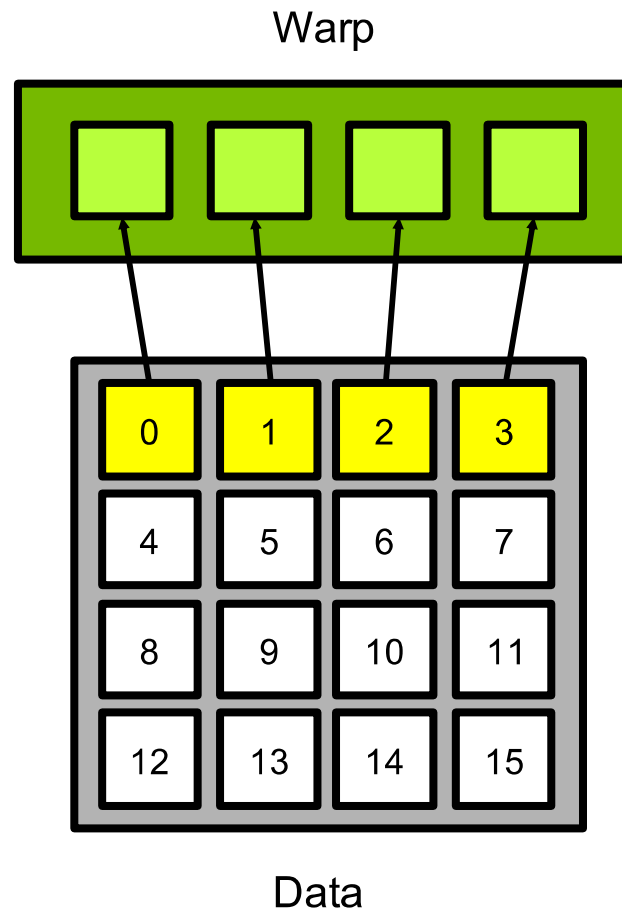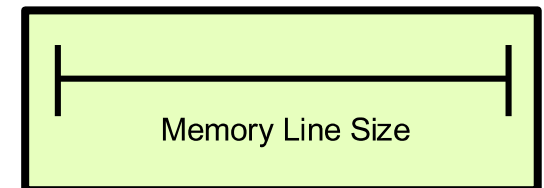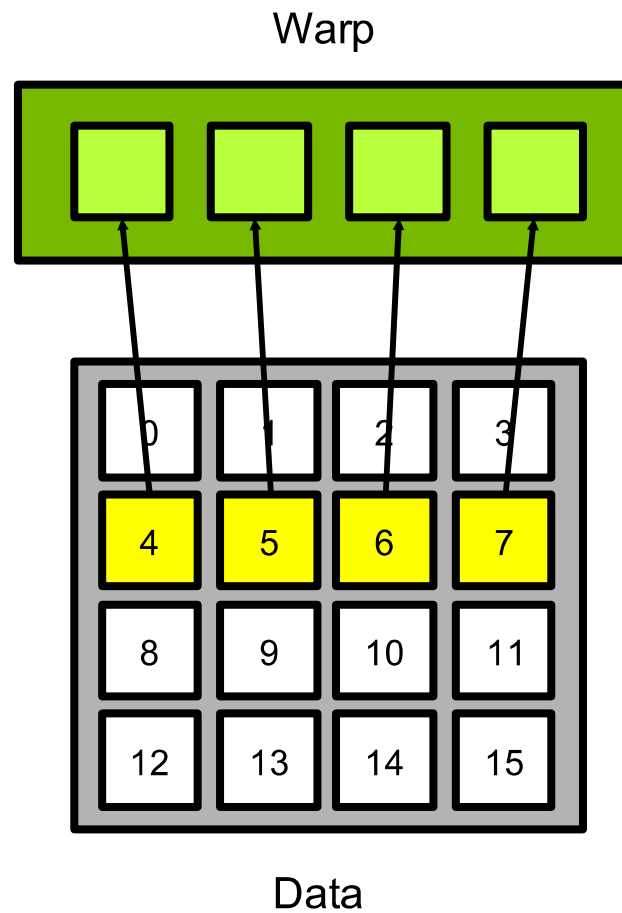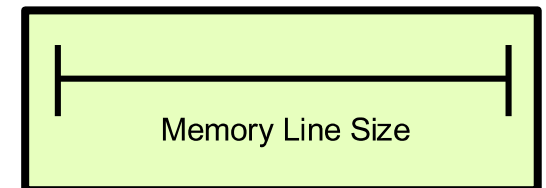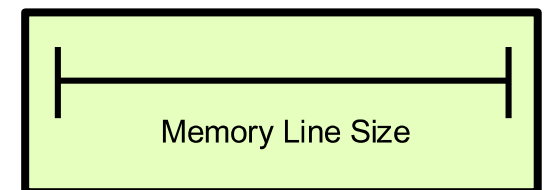| ? | ? | ? | ? |

Sum = 24

Data

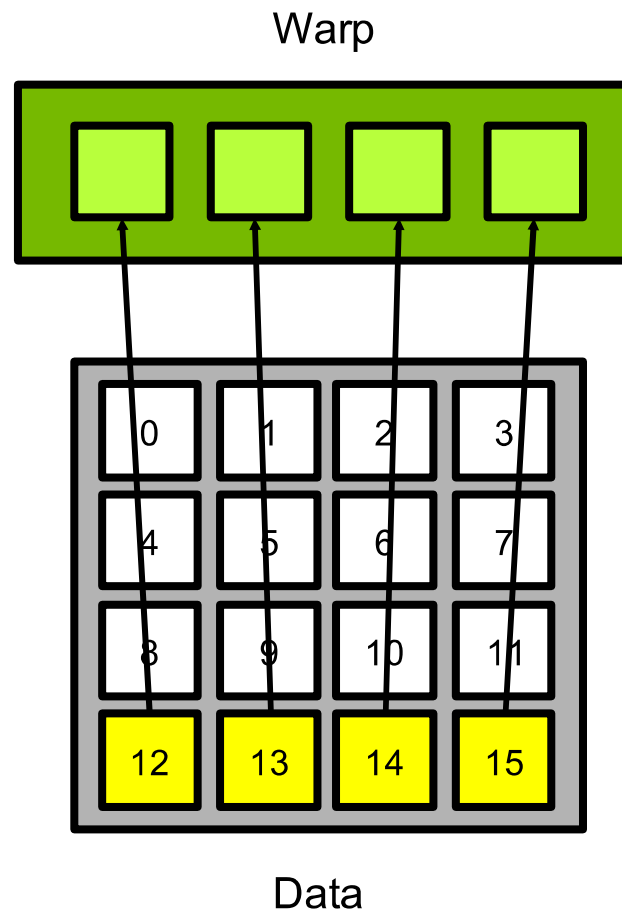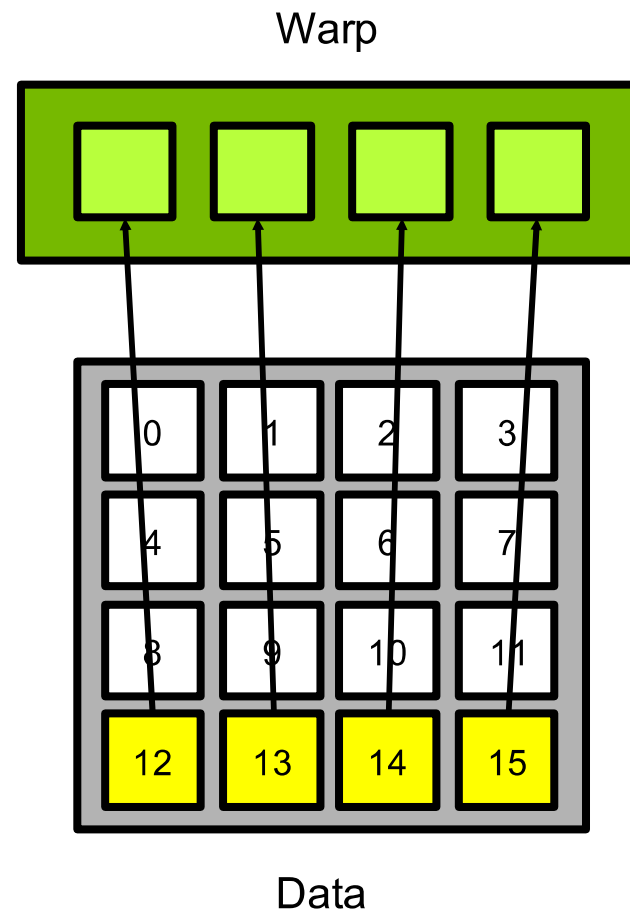Here when we consider the parallel execution, we see that the warp's memory access is coalesced

Warp

Data

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

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



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

Data