



Accelerating Applications with CUDA C/C++

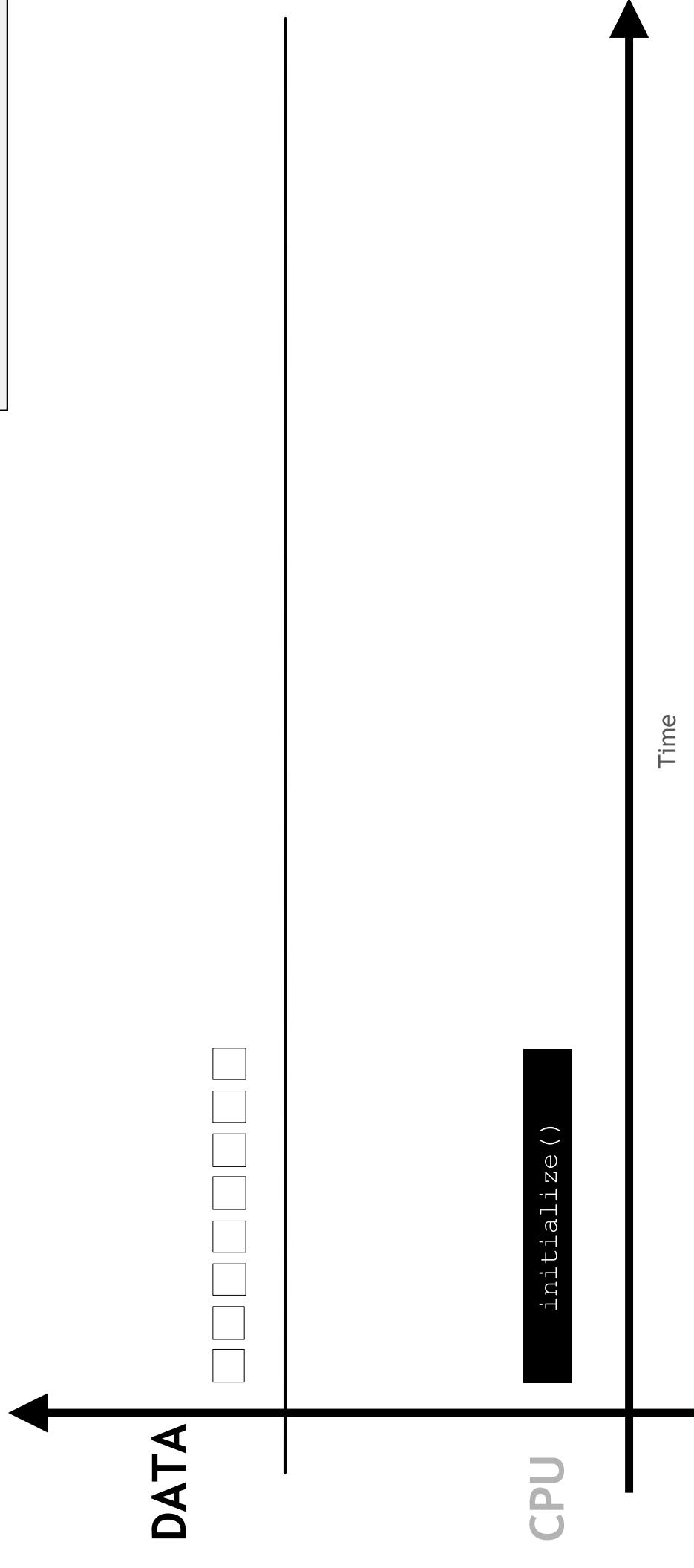
TOPICS

- GPU-accelerated vs. CPU-only Applications
- CUDA Kernel Execution
- Parallel Memory Access
- Appendix: Glossary

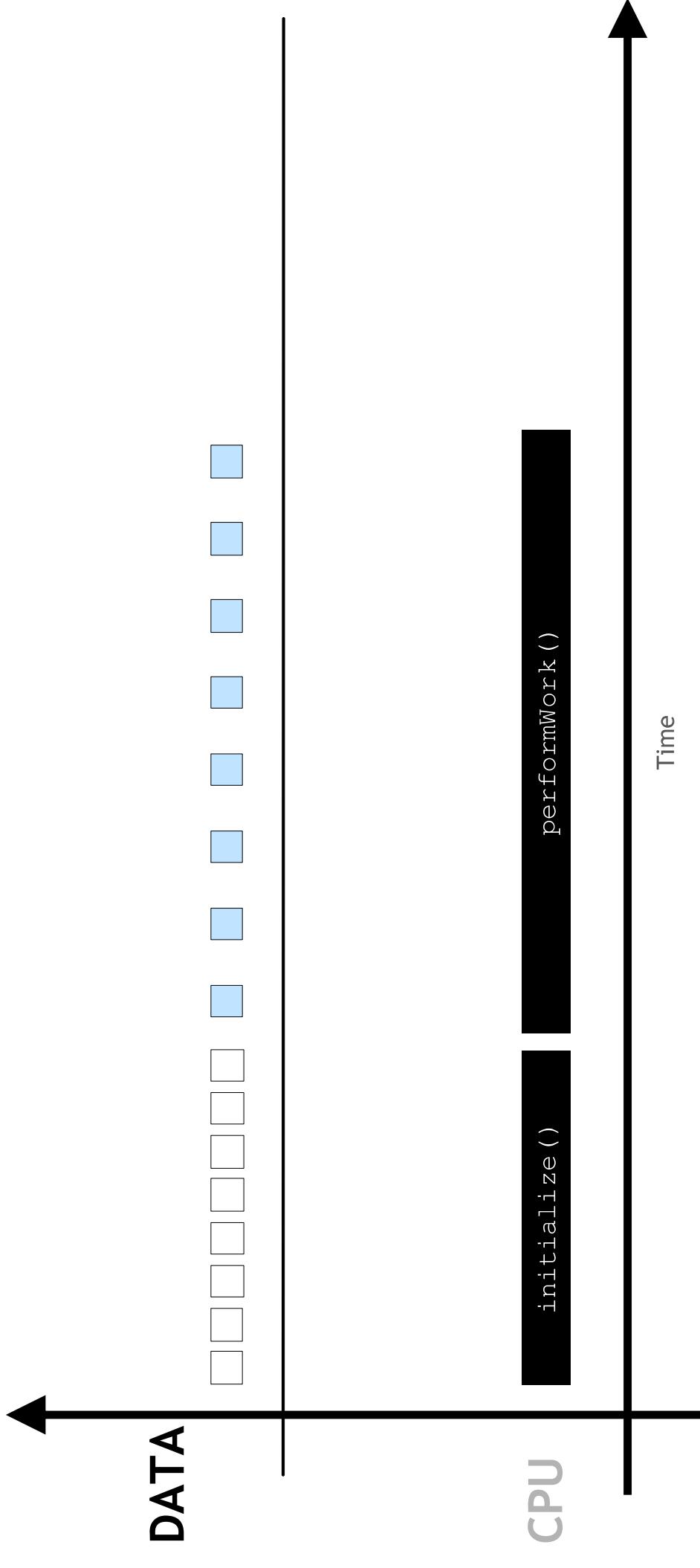
GPU-accelerated vs. CPU-only Applications



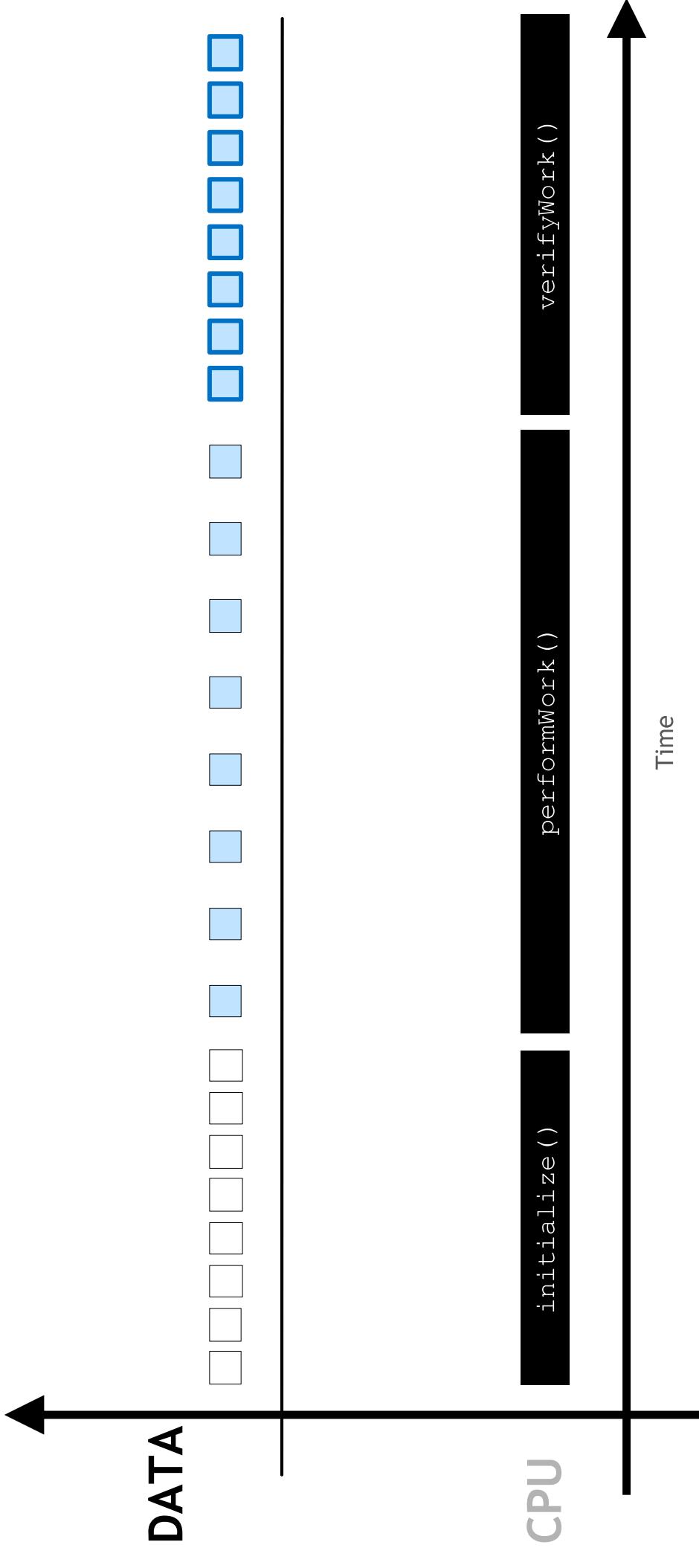
In CPU-only applications data
allocated on CPU

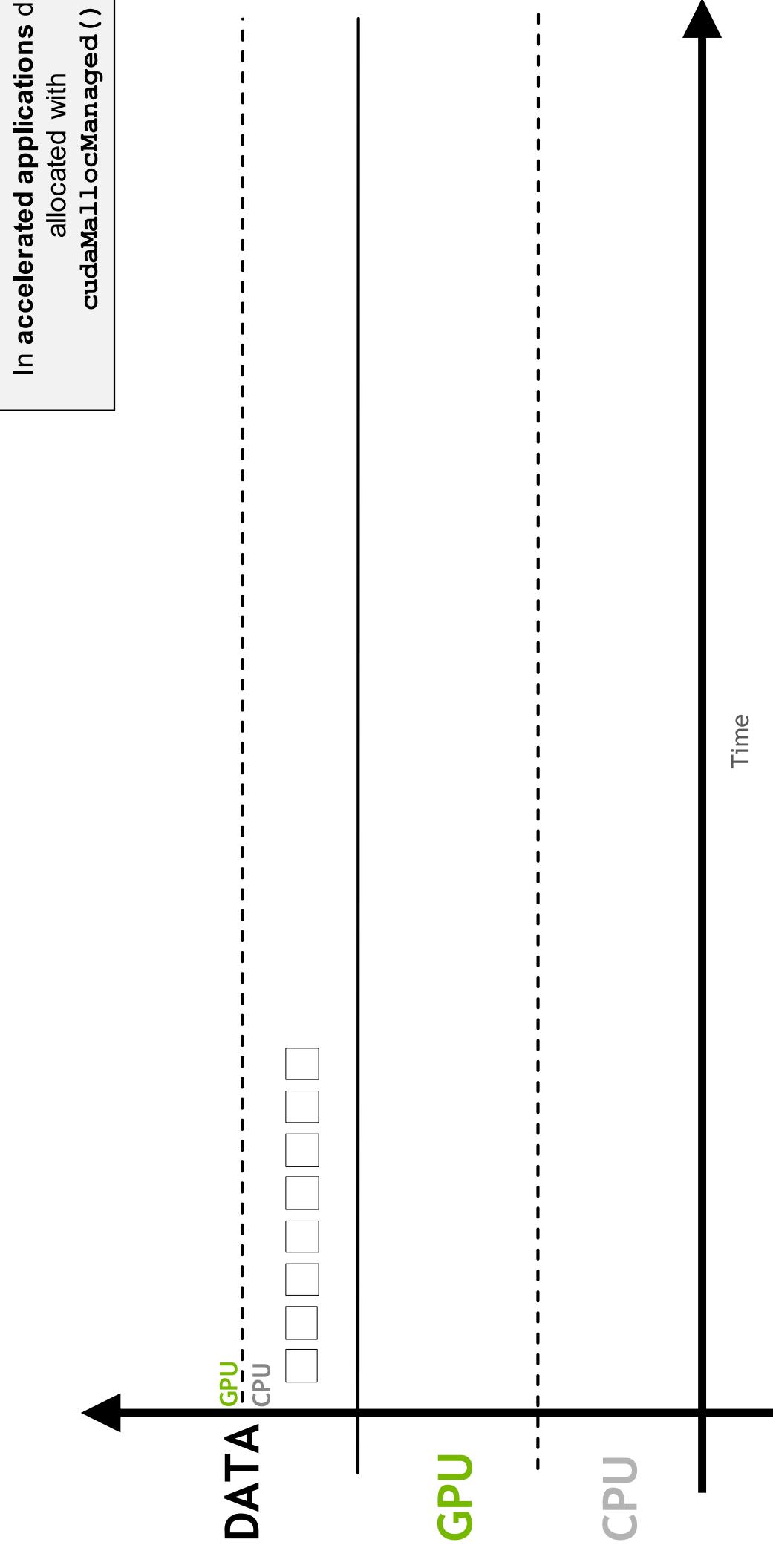


...and all work is performed on

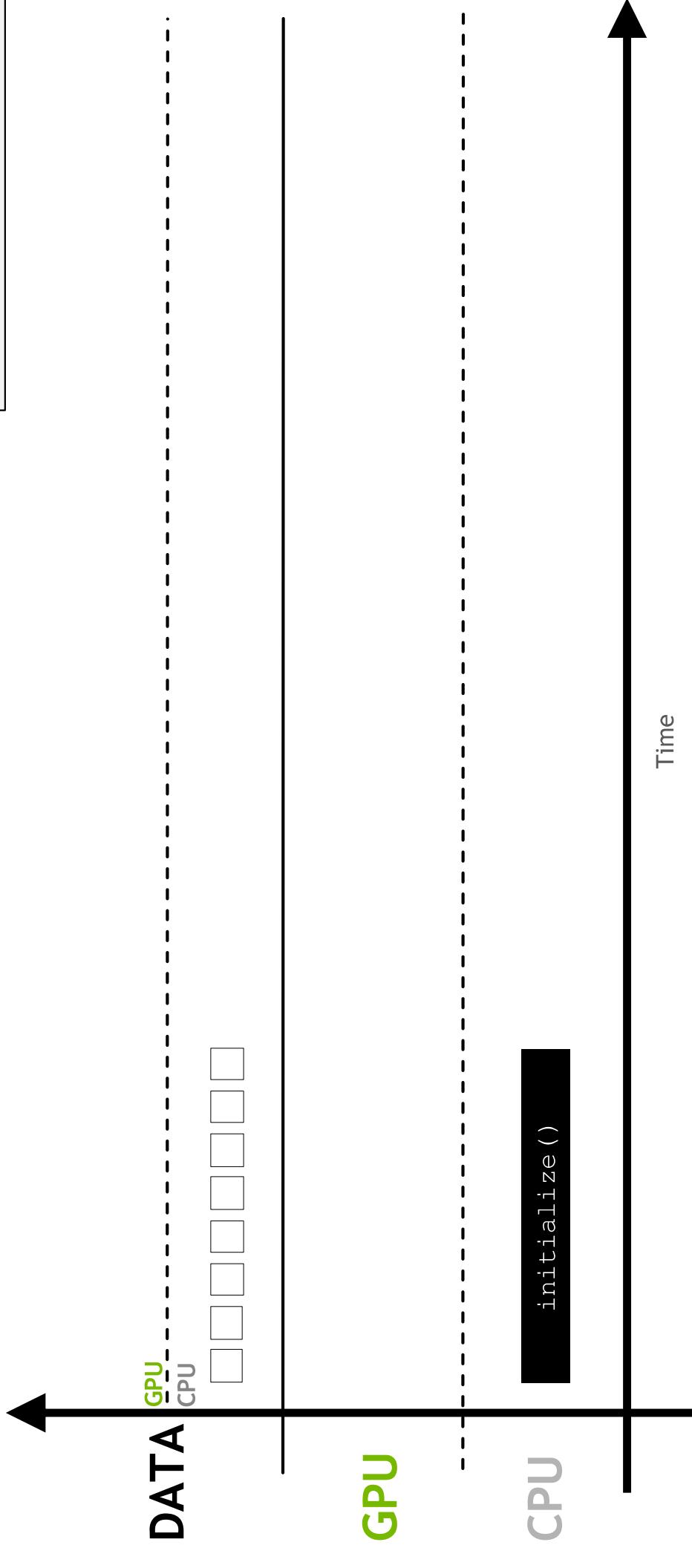


...and all work is performed on

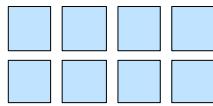




... where it can be accessed
worked on by the CPU



... and automatically migrated to GPU where parallel work can be done



DATA

GPU

CPU



GPU

`performWork()`

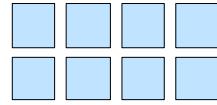
CPU

`initialize()`

Time



Work on the GPU is **asynchronous**
and CPU can work at the same time



DATA

GPU

CPU



GPU

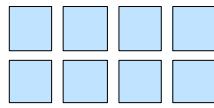
performWork ()

CPU

initialize () cpuWork ()

Time

CPU code can sync with the asynchronous GPU work, waiting to complete, with `cudaDeviceSynchronize`



DATA

GPU

CPU



GPU

CPU

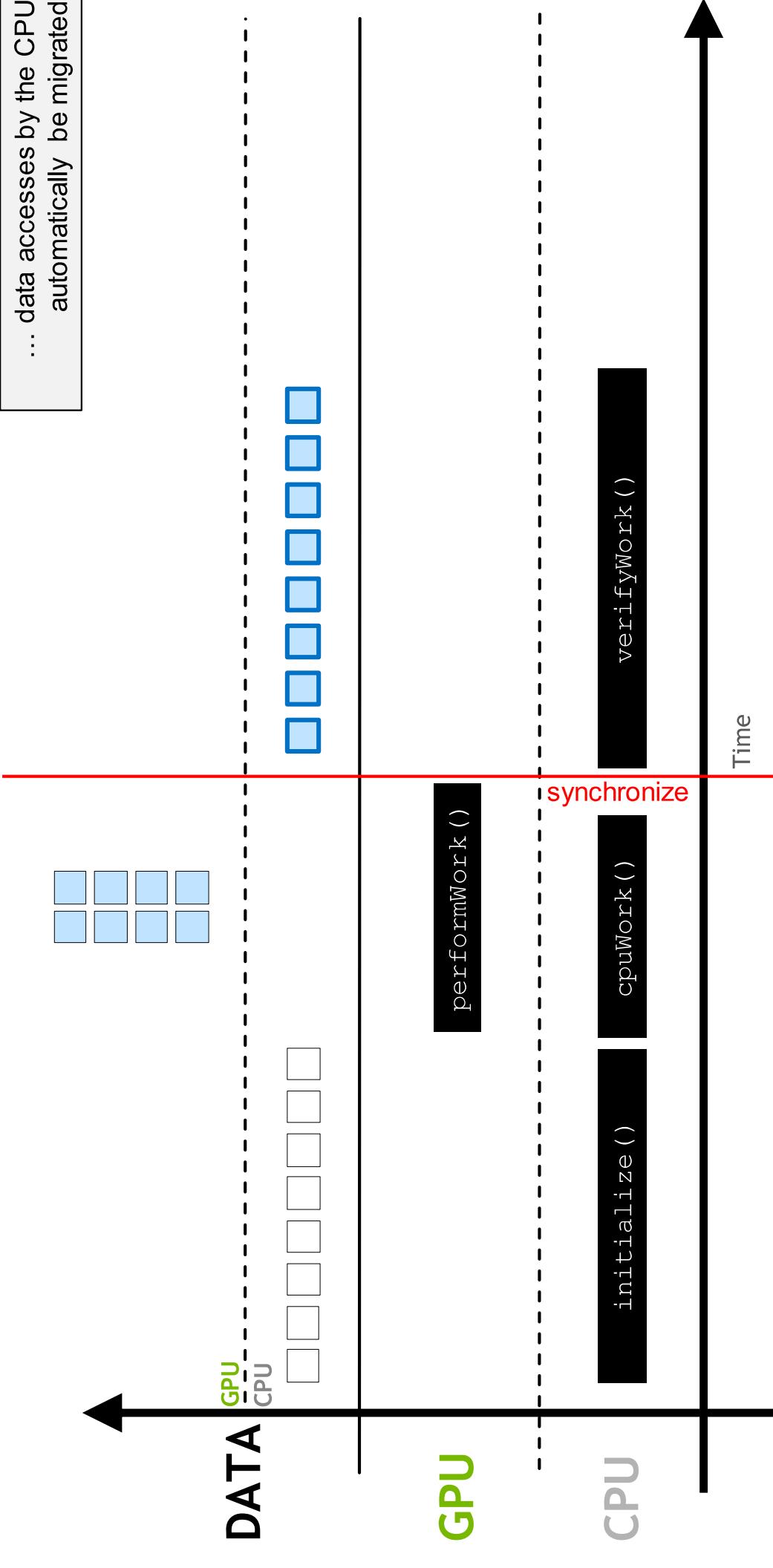
`performWork ()`

`initialize ()` `cpuWork ()`

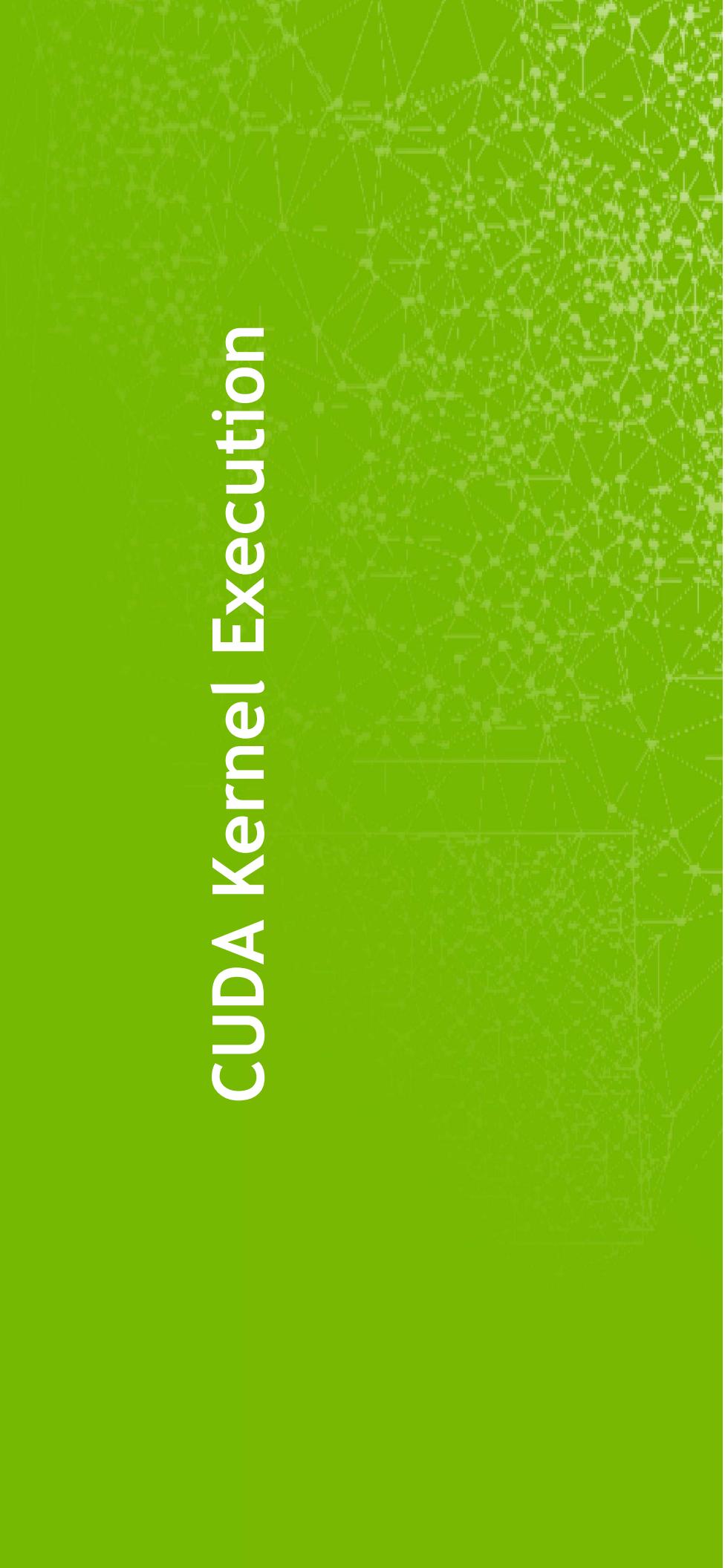
`synchronize`

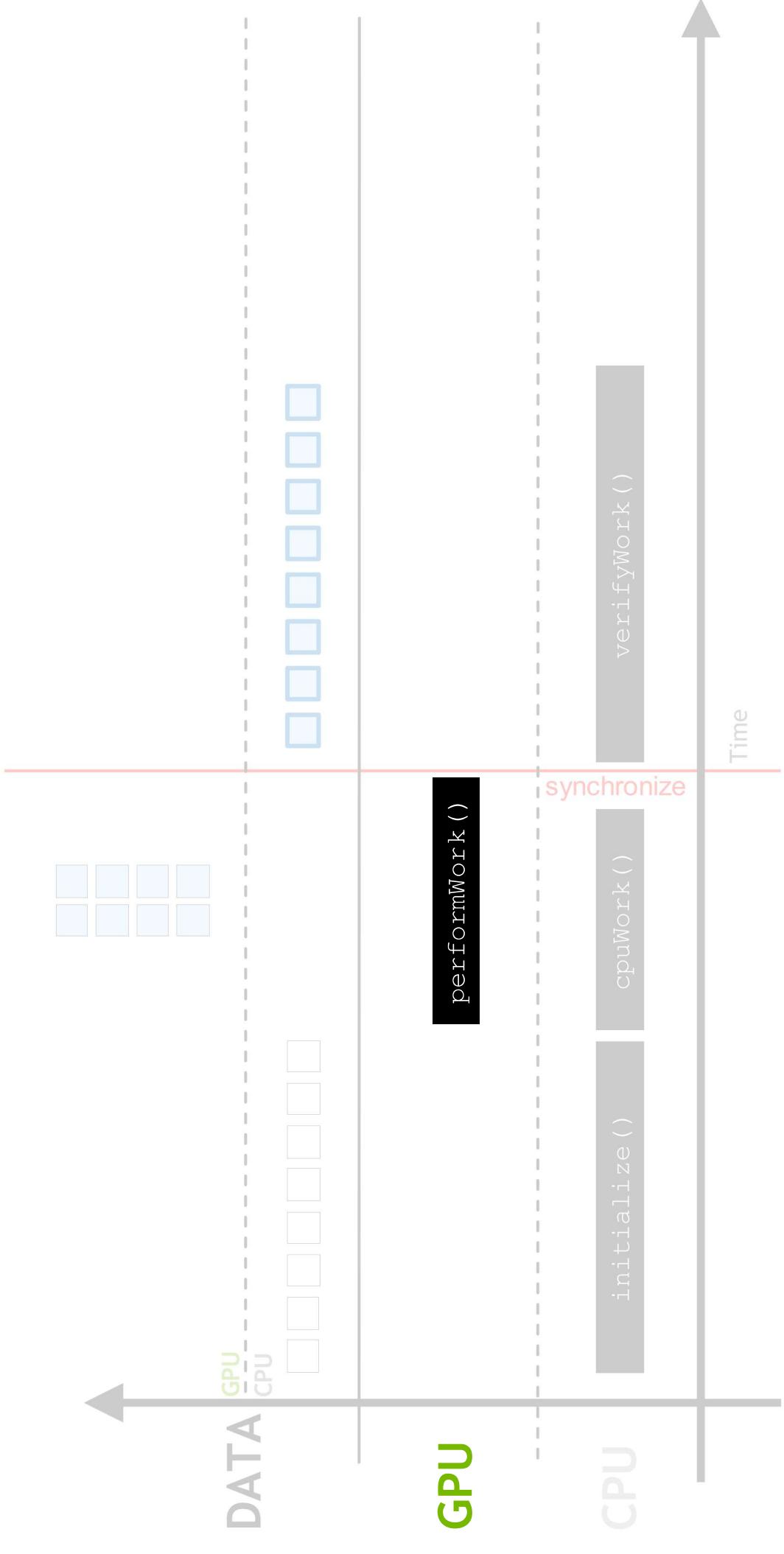
Time

... data accesses by the CPU
automatically be migrated

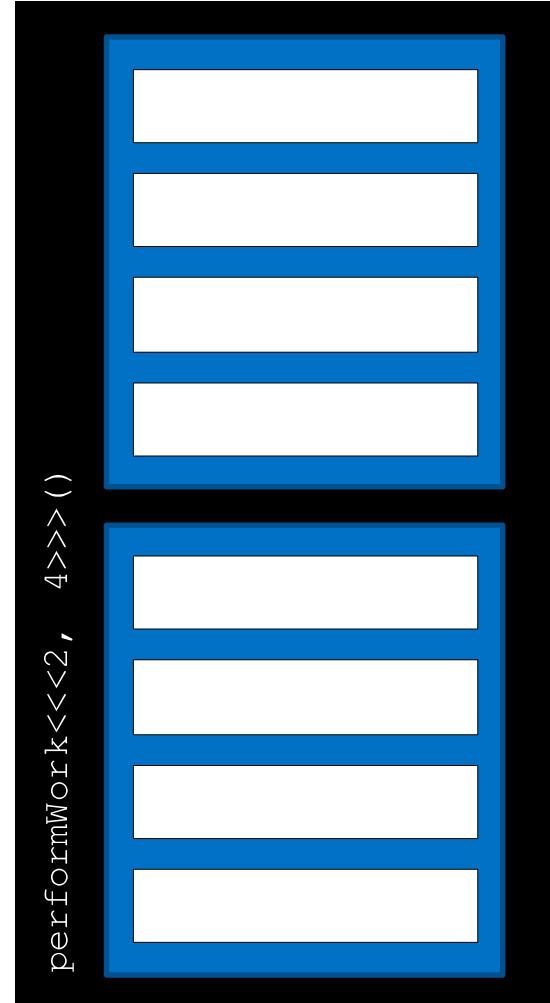


CUDA Kernel Execution



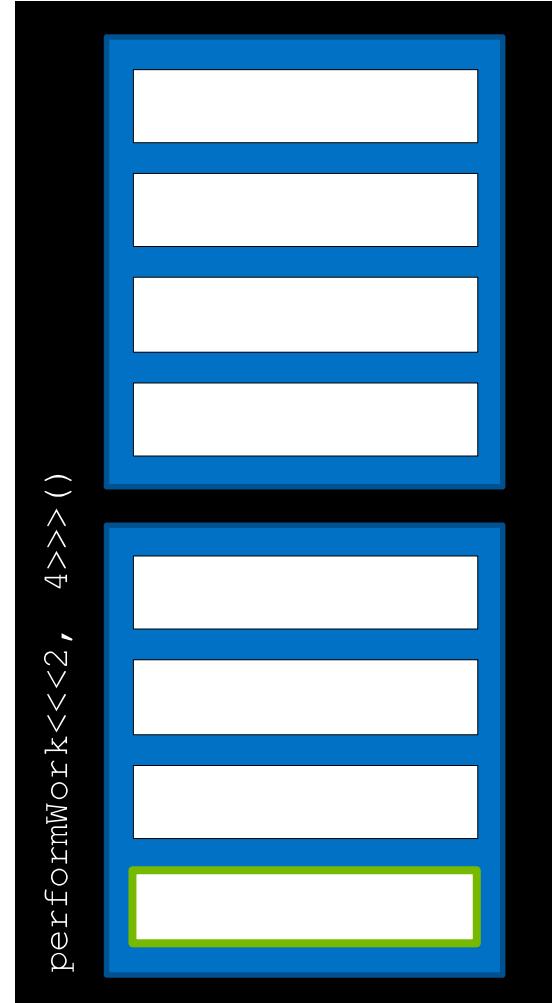


GPUs do work in parallel



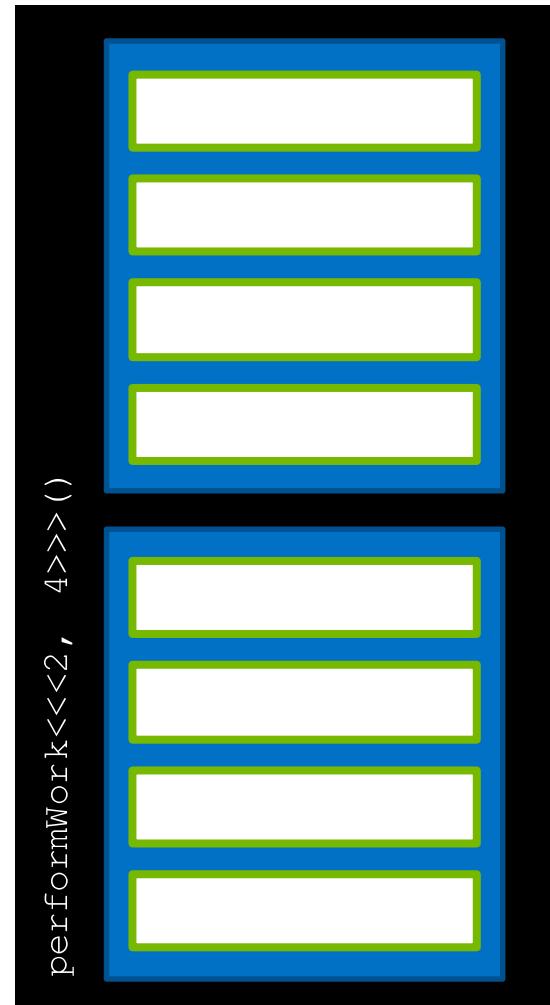
GPU

GPU work is done in a three



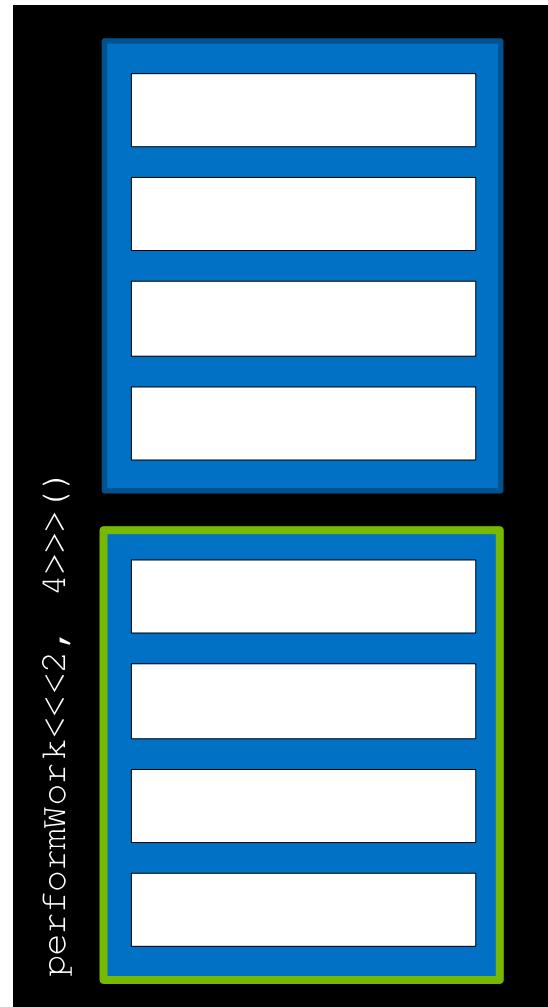
GPU

Many threads run in parallel



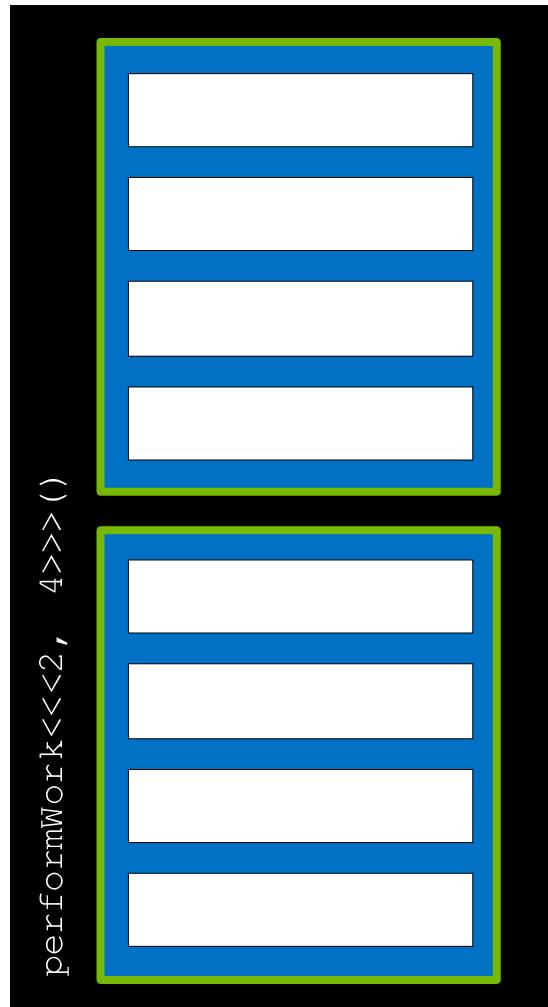
GPU

A collection of threads is a **block**



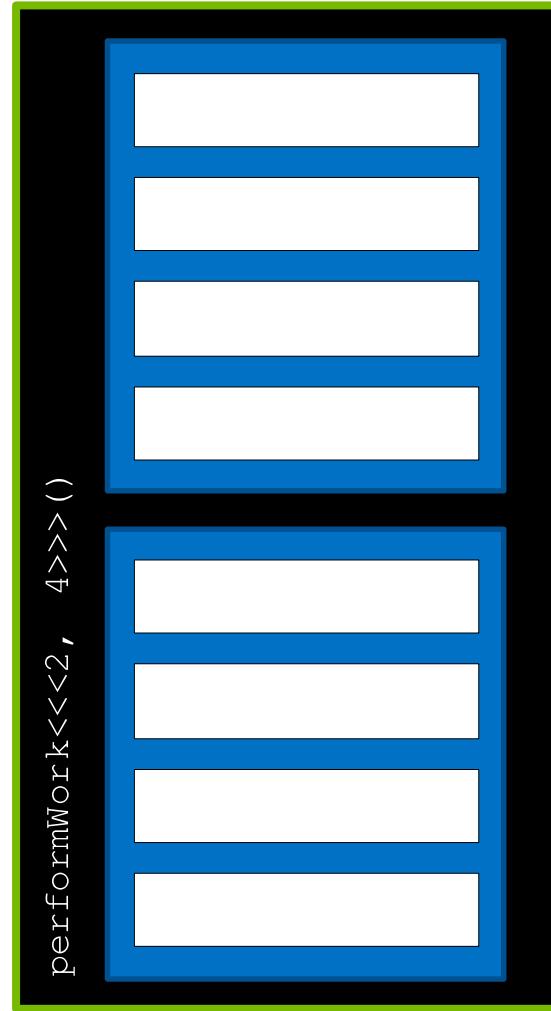
GPU

There are many blocks



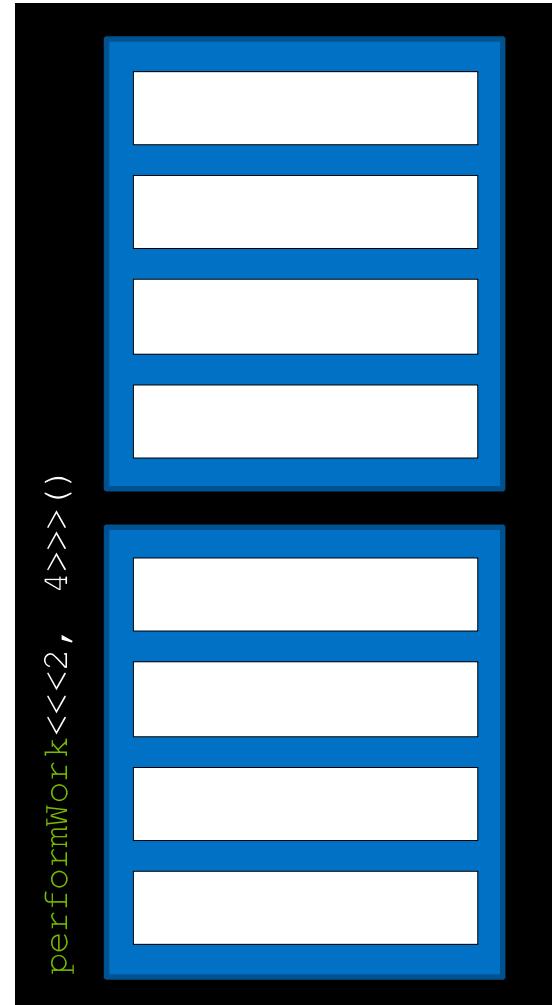
GPU

A collection of blocks is a group



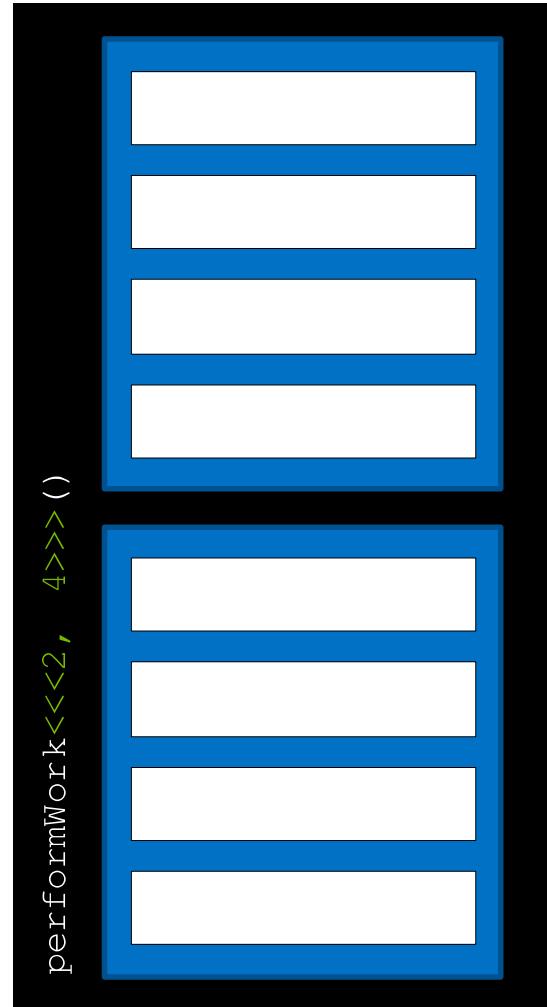
GPU

GPU functions are called **kernel**



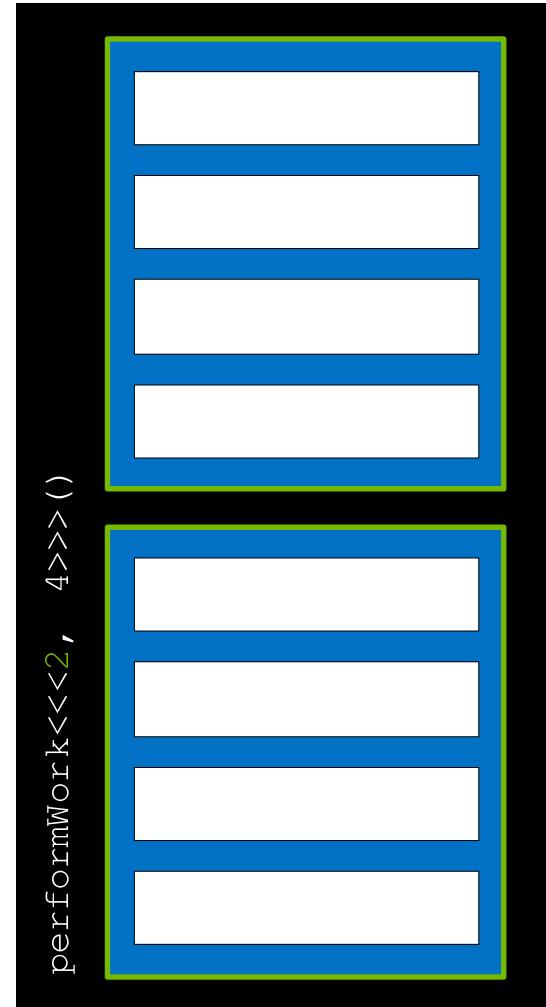
GPU

Kernels are launched with
execution configuration



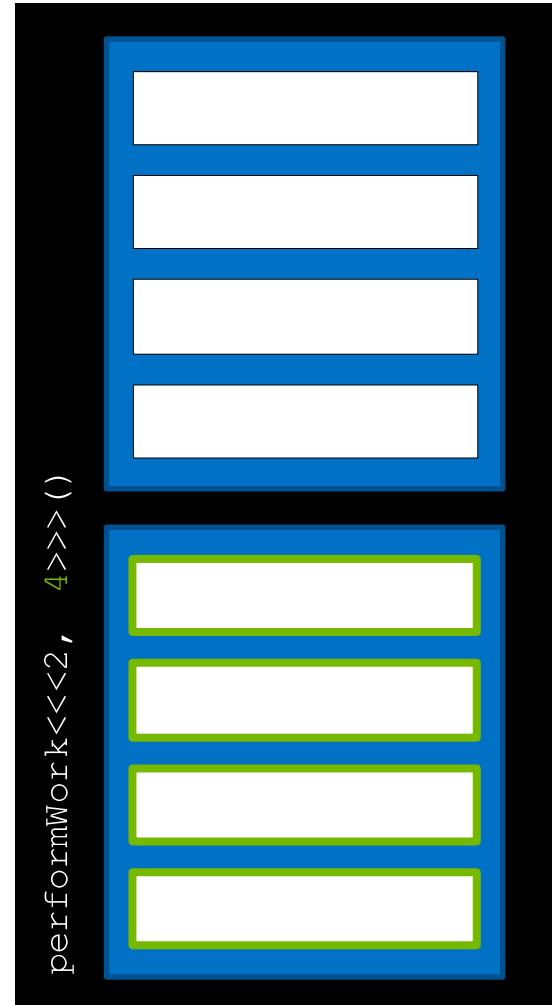
GPU

The execution configuration depends on the number of blocks in the grid



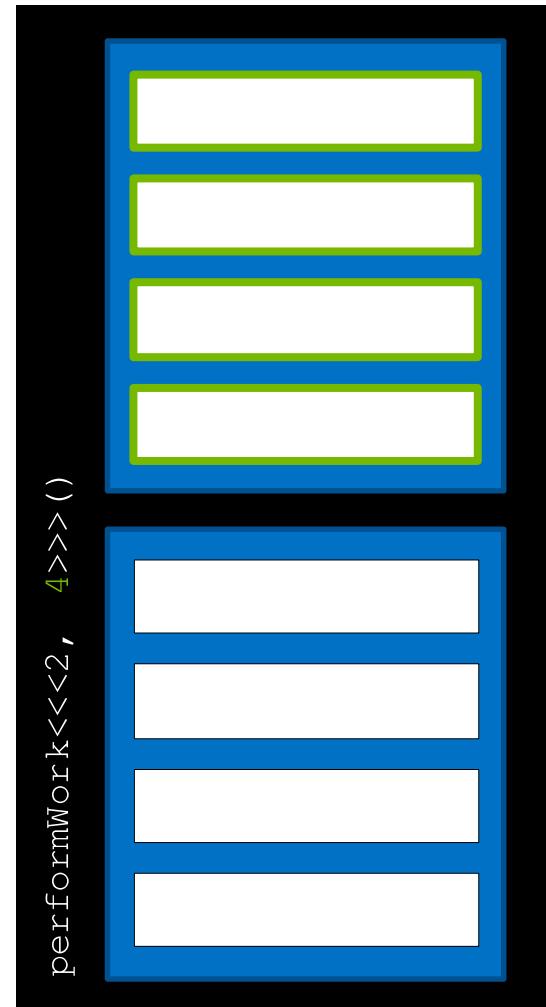
GPU

... as well as the number of threads
each block



GPU

Every block in the grid contains
same number of threads

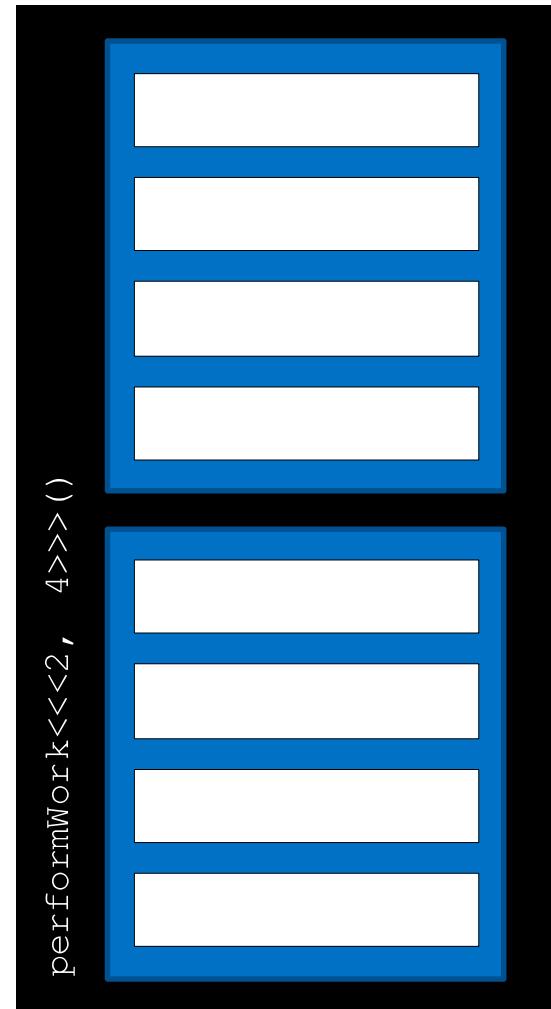


GPU

CUDA-Provided Thread Hierarchy Variables

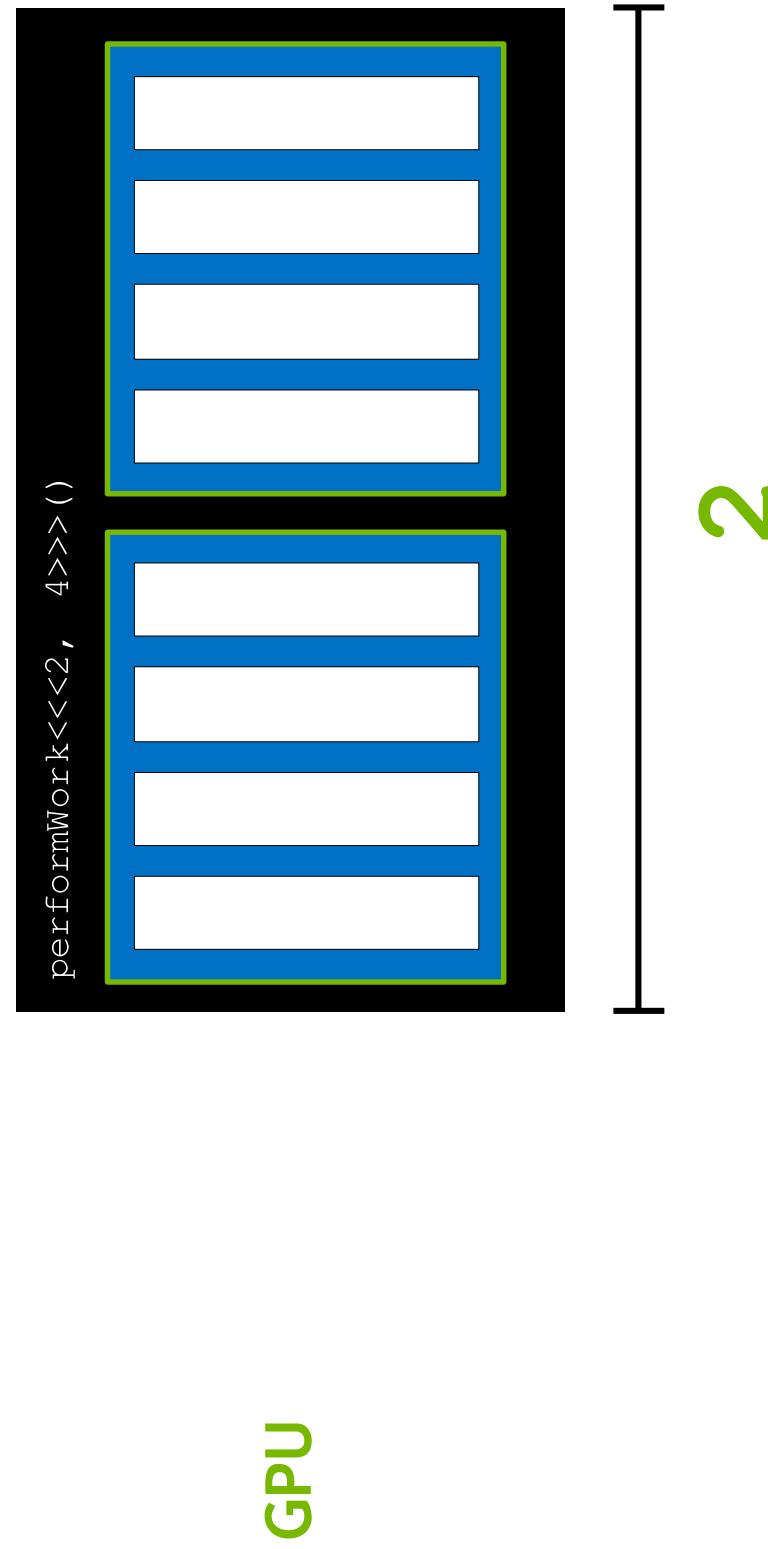


Inside kernels definitions, CU provided variables describe executing thread, block, and (

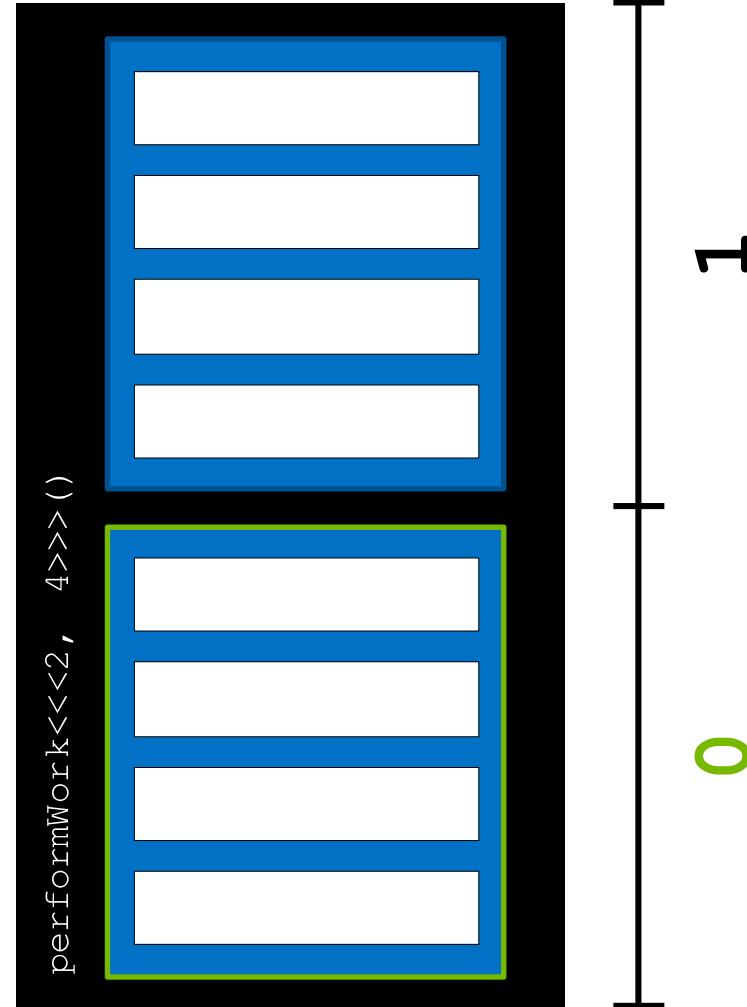


GPU

`gridDim.x` is the number of blocks
the grid, in this case 2

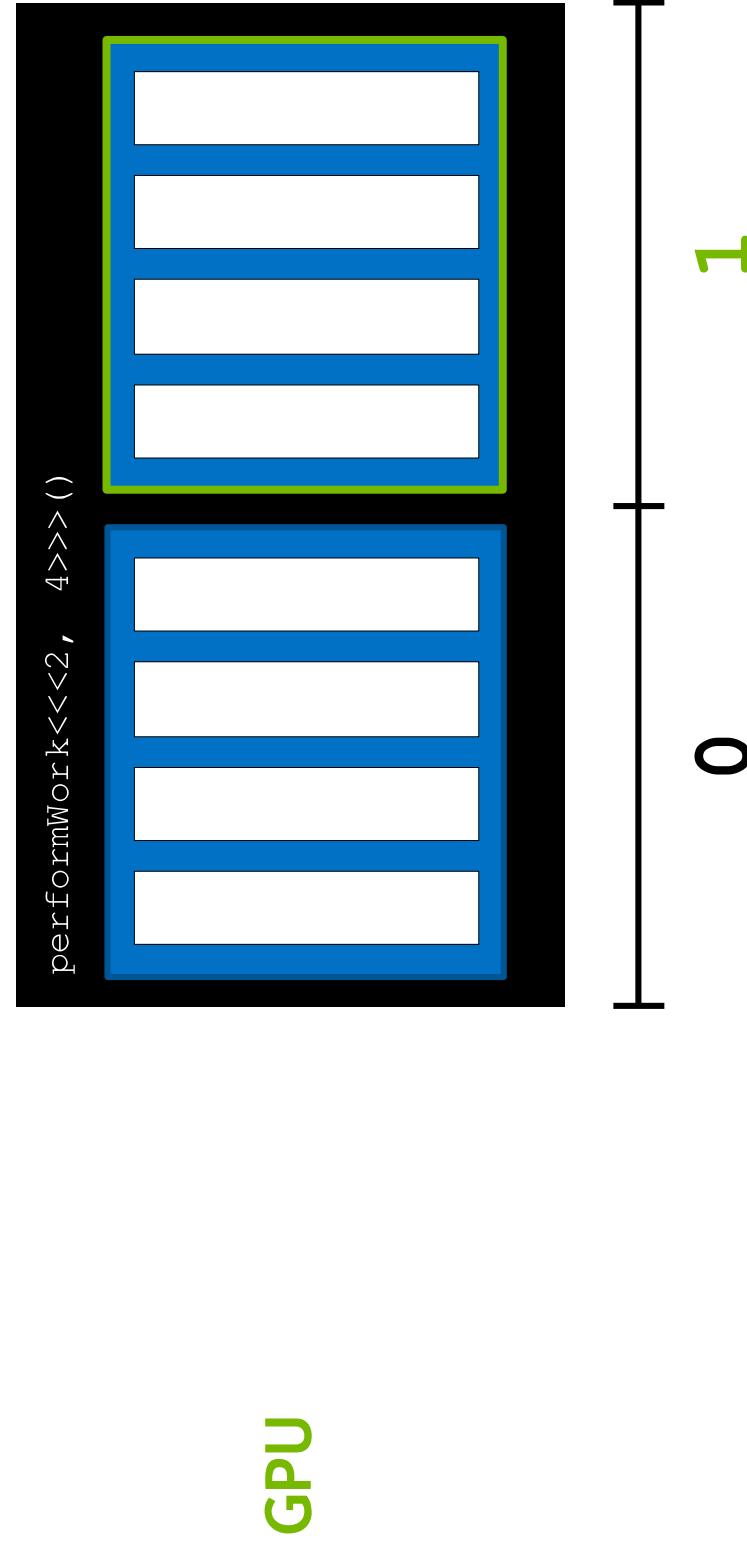


blockIdx.x is the index of current block within the grid, in case 0

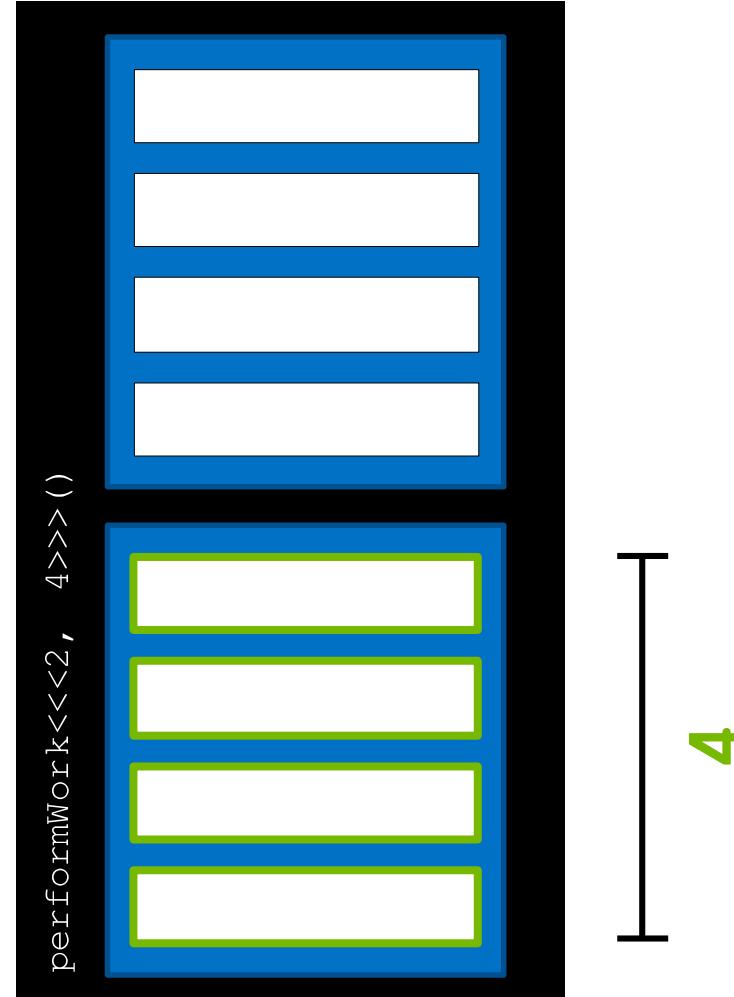


GPU

blockIdx.x is the index of
current block within the grid, in
case 1

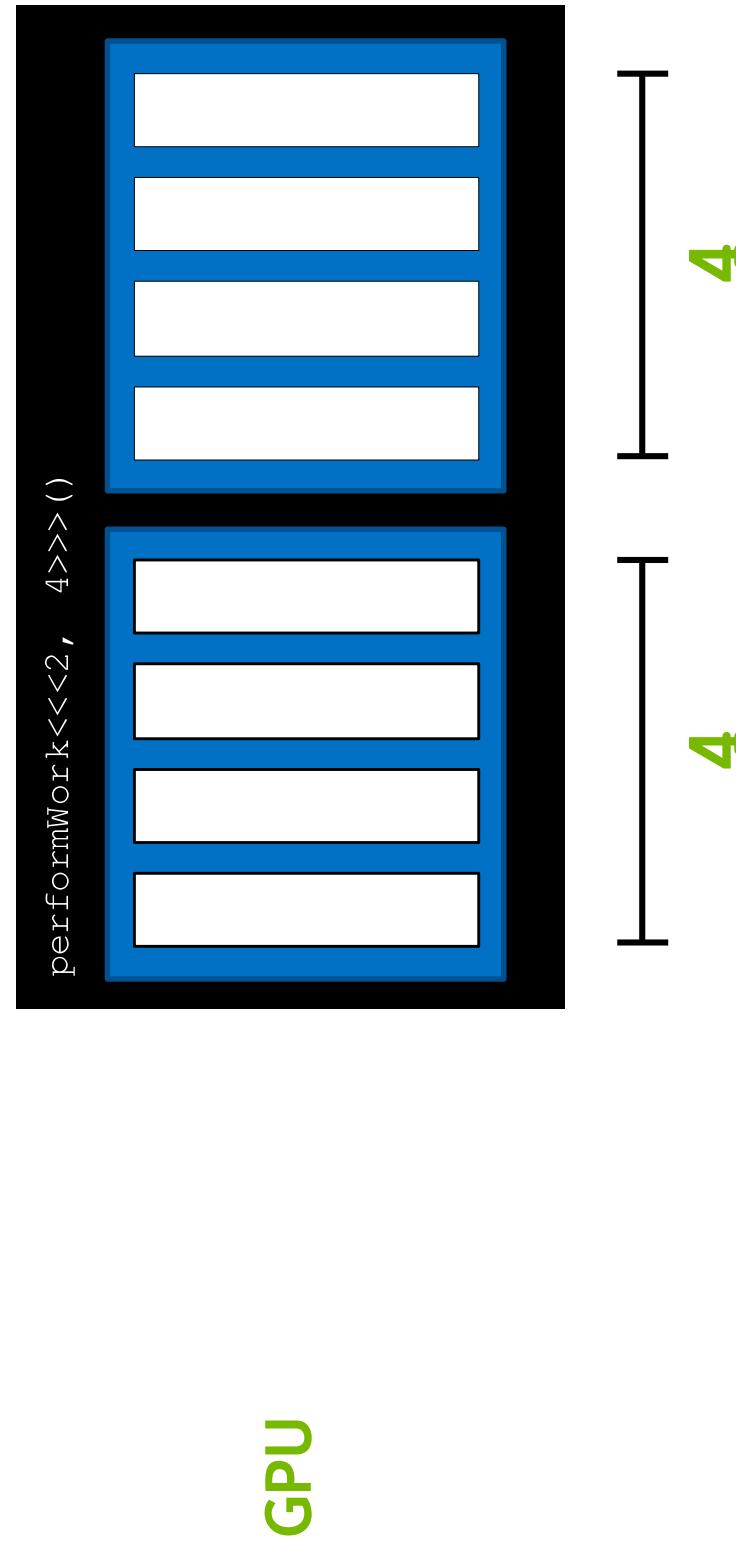


Inside a kernel `blockDim.x` denotes
the number of threads in a block
this case 4

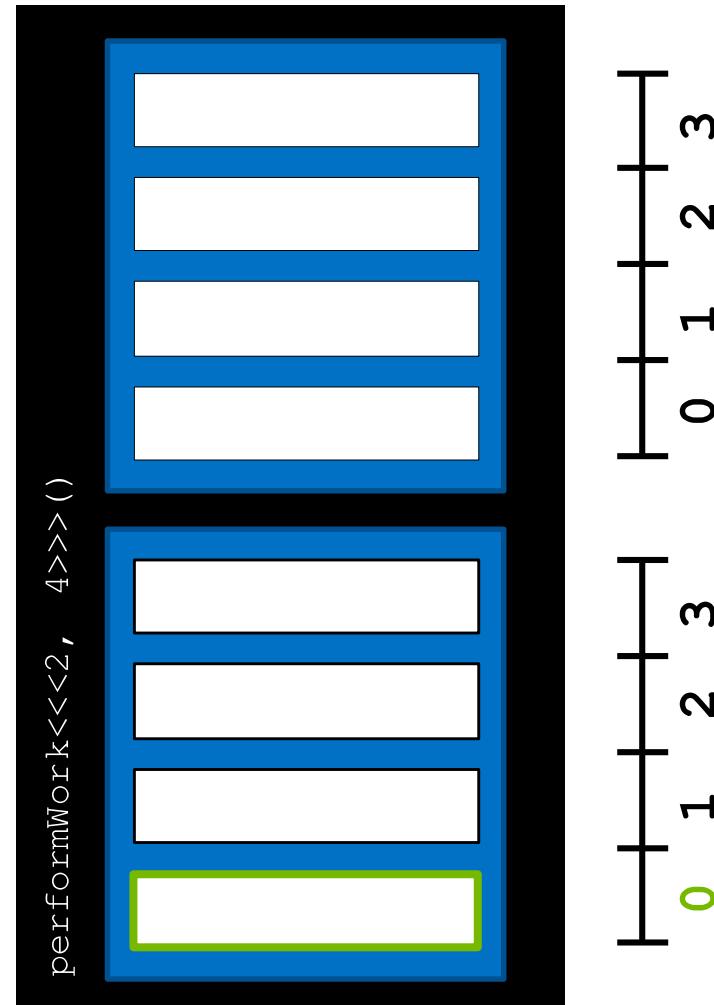


GPU

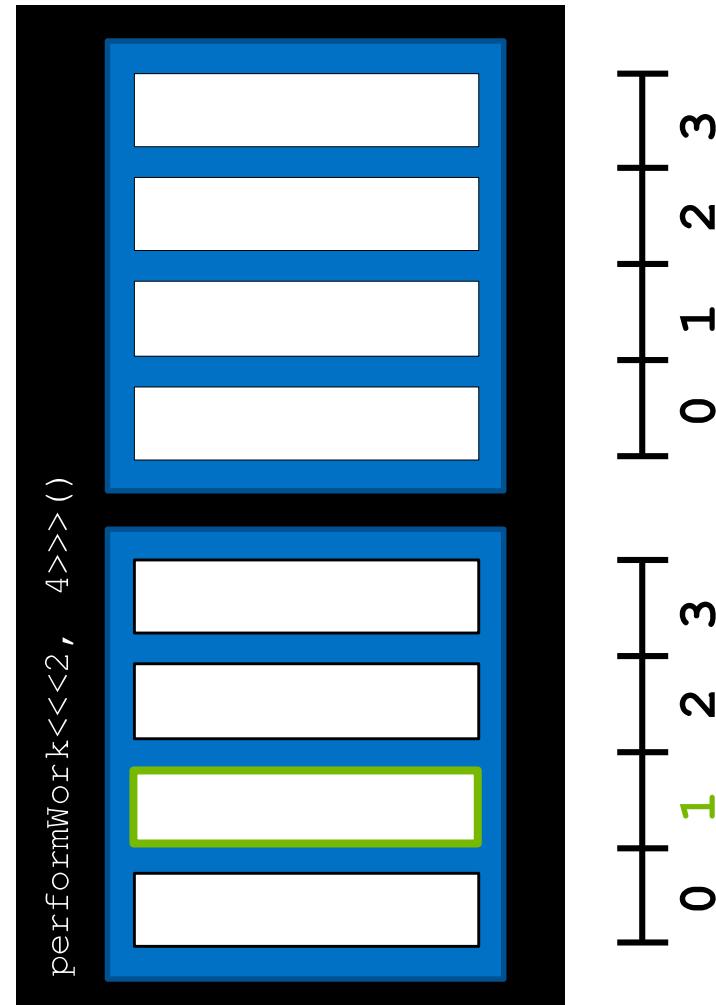
All blocks in a grid contain the same number of threads



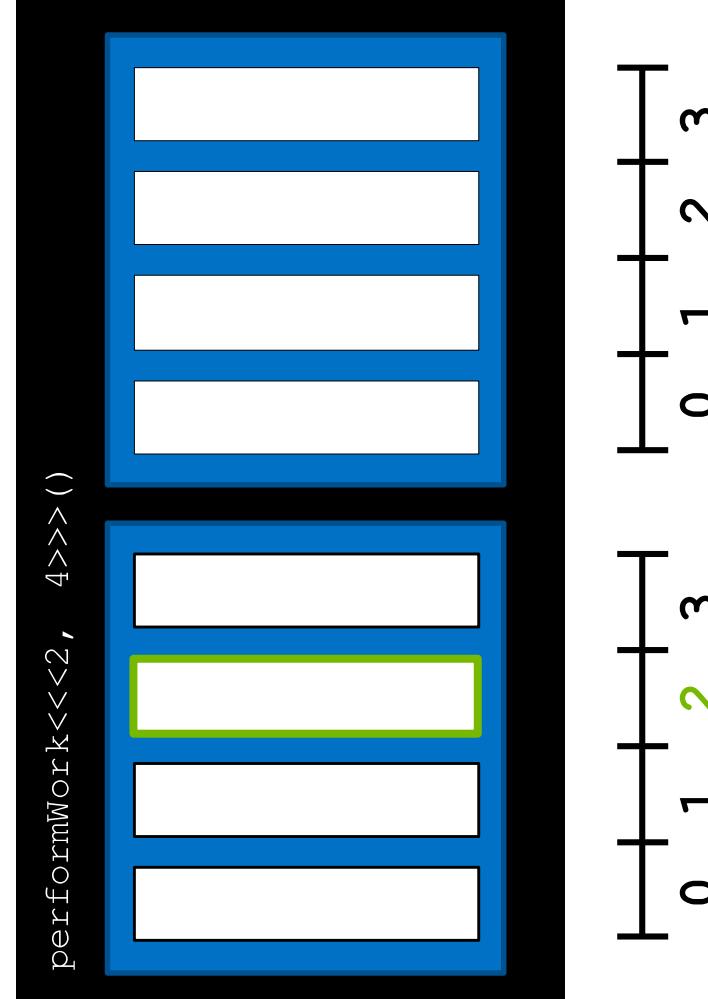
Inside a kernel **threadIdx**
describes the index of the thread
a block. In this case 0



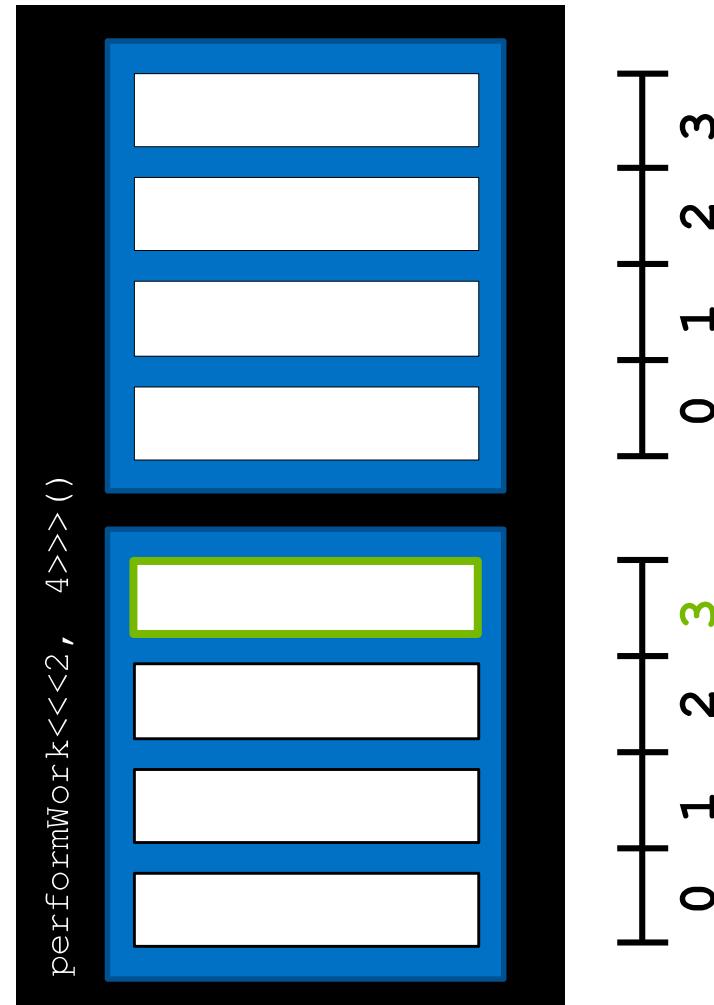
Inside a kernel **threadIdx**
describes the index of the thread
a block. In this case 1



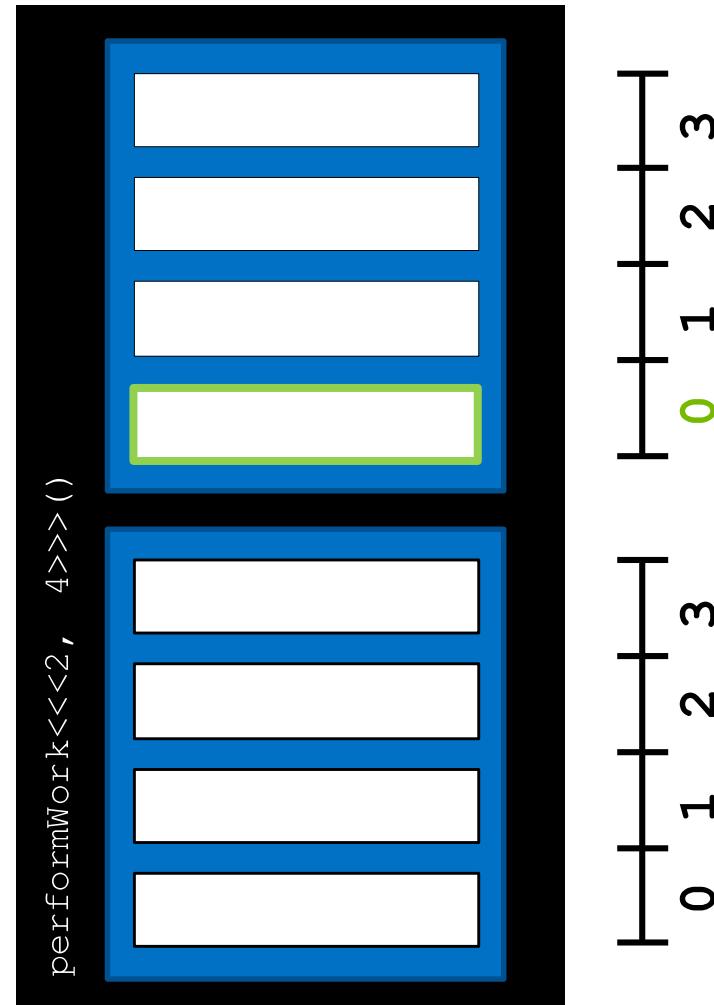
Inside a kernel **threadIdx**
describes the index of the thread
a block. In this case 2



Inside a kernel **threadIdx**
describes the index of the thread
a block. In this case 3

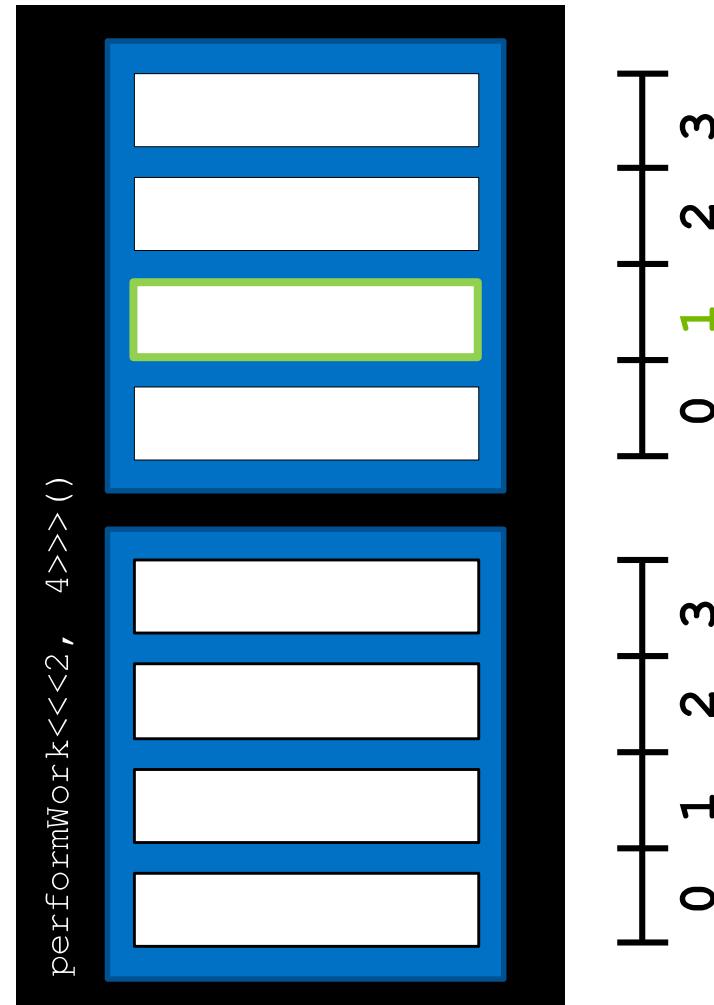


Inside a kernel **threadIdx**
describes the index of the thread
a block. In this case 0

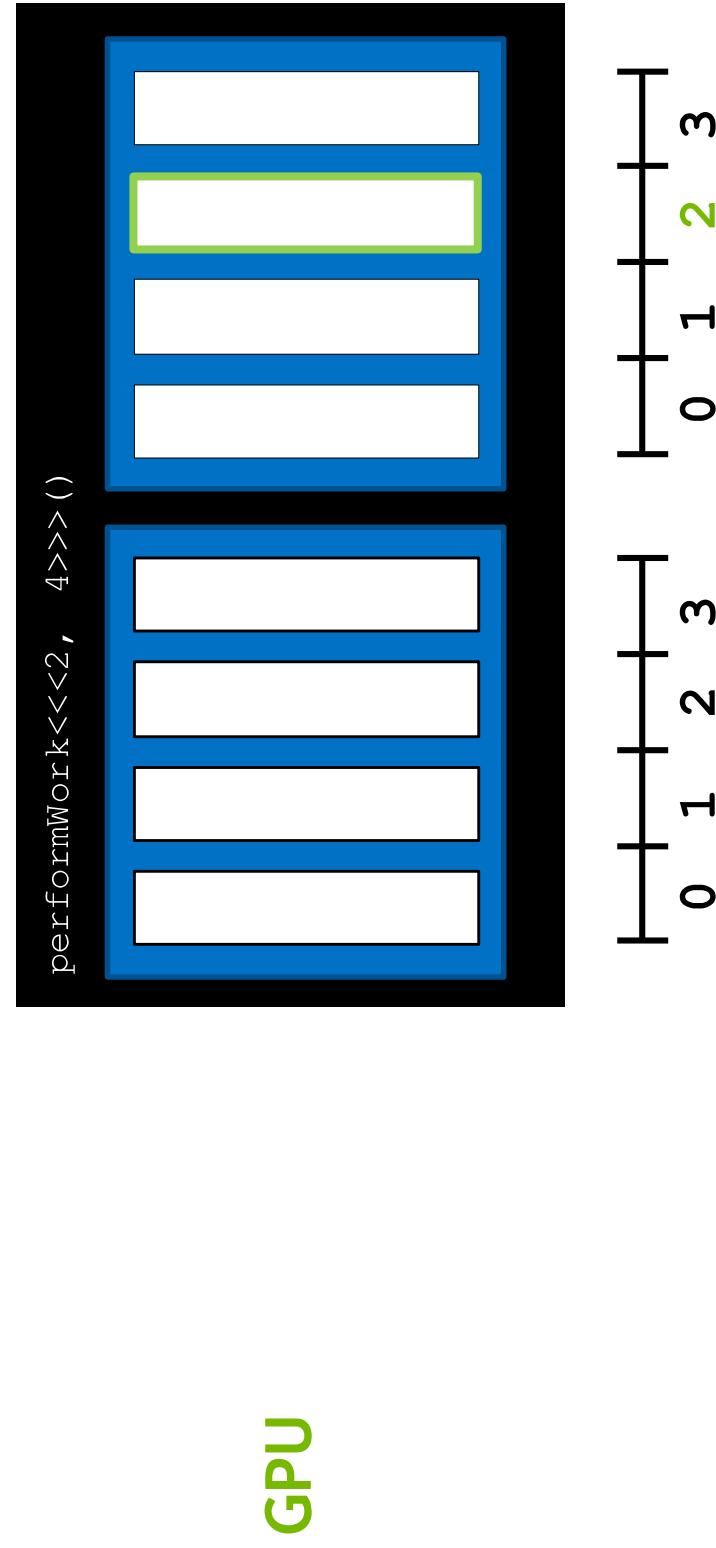


GPU

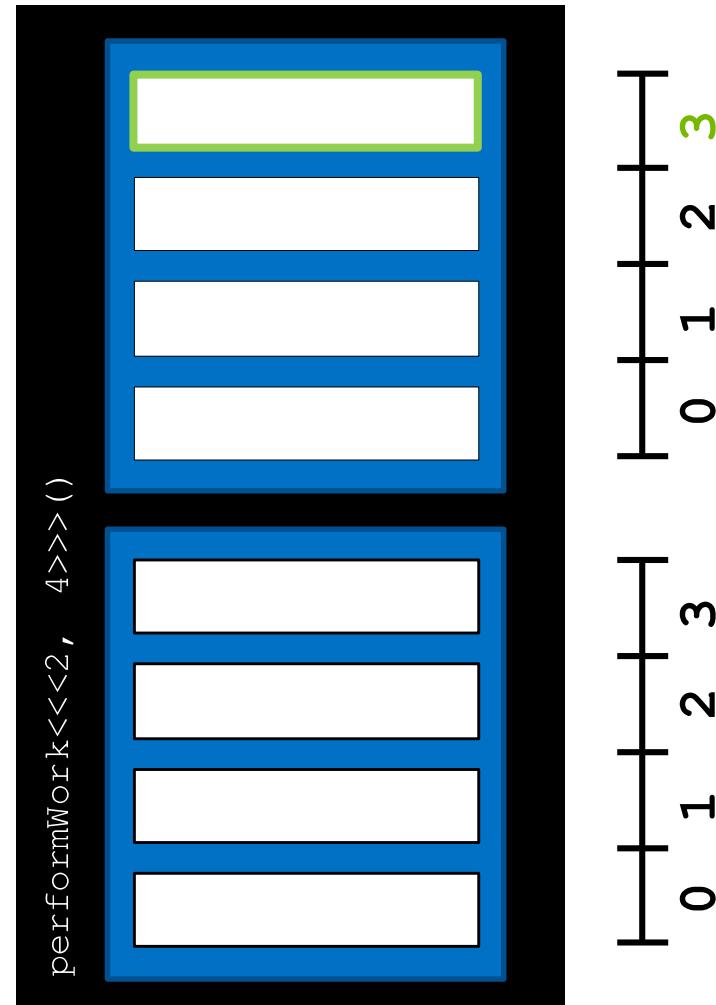
Inside a kernel **threadIdx**
describes the index of the thread
a block. In this case 1



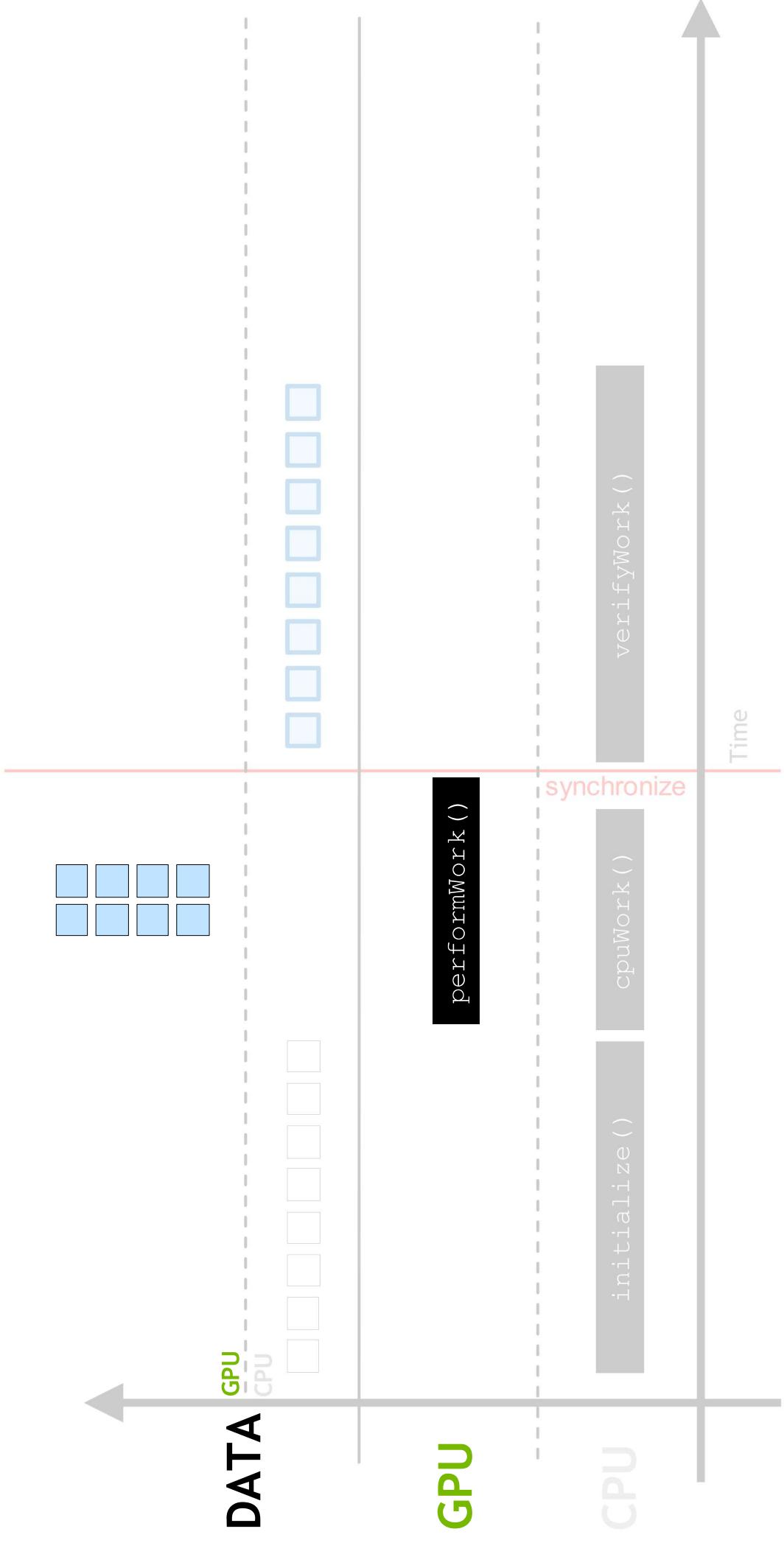
Inside a kernel **threadIdx**
describes the index of the thread
a block. In this case 2



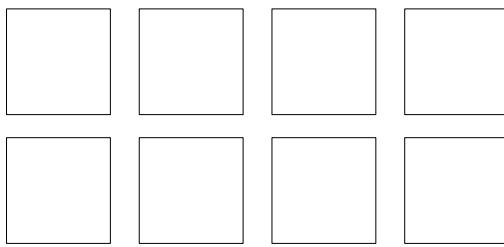
Inside a kernel **threadIdx**
describes the index of the thread
a block. In this case 3



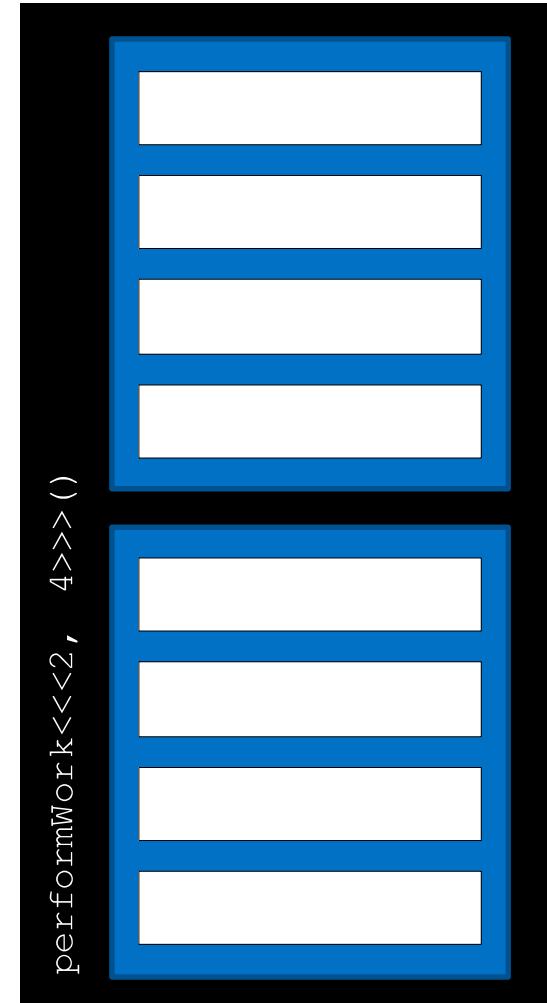
Coordinating Parallel Threads



Assume data is in a 0 indexed vector



GPU DATA

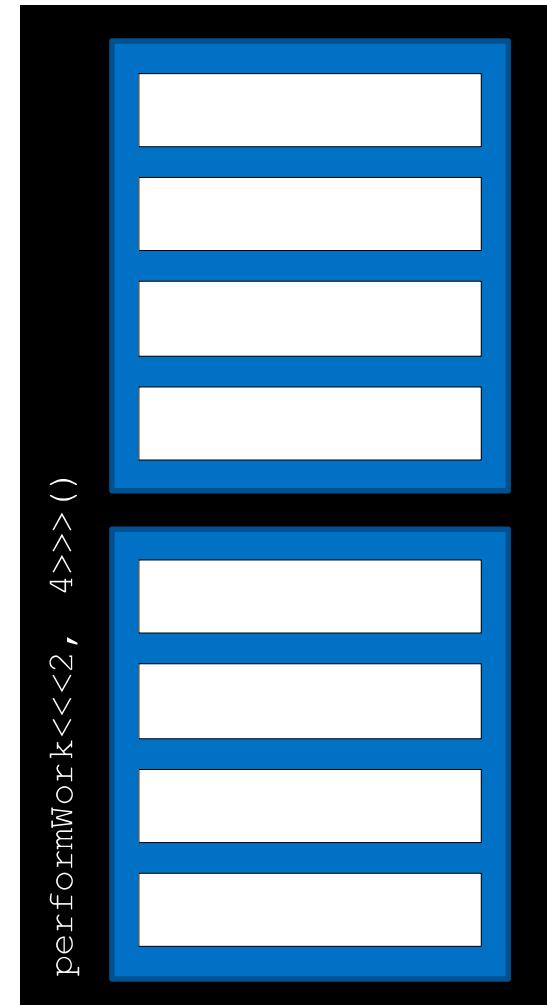


GPU

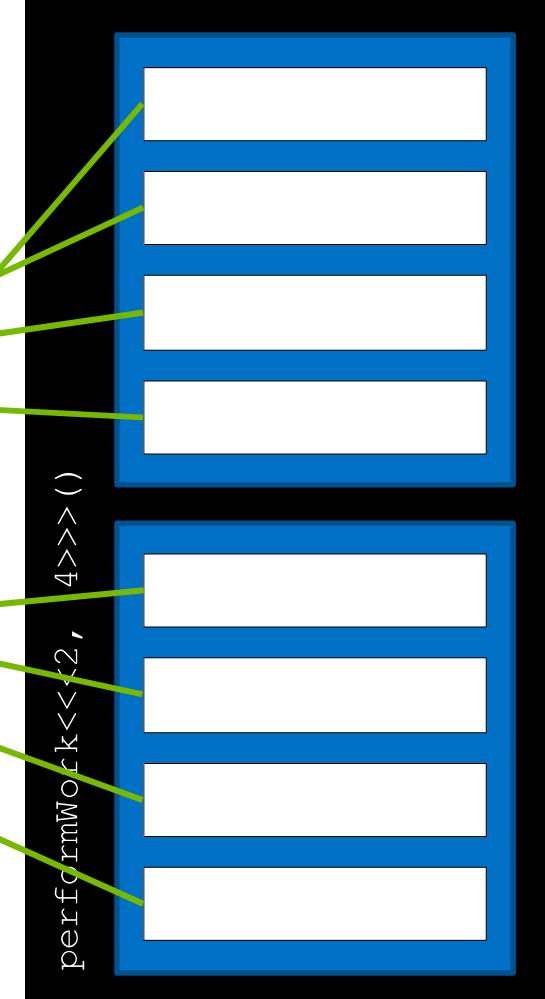
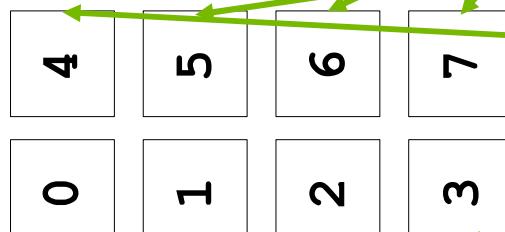
Assume data is in a 0 indexed v

0	4
1	5
2	6
3	7

GPU
DATA



Somehow, each thread must mapped to work on an element vector



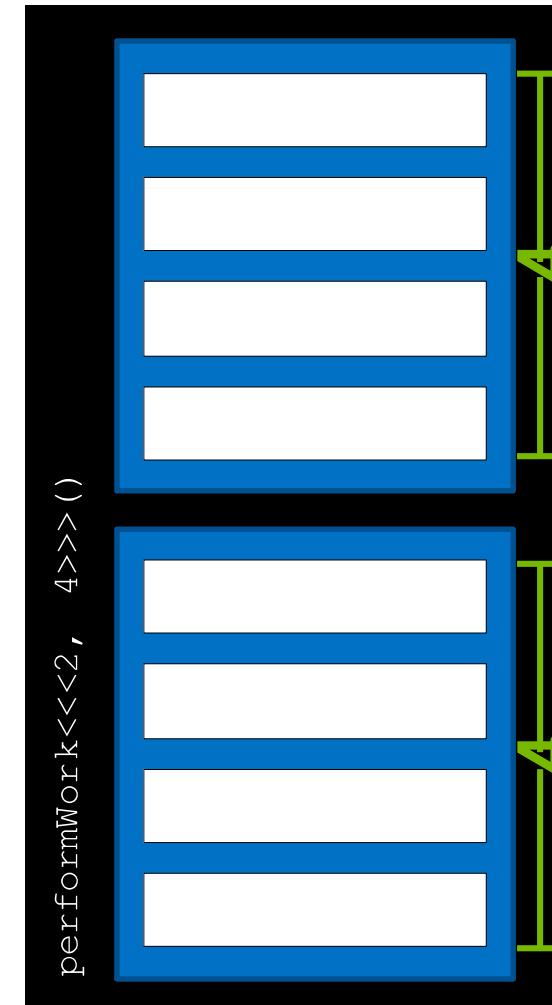
GPU
DATA

GPU

Recall that each thread has access to the size of its block via `blockD`

0	4
1	5
2	6
3	7

GPU
DATA

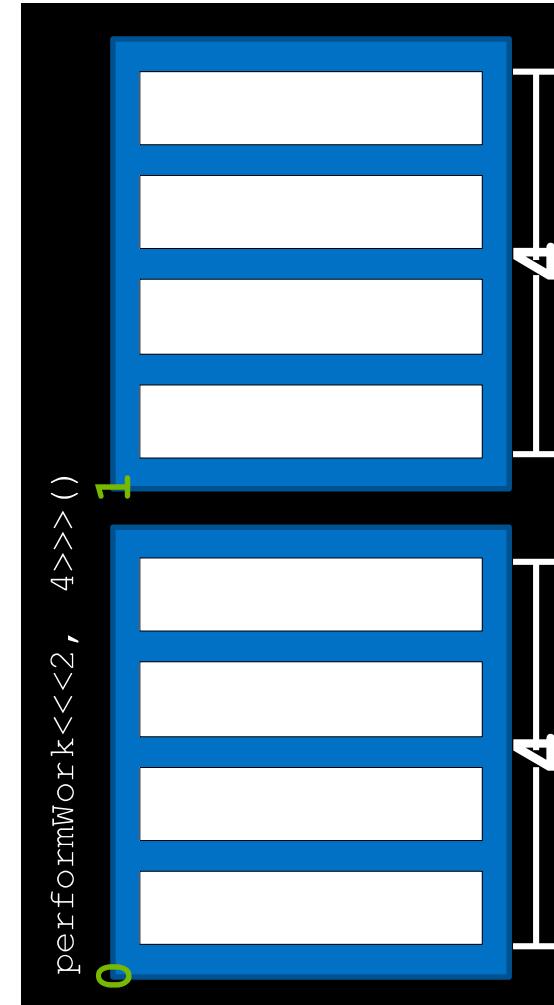


GPU

GPU DATA

0	4
1	5
2	6
3	7

...and the index of its block with
grid via `blockIdx.x`

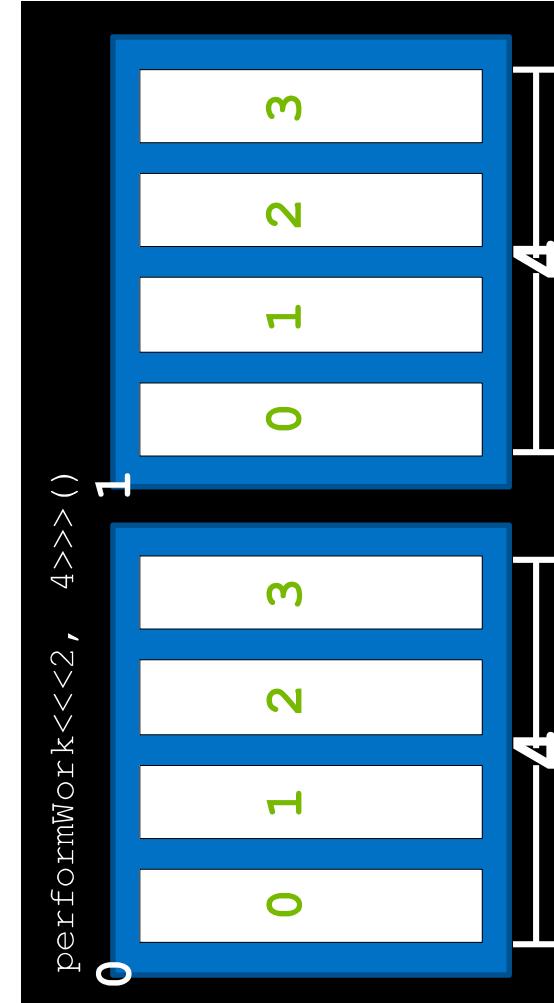


GPU

GPU DATA

0	4
1	5
2	6
3	7

...and its own index within its block
`threadIdx.x`

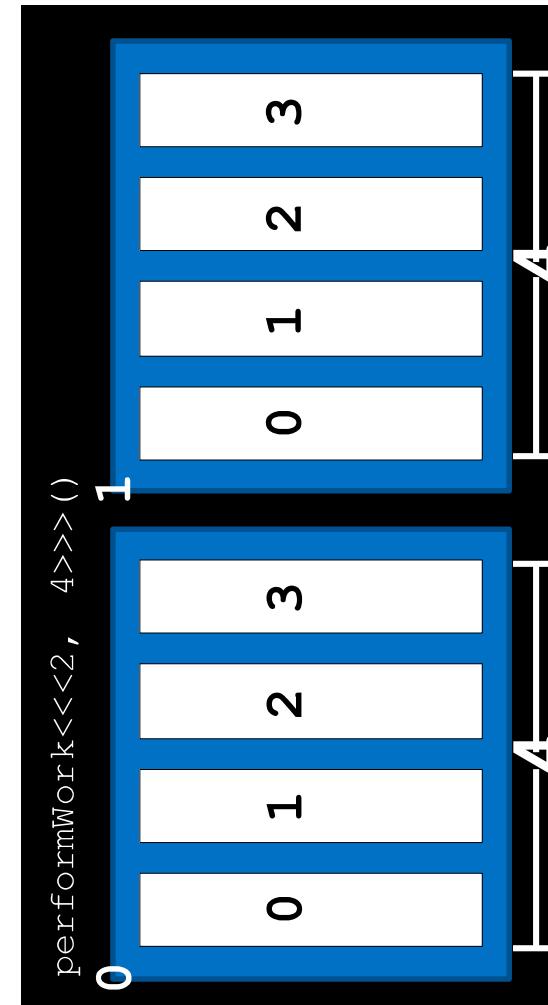


GPU

Using these variables, the formula `threadIdx.x + blockIdx.x * blockDim.x` will map each thread one element in the vector.

0	4
1	5
2	6
3	7

GPU
DATA

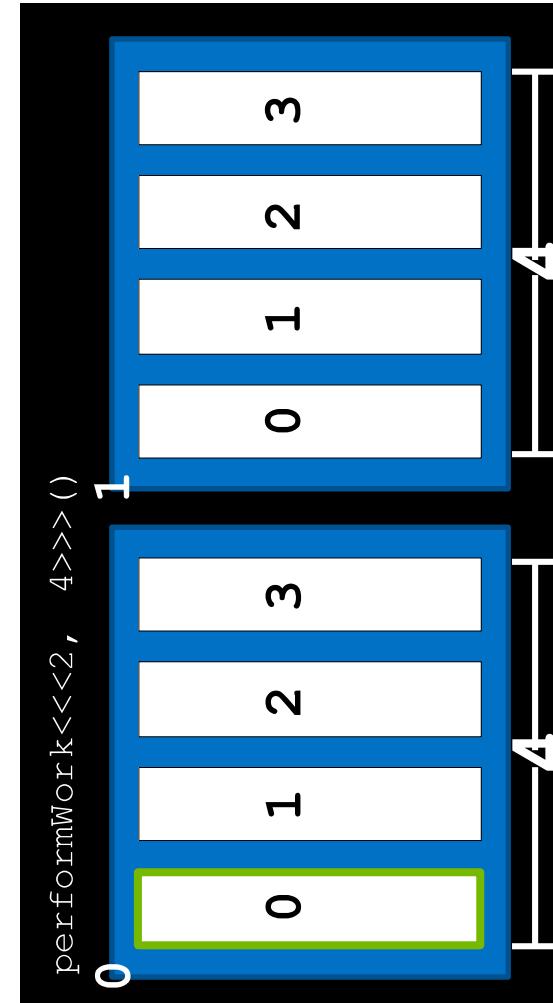


GPU DATA

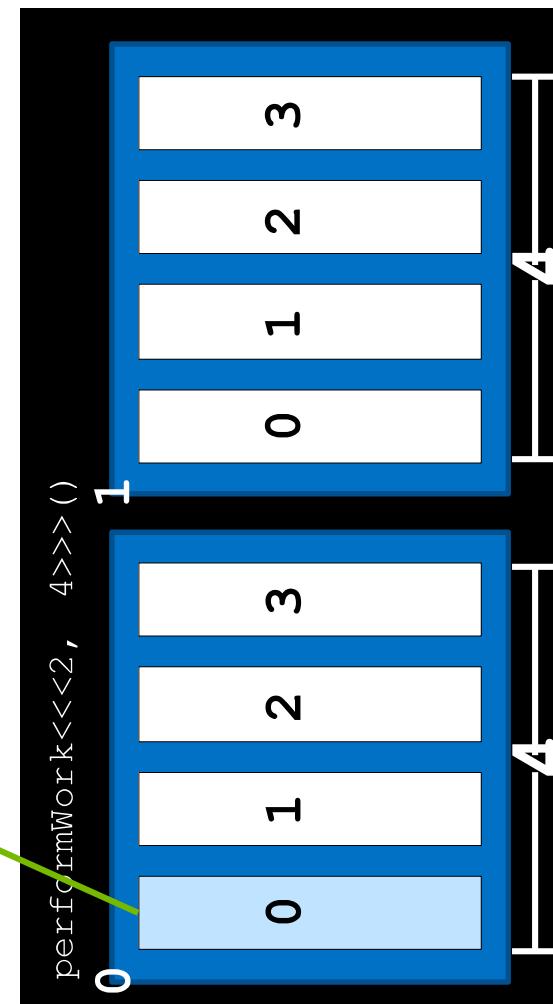
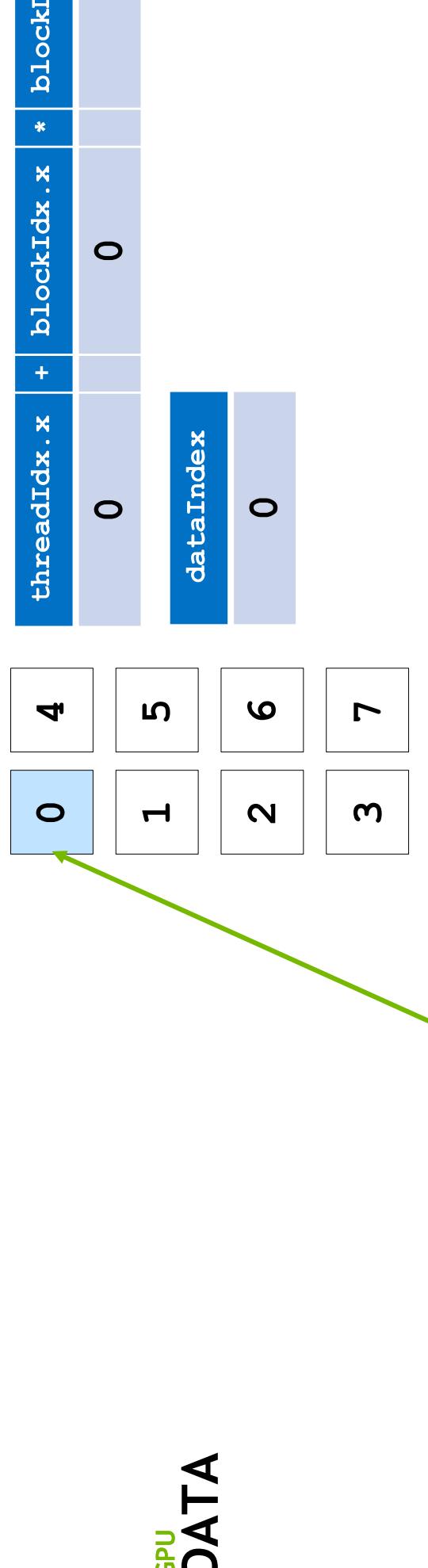
0	4	5	6	7
1	2	3		

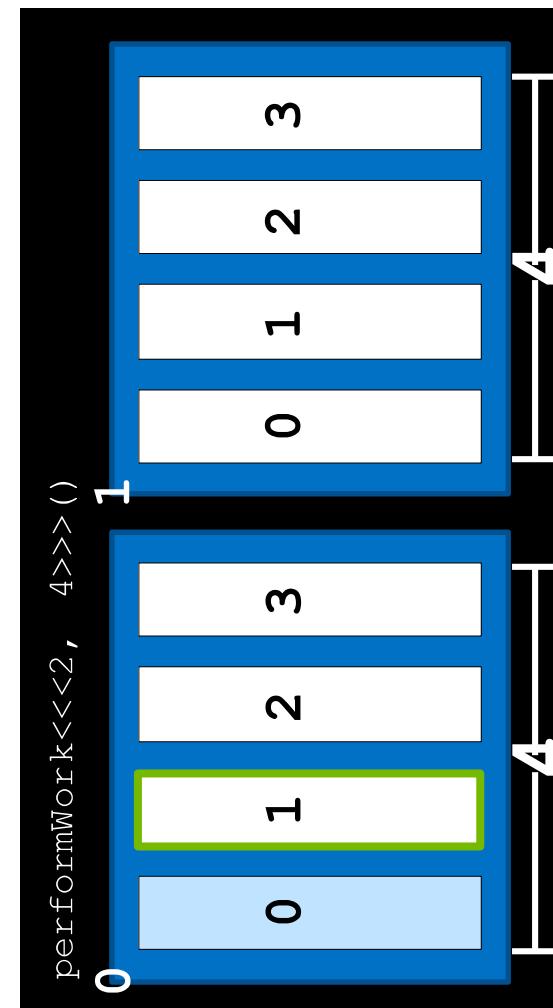
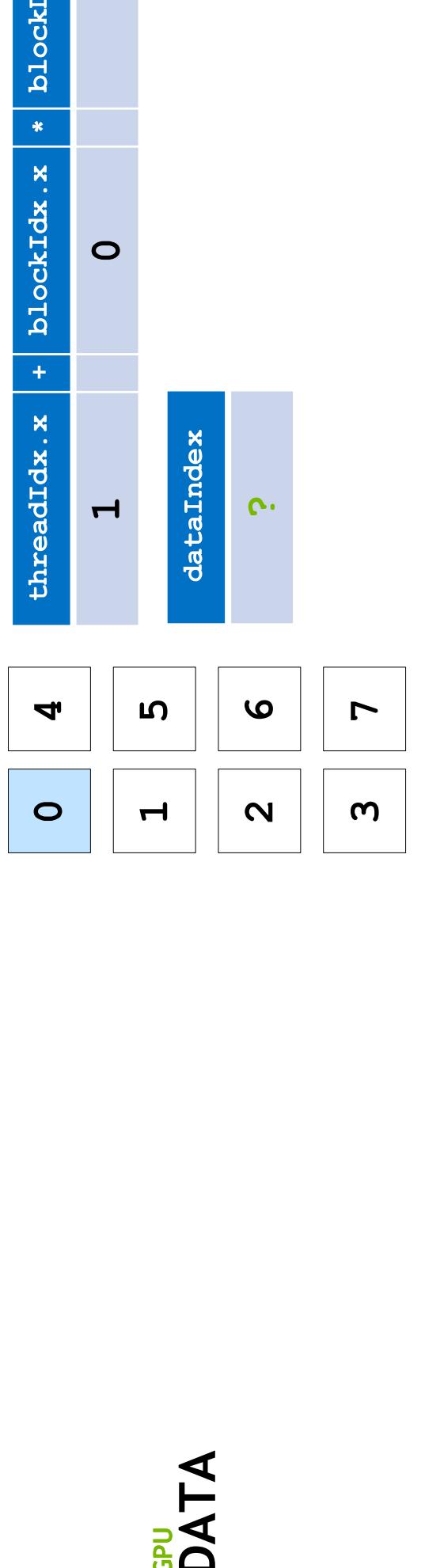
threadIdx.x	*	blockIdx.x	*	blockIdx.y
0		0		

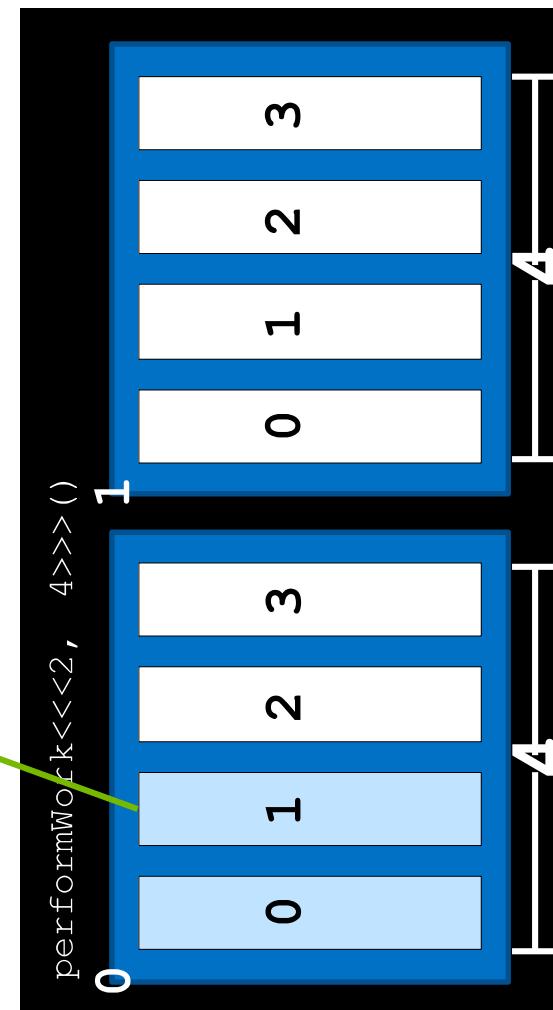
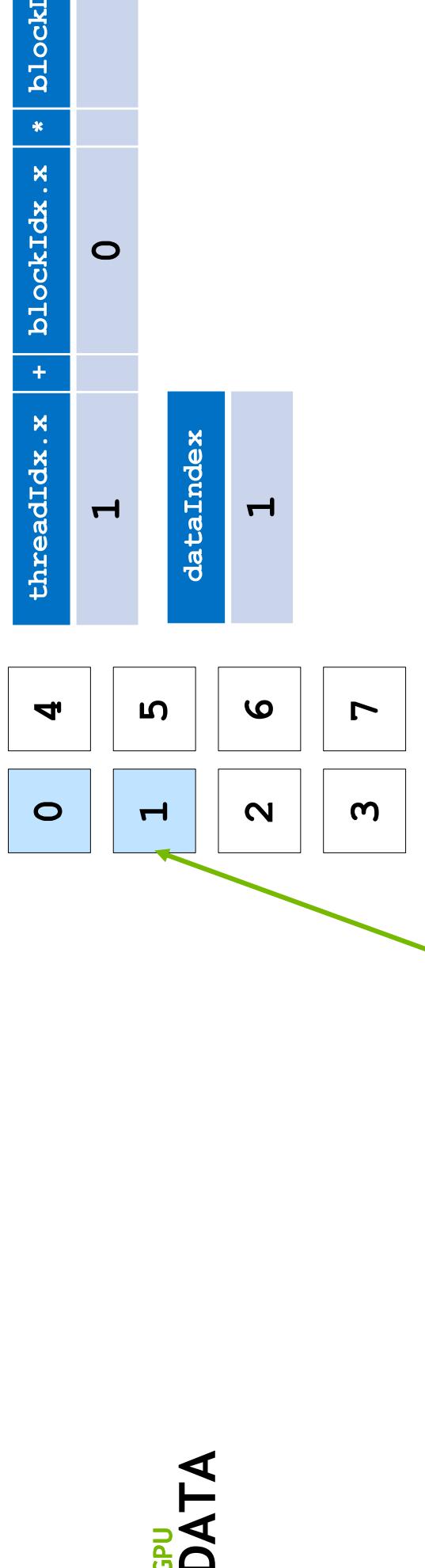
dataIndex
?



GPU







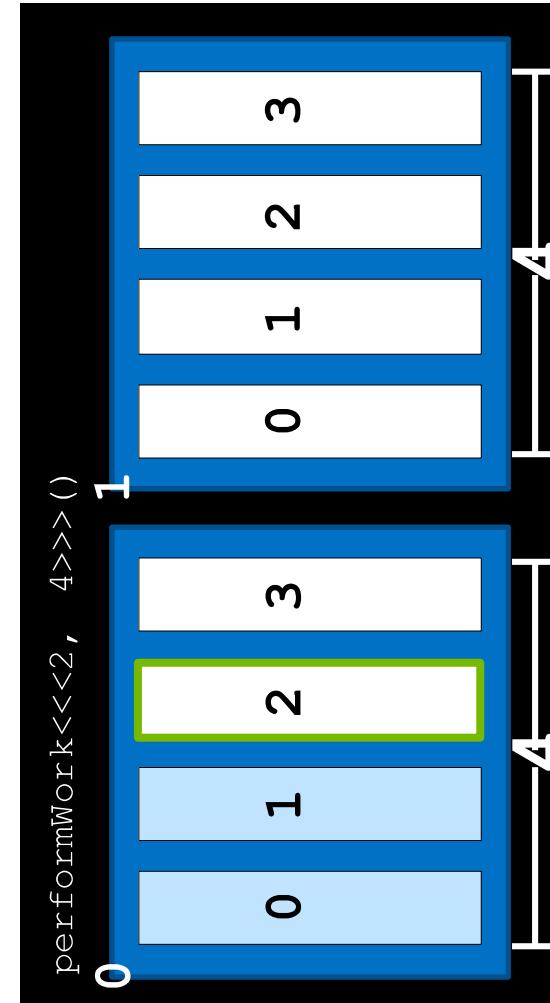
GPU DATA

0	4	5	6	7
1	2	3		

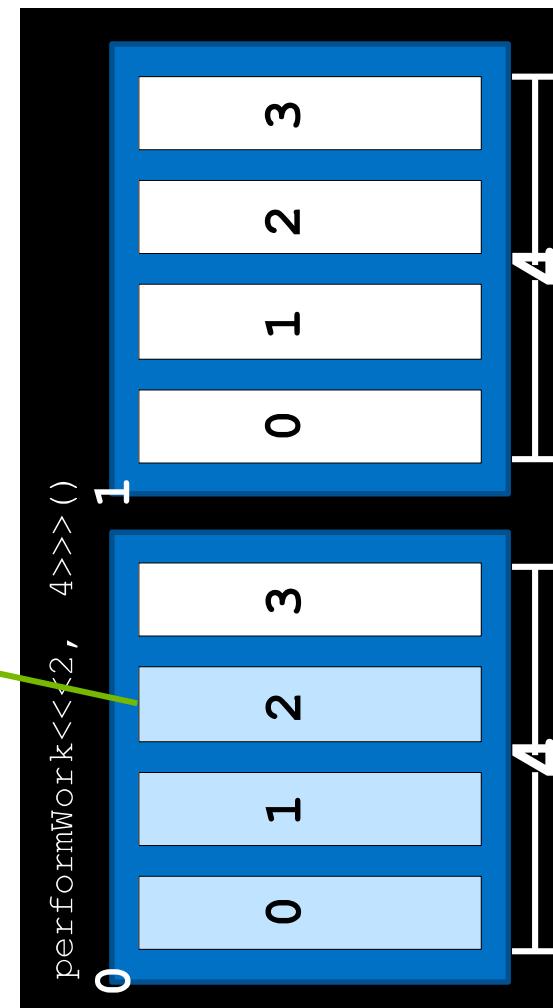
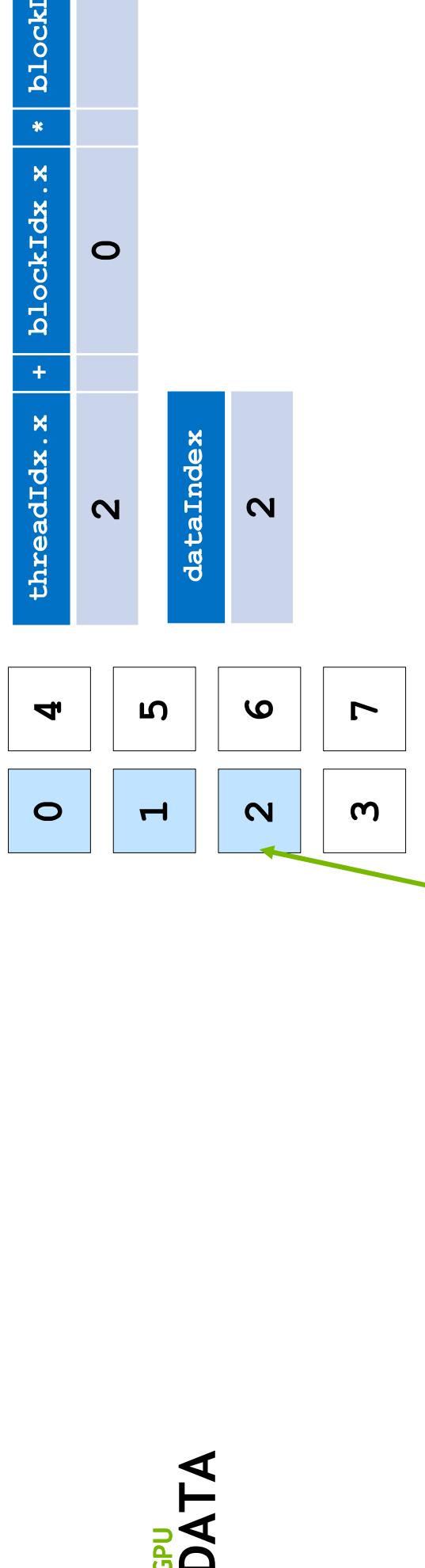
threadIdx.x	+	blockIdx.x	*	blockIdx.y
0	2	0		

dataIndex

?



GPU



GPU DATA

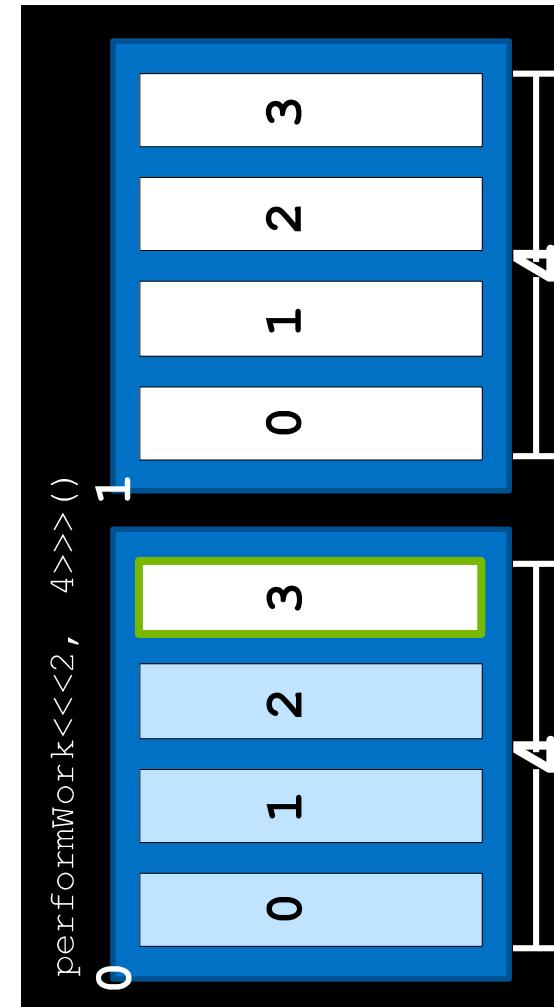
0	4	5	6	7	3
1	2	3			

threadIdx.x	+	blockIdx.x	*	blockIdx.y
3		0		

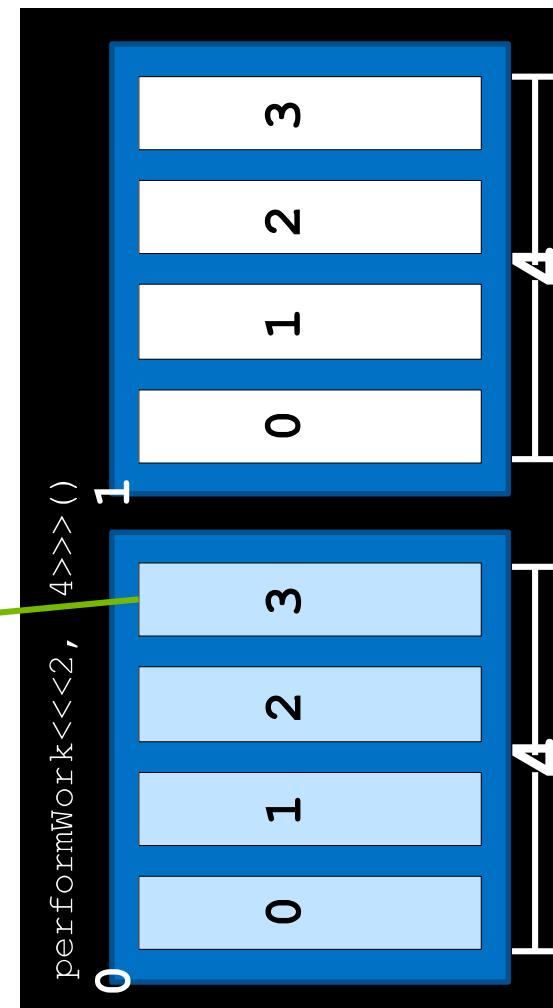
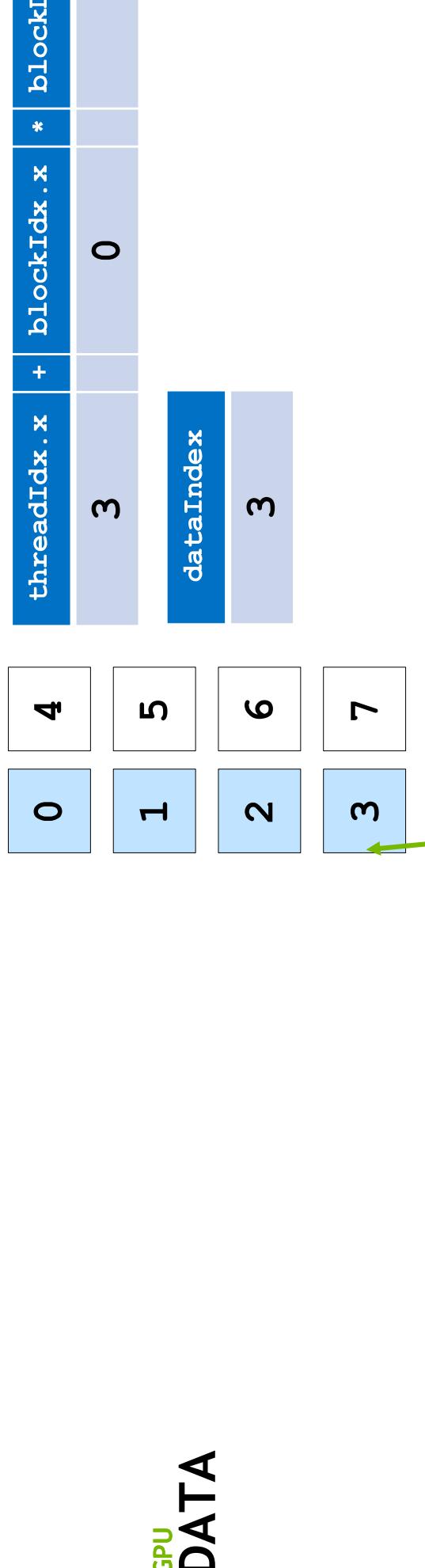
dataIndex

?

performWork<<<2, 4>>>()



GPU

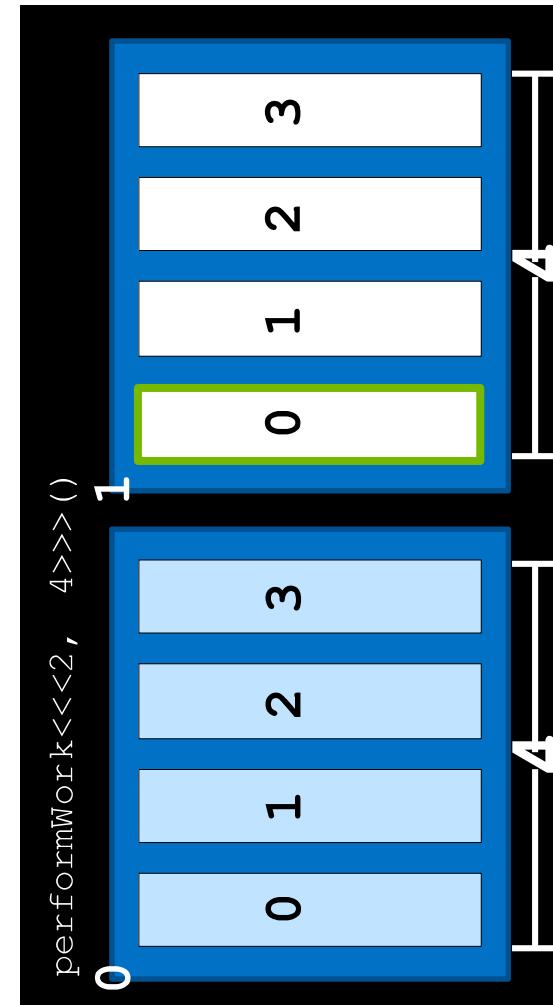


GPU DATA

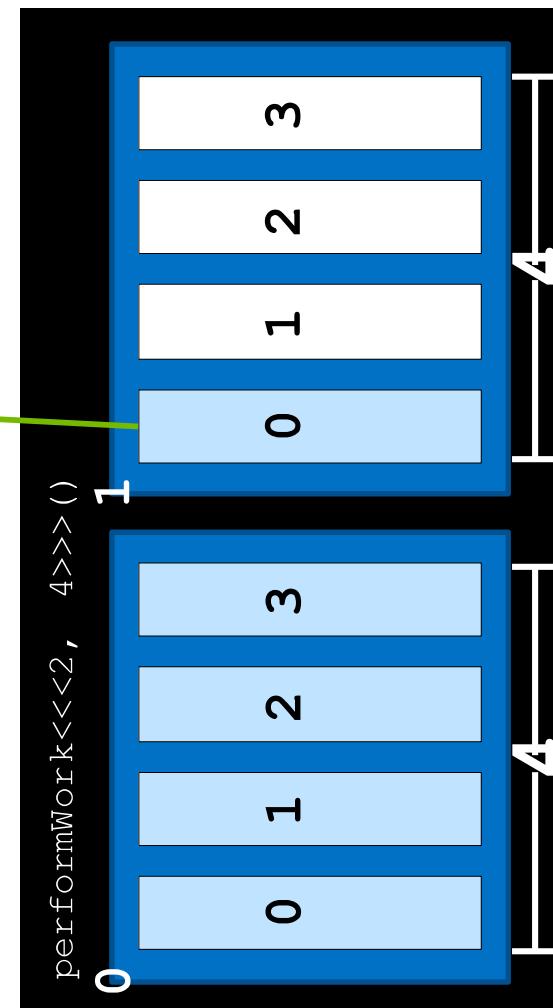
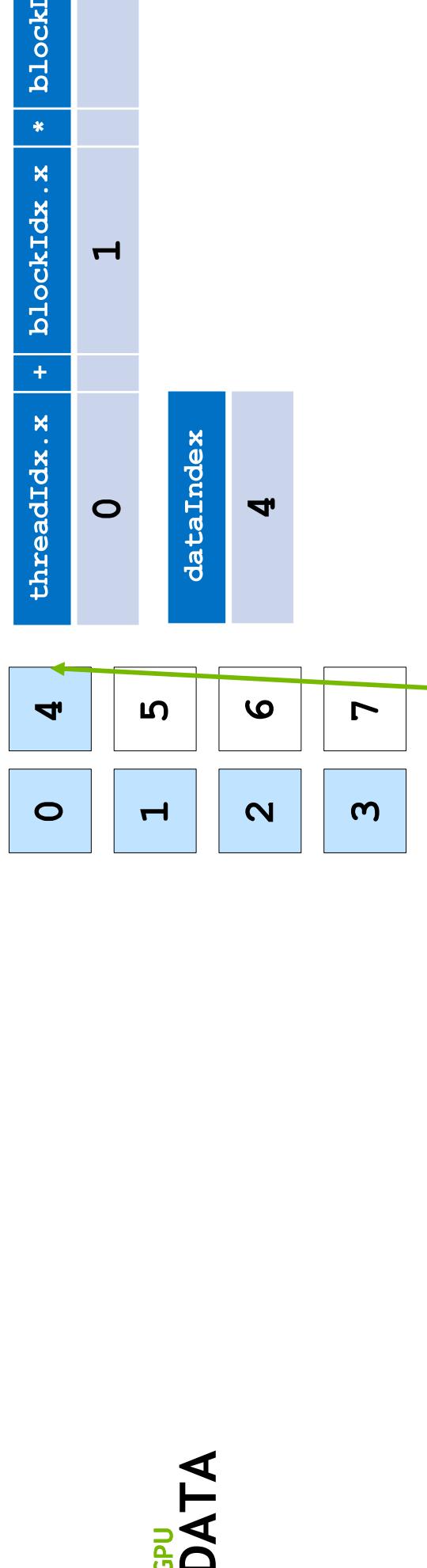
0	4	5	6	7
1	2	3		

threadIdx.x	*	blockIdx.x	*	blockIdx.y
0		1		

dataIndex
?



GPU

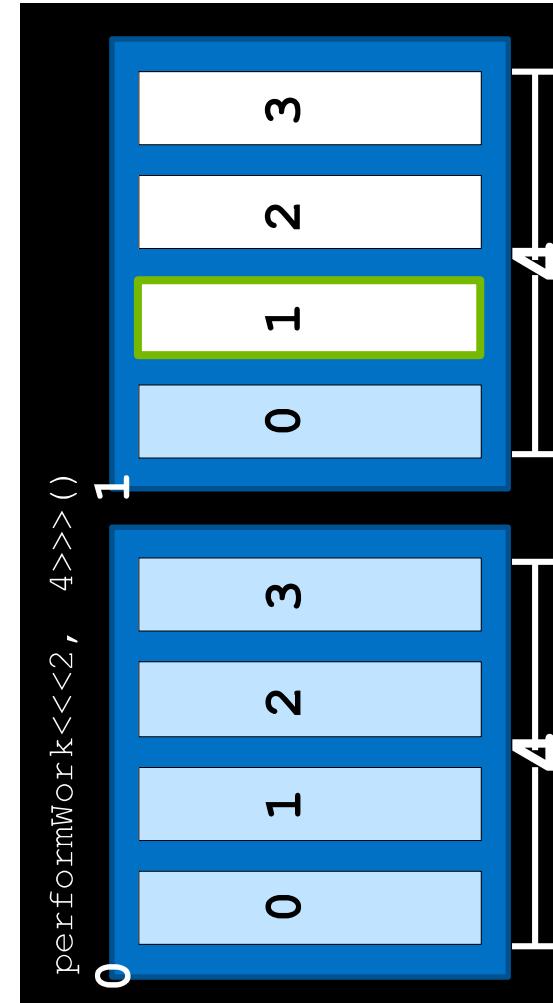


GPU DATA

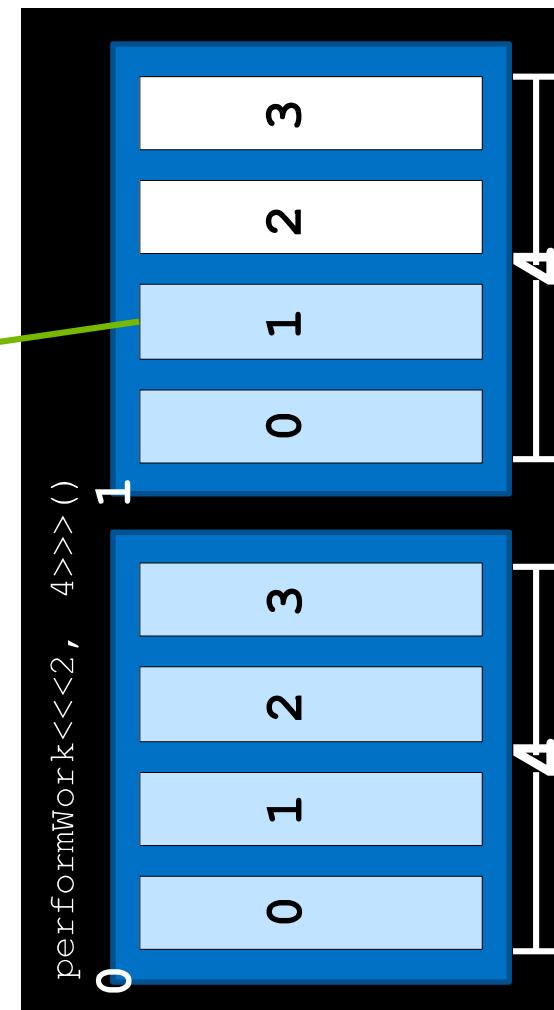
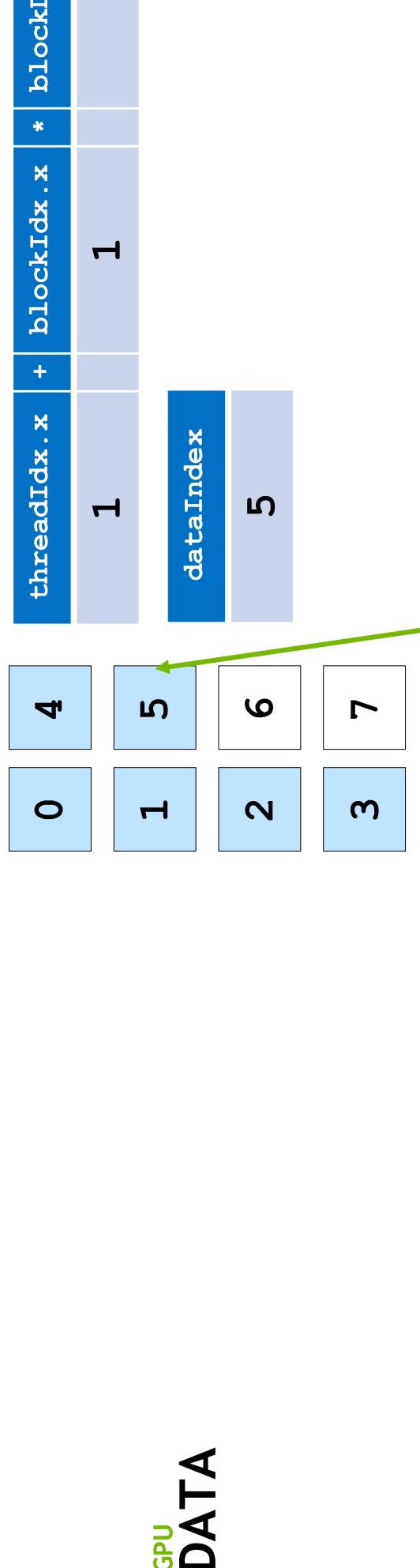
0	4	5	6	7
1	2	3		

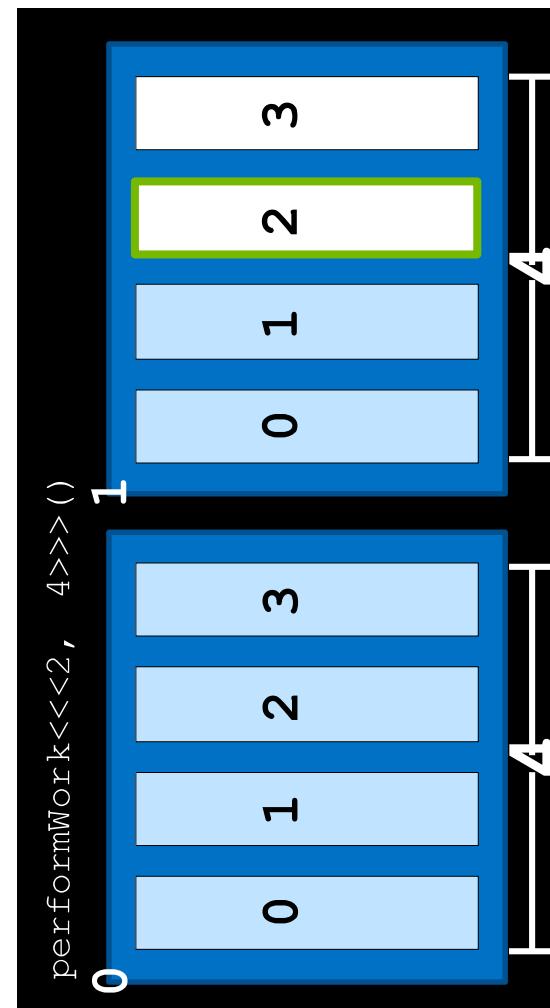
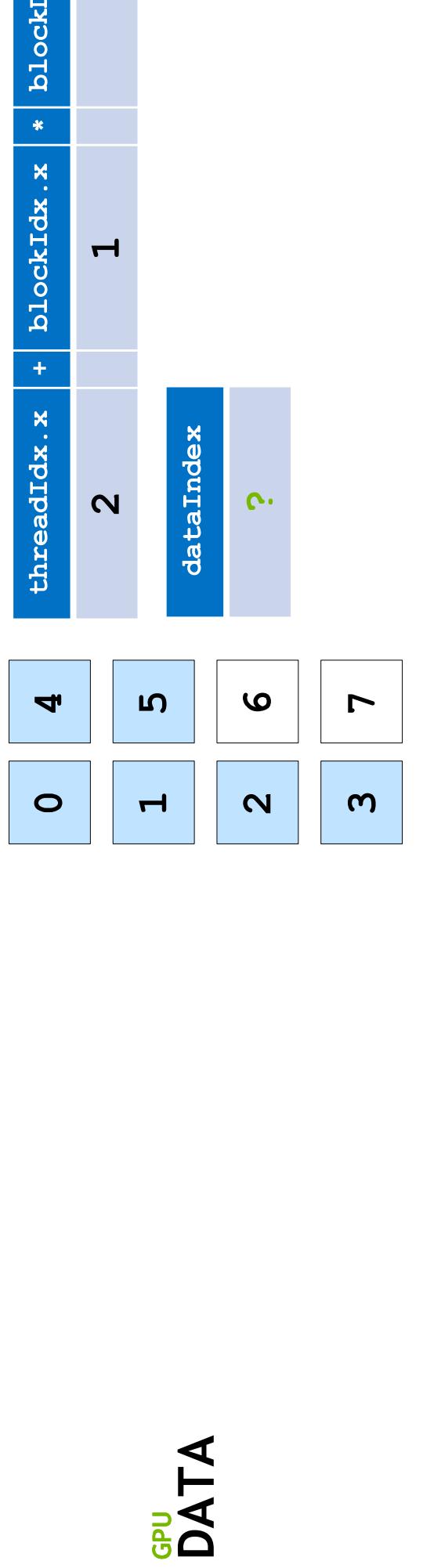
threadIdx.x	*	blockIdx.x	*	blockIdx.y
1		1		

dataIndex
?

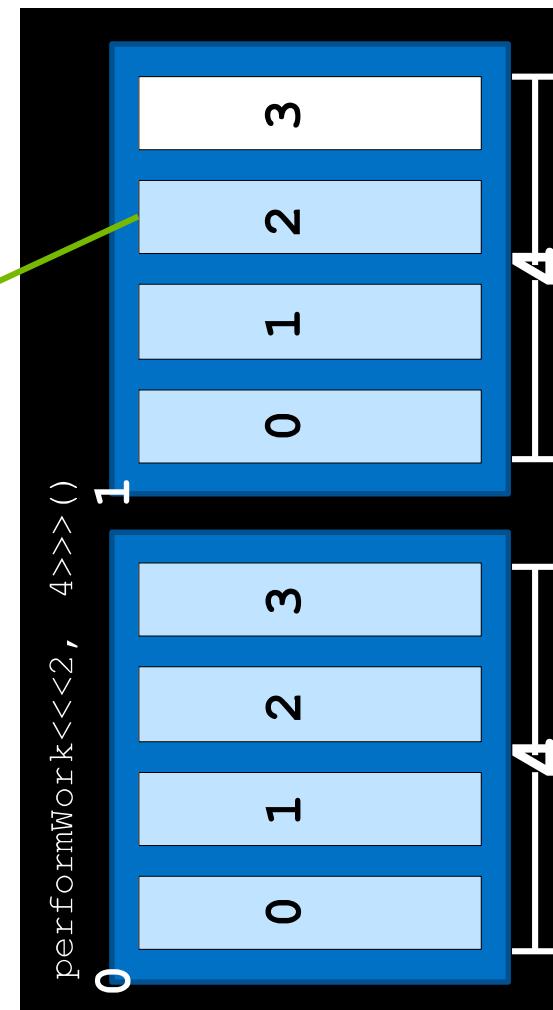
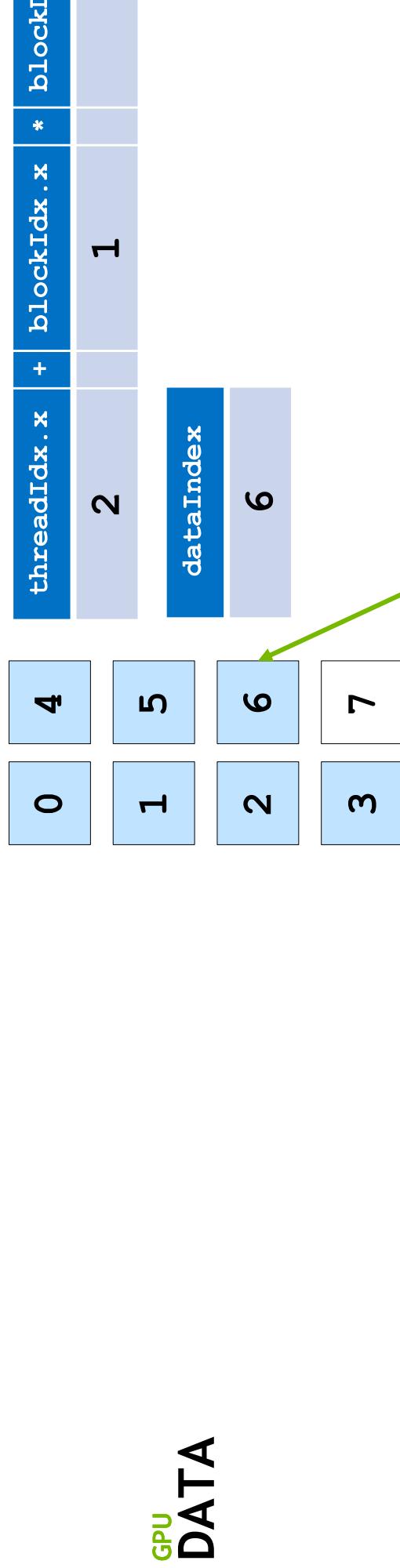


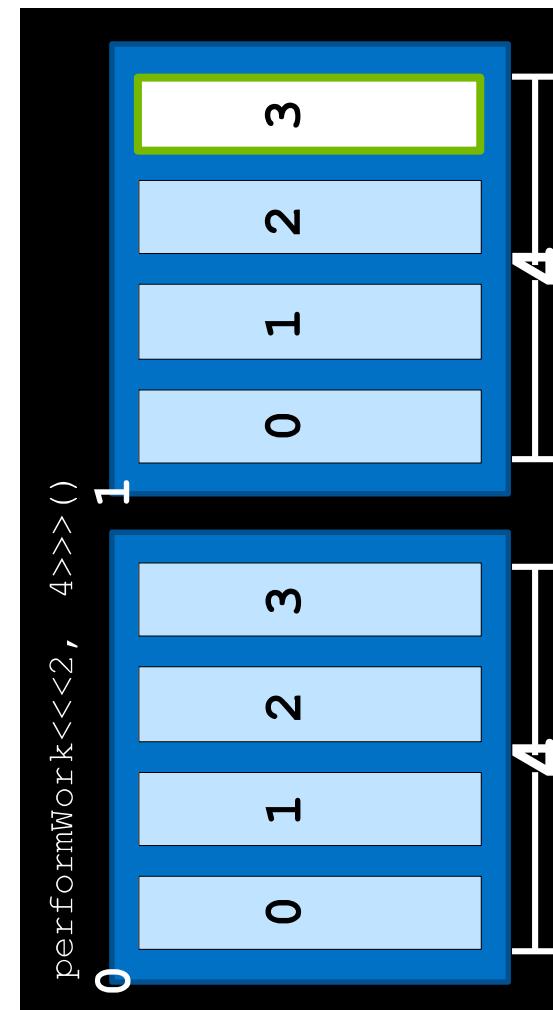
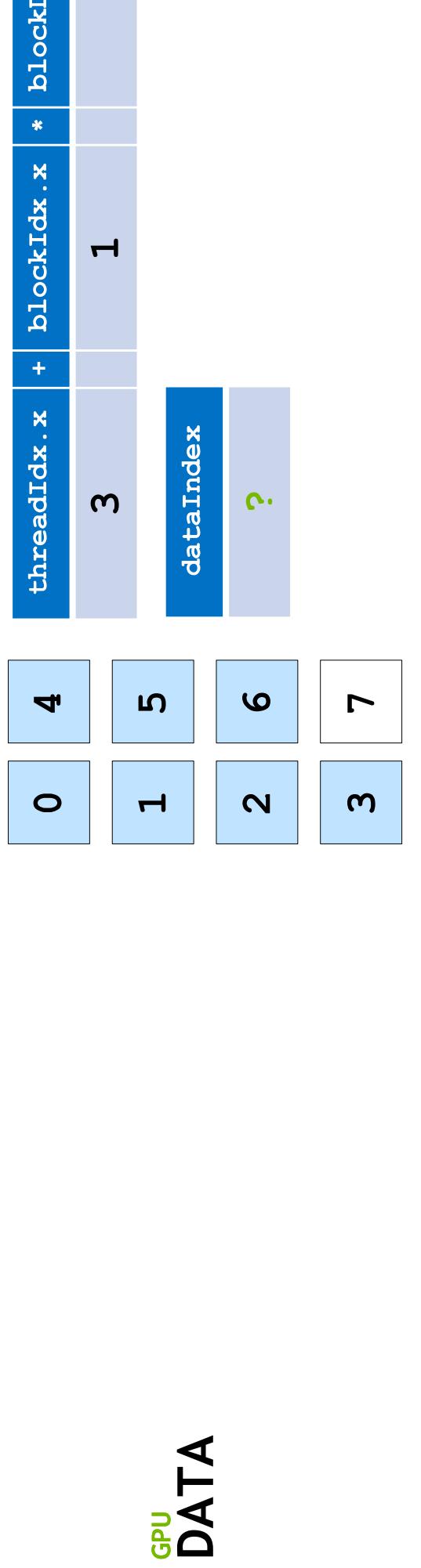
GPU

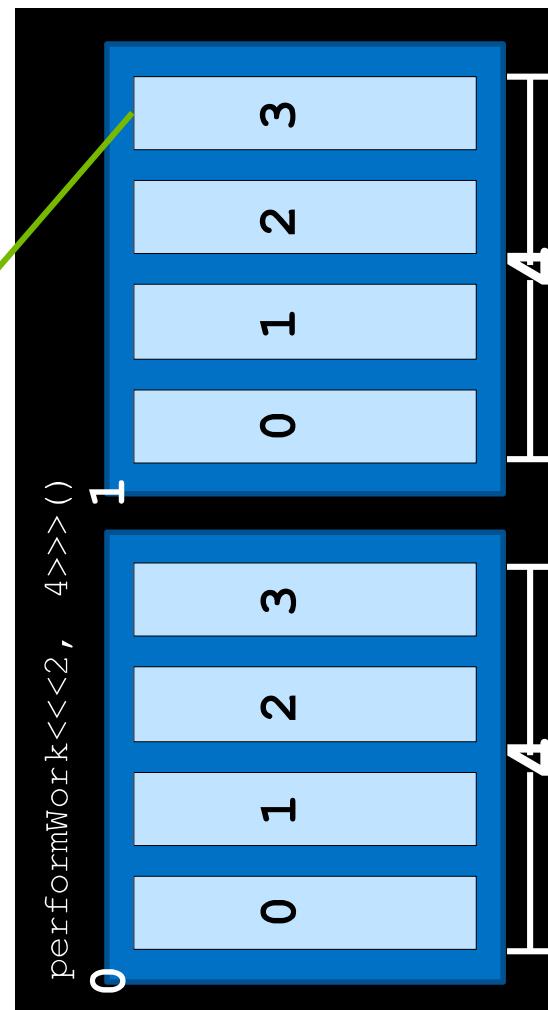
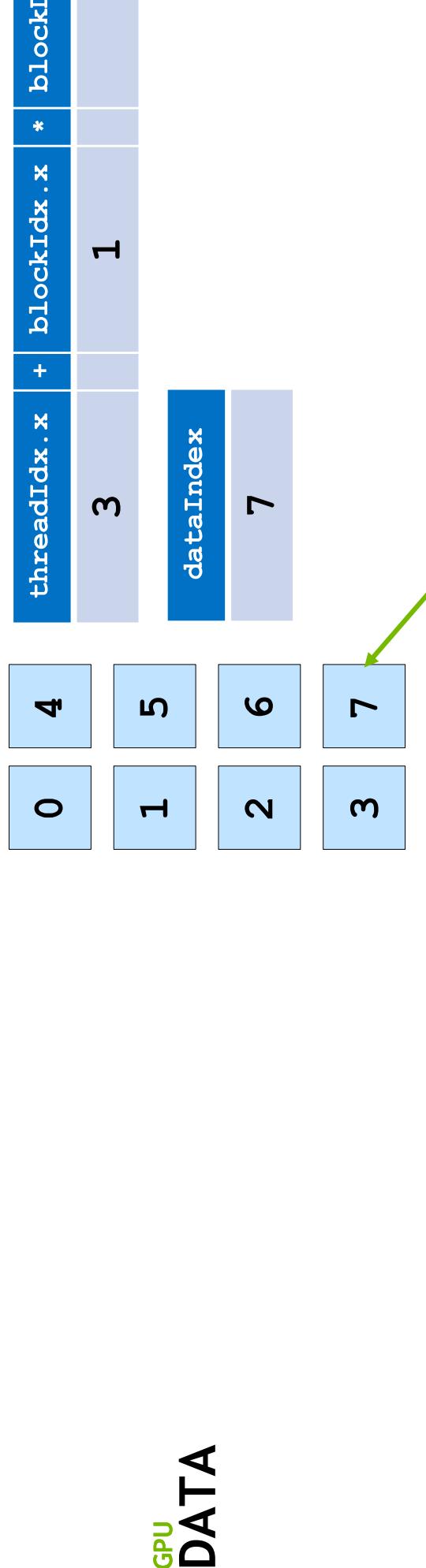




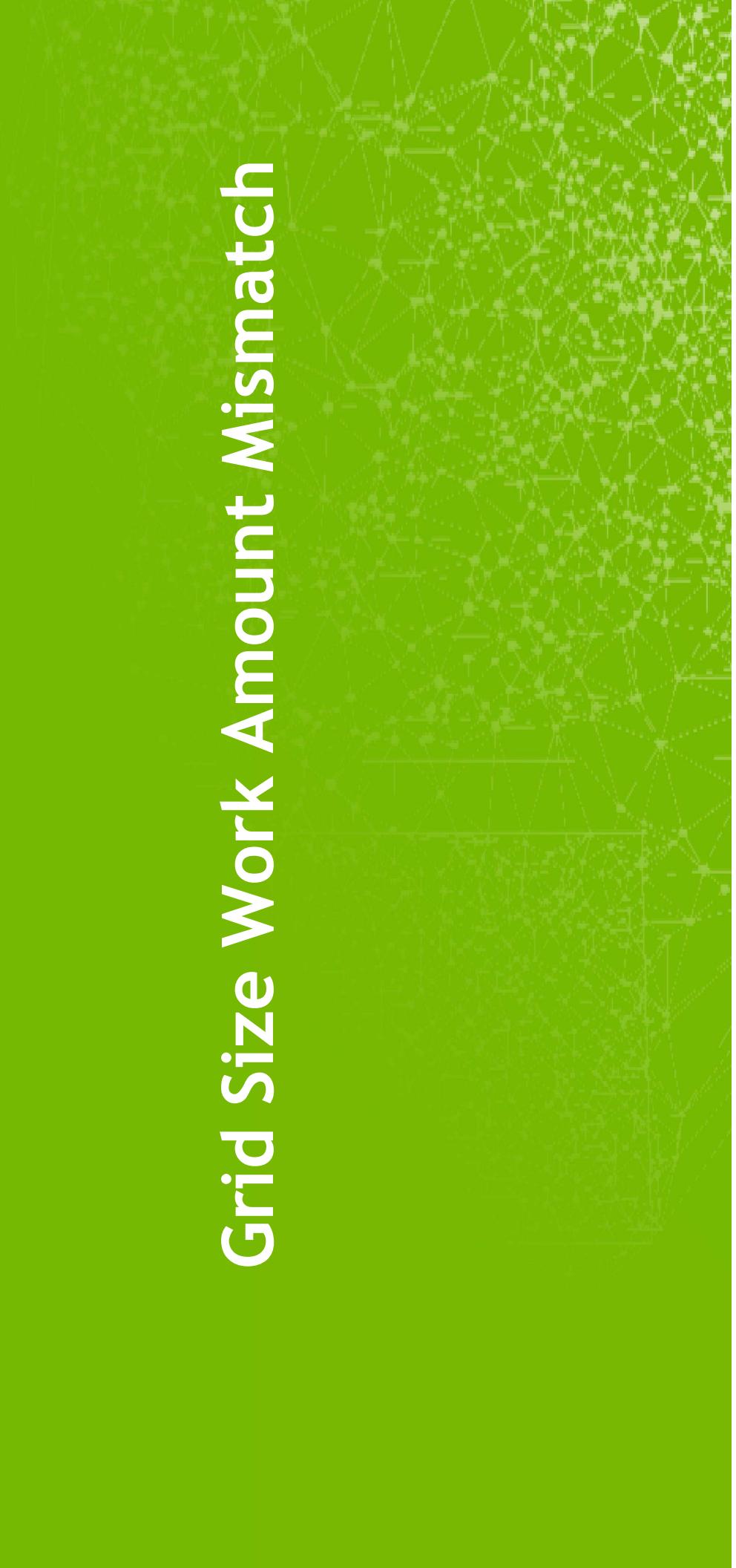
GPU



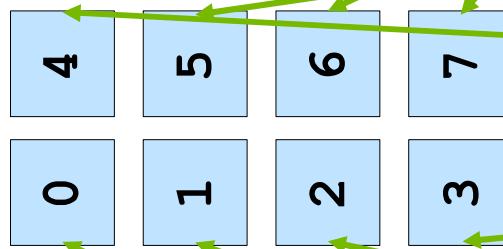




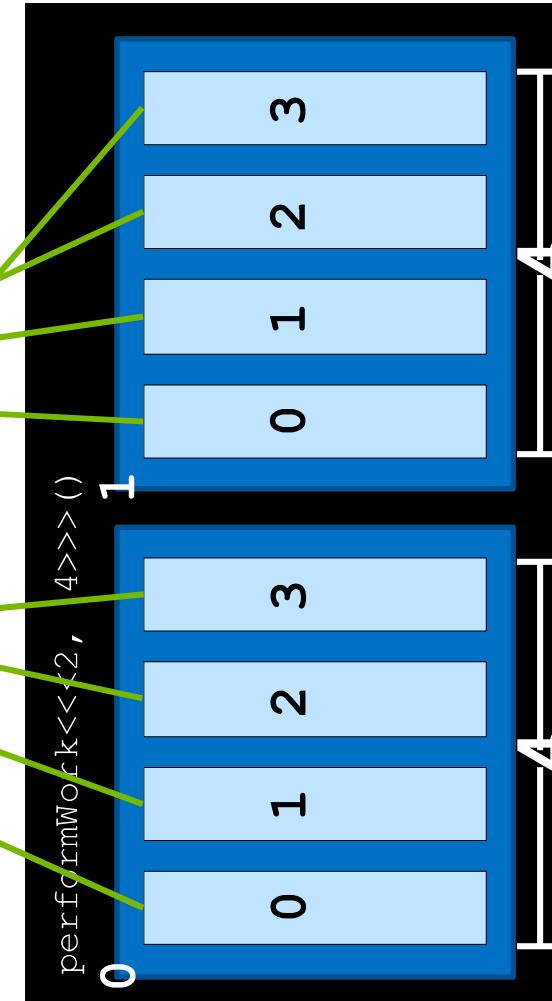
Grid Size Work Amount Mismatch



In previous scenarios, the number of threads in the grid matched the number of elements exactly

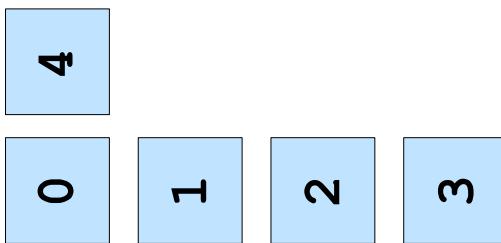


GPU DATA

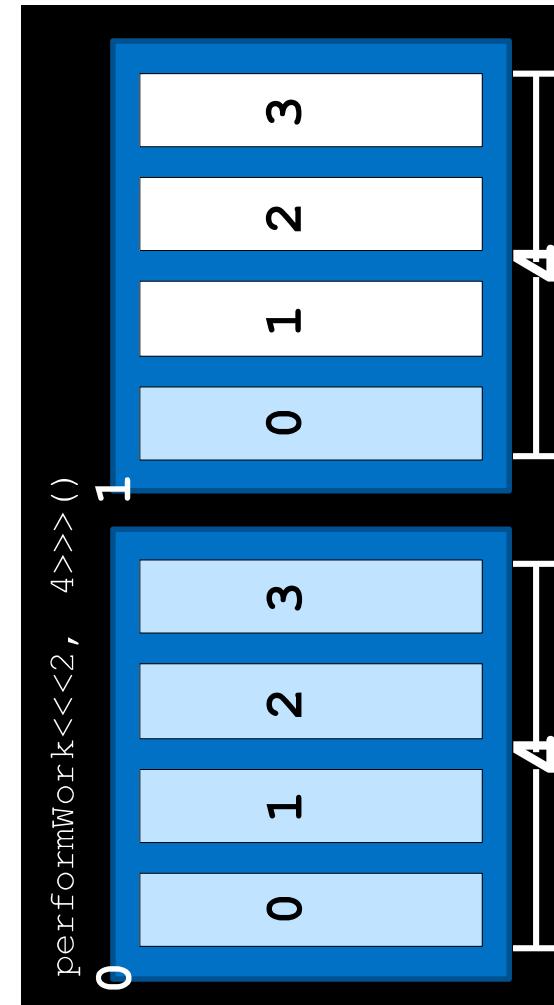


GPU

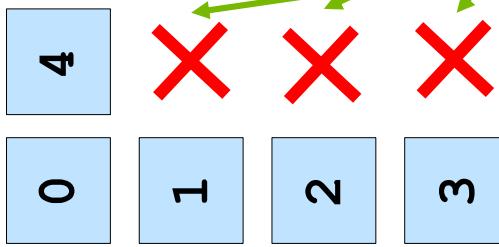
GPU DATA



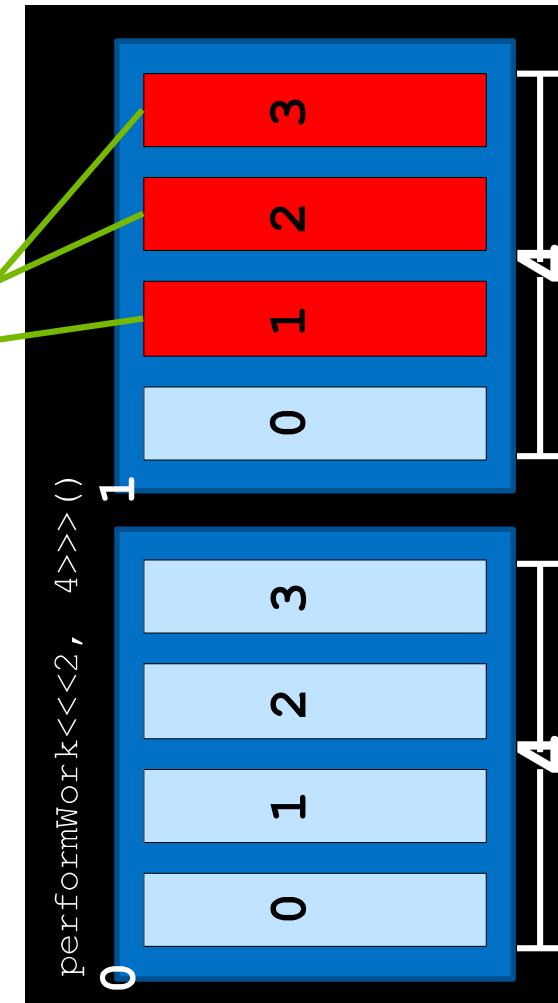
What if there are more threads
work to be done?



Attempting to access non-existent elements can result in a runtime error.

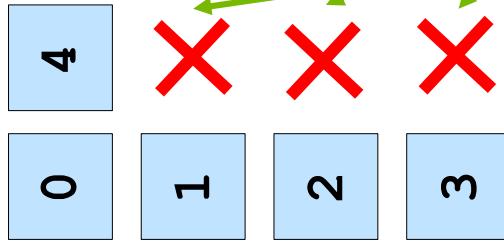


GPU
DATA

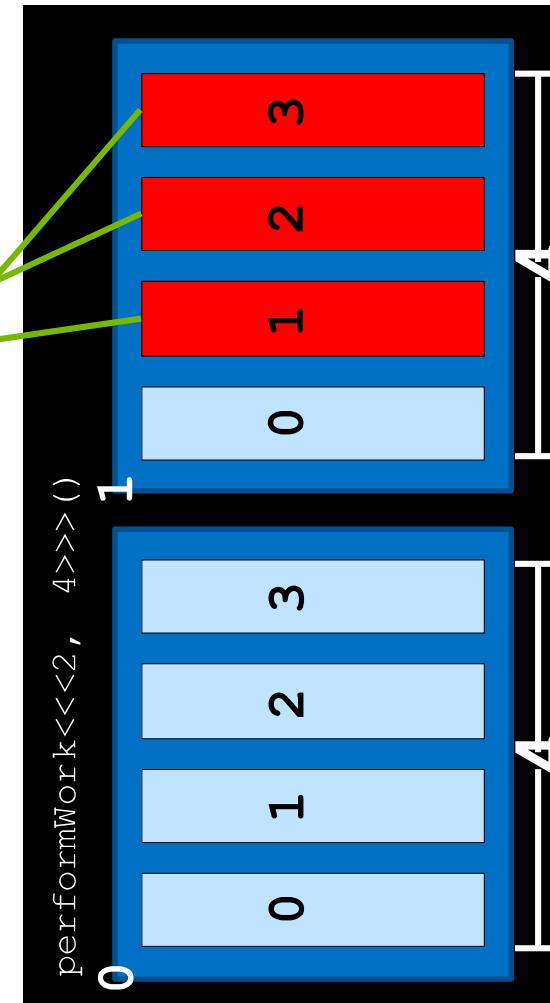


GPU

Code must check that the data calculated by `threadIdx.x * blockIdx.x * blockDim.x` than N, the number of data elements.



GPU
DATA

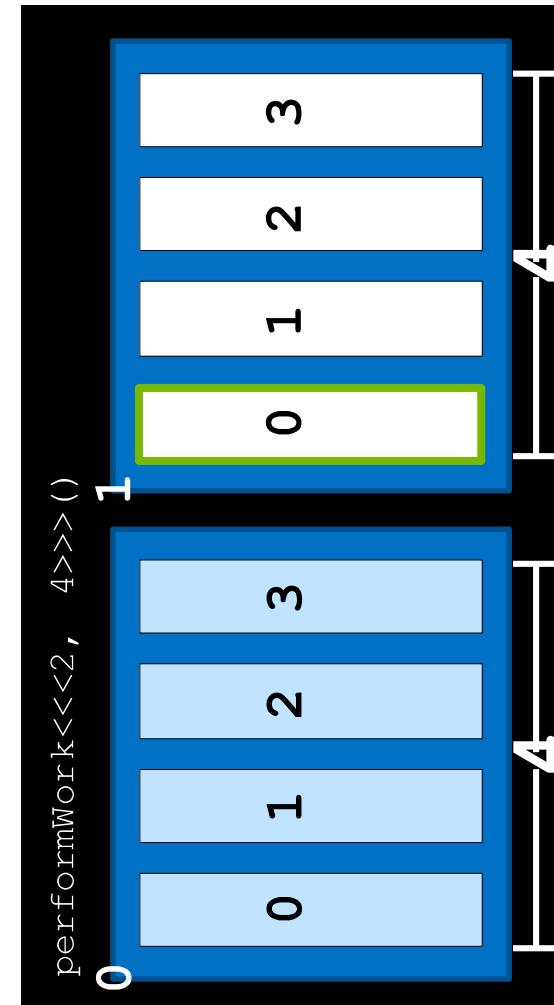


GPU

	threadIdx.x	+	blockIdx.x	*	blockIdx.y
	0		1		
	dataIndex <	N	=	Can	
	4		5		

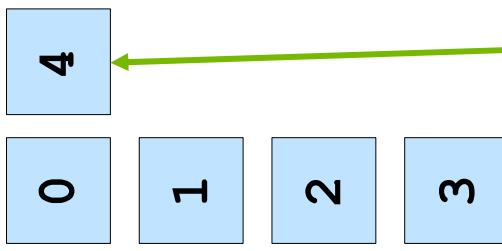
0 4
1 2
2 3

GPU
DATA

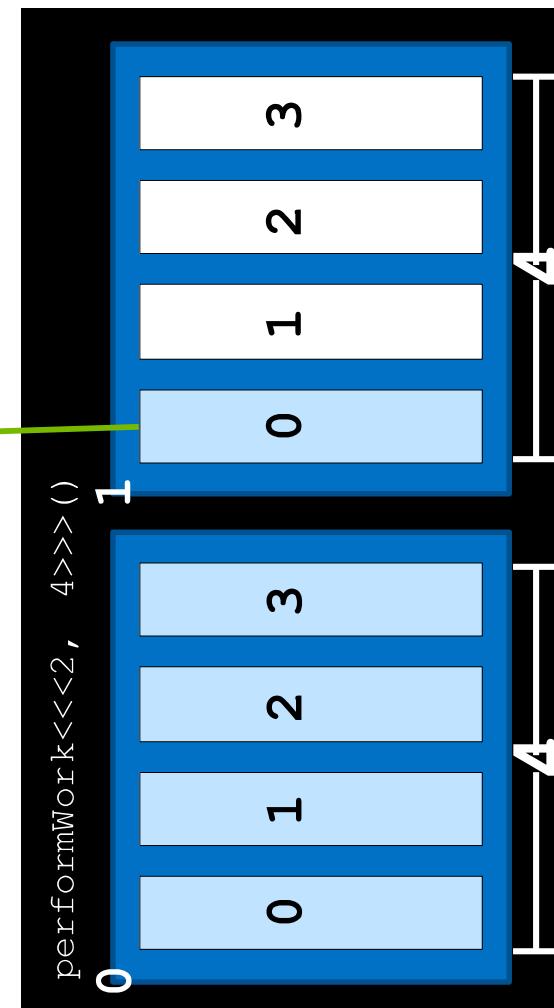


GPU

threadIdx.x	+	blockIdx.x	*	blockIdx.y
0		1		
dataIndex <	N	=		Can



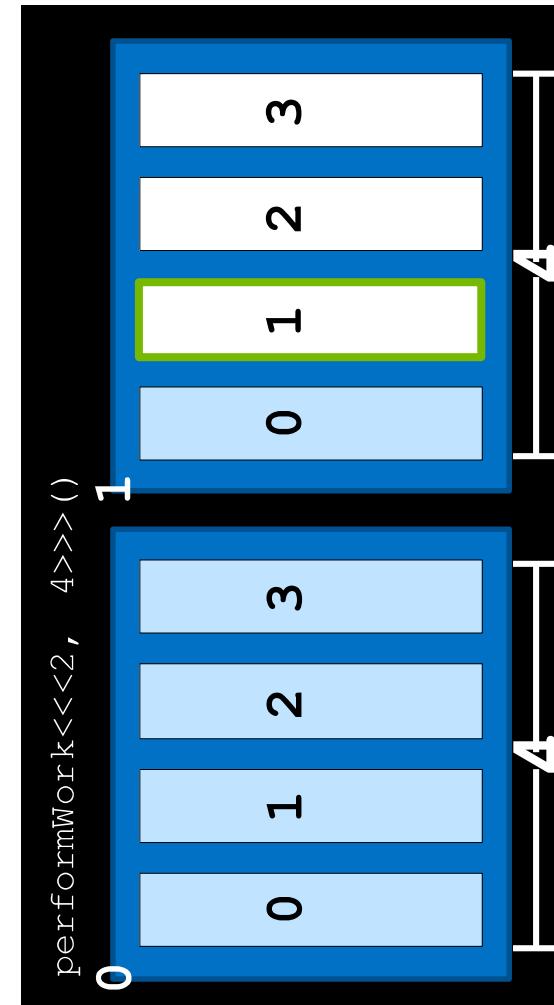
GPU
DATA



	threadIdx.x	+	blockIdx.x	*	blockIdx.y
	1		1		1
	5		5		5

0 4
1 2
2 3

GPU
DATA

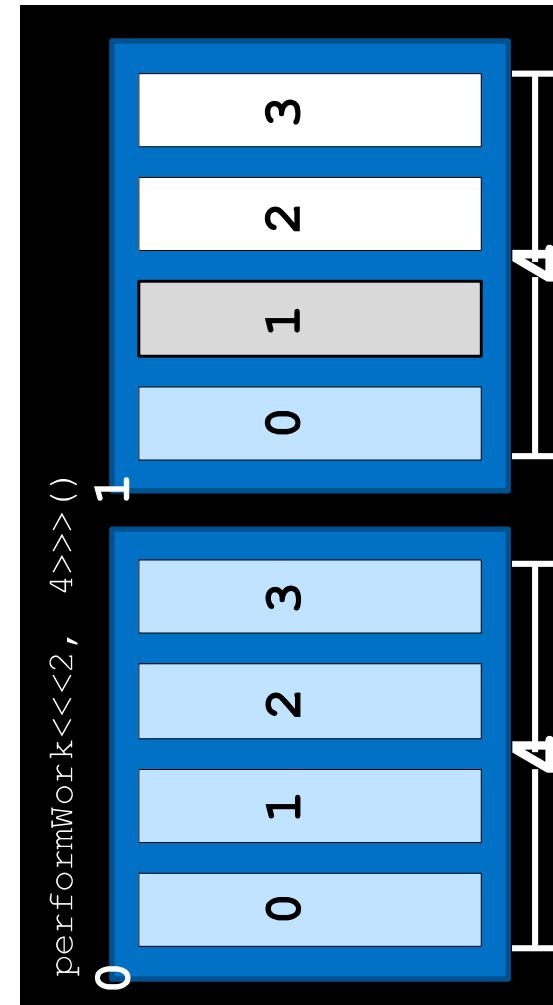


GPU

	threadIdx.x	+	blockIdx.x	*	blockIdx.y	*	blockIdx.z	
	1		1		1		1	
	dataIndex	<	N	=	Can		fa	
	5		5					

0	4
1	
2	
3	

GPU
DATA

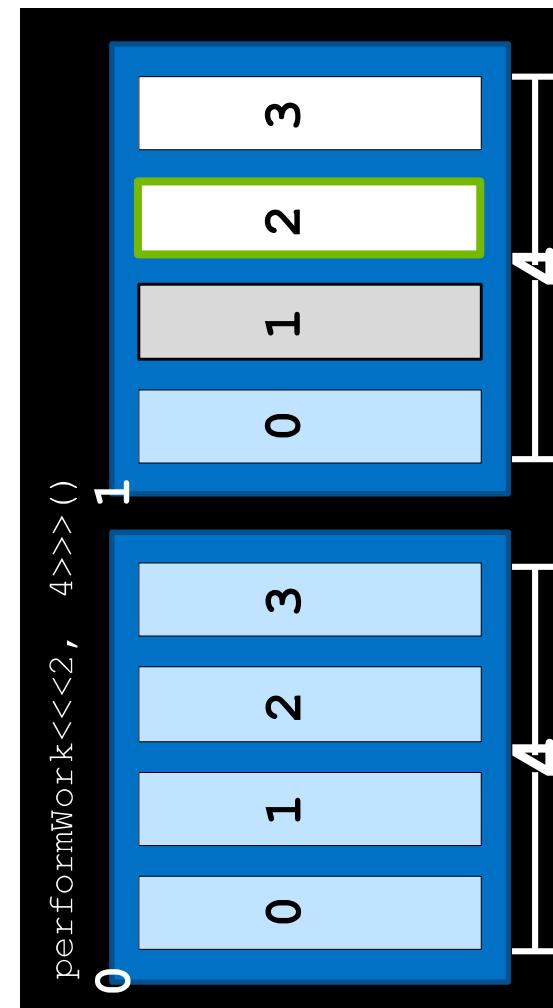


GPU

	threadIdx.x	+	blockIdx.x	*	blockIdx.y	*	blockIdx.z	
	0		2		1			
	1		dataIndex	<	N	=	Can	
	2		6		5			
	3							

0	4
1	
2	
3	

GPU
DATA

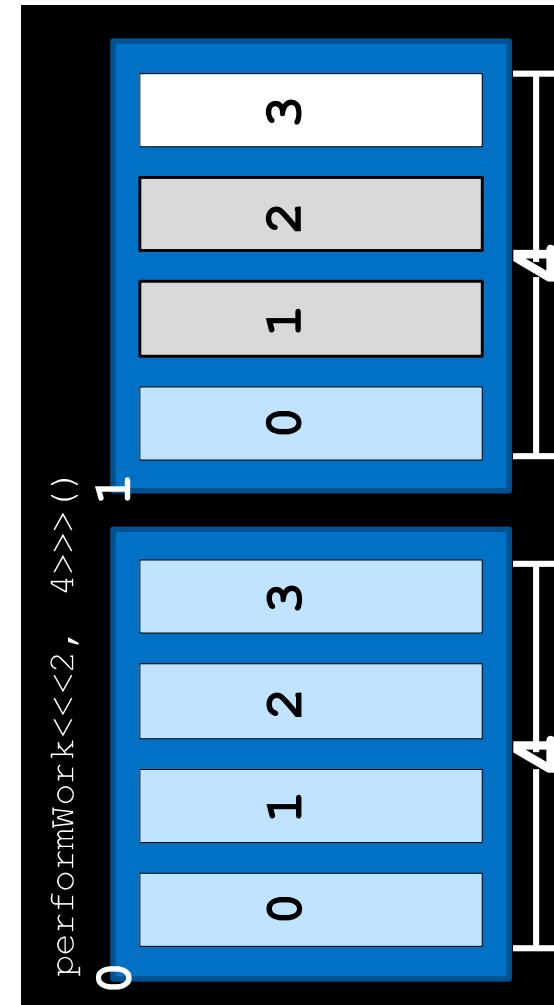


GPU

threadIdx.x	+	blockIdx.x	*	blockIdx.y
1				
2				
		dataIndex	<	N
		6		5
				fa

0	4
1	
2	
3	

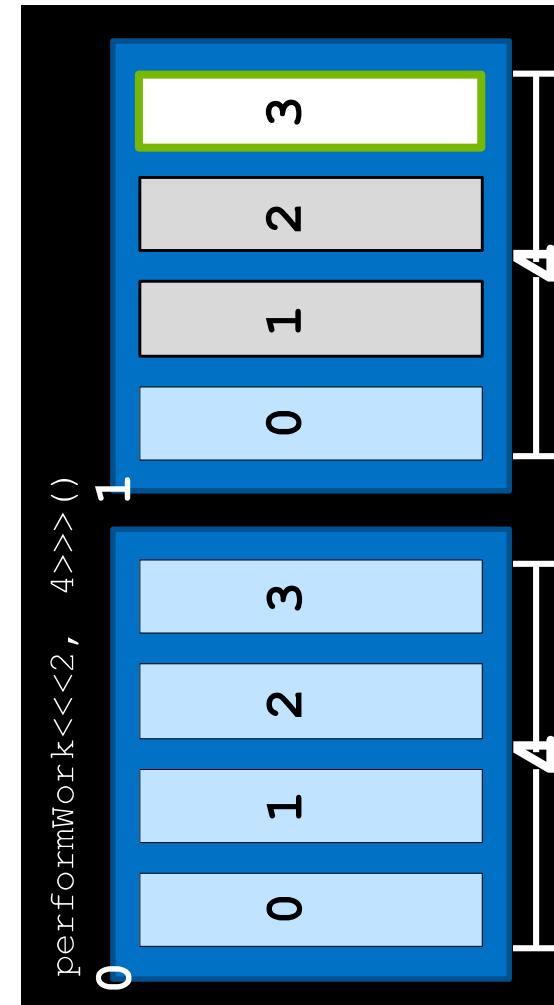
GPU
DATA



GPU

	threadIdx.x + blockIdx.x * blockDim
0	4
1	2
2	6
3	5

GPU
DATA

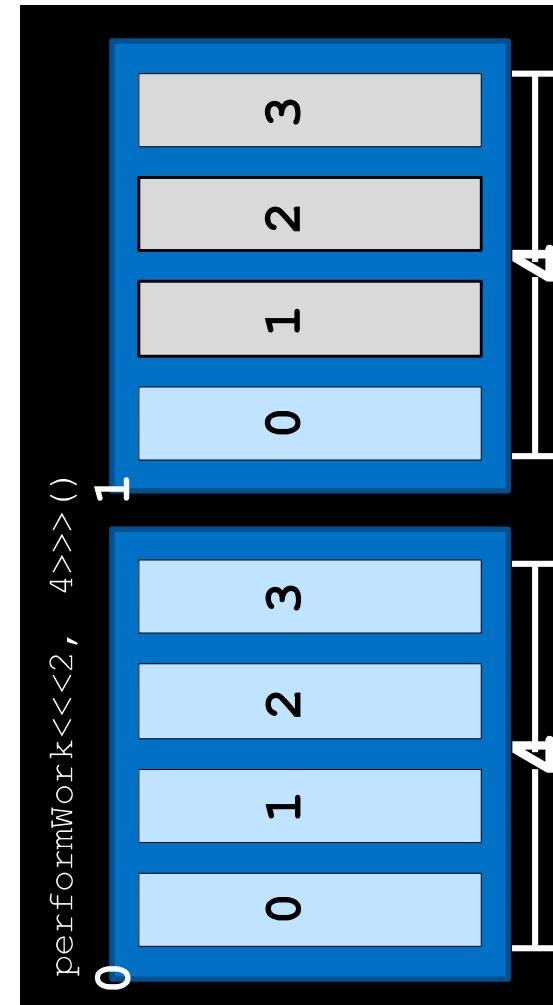


GPU

threadIdx.x	+	blockIdx.x	*	blockIdx.y
1				
2				
		dataIndex	<	N
		6		5
				fa

0	4
1	
2	
3	

GPU
DATA



GPU

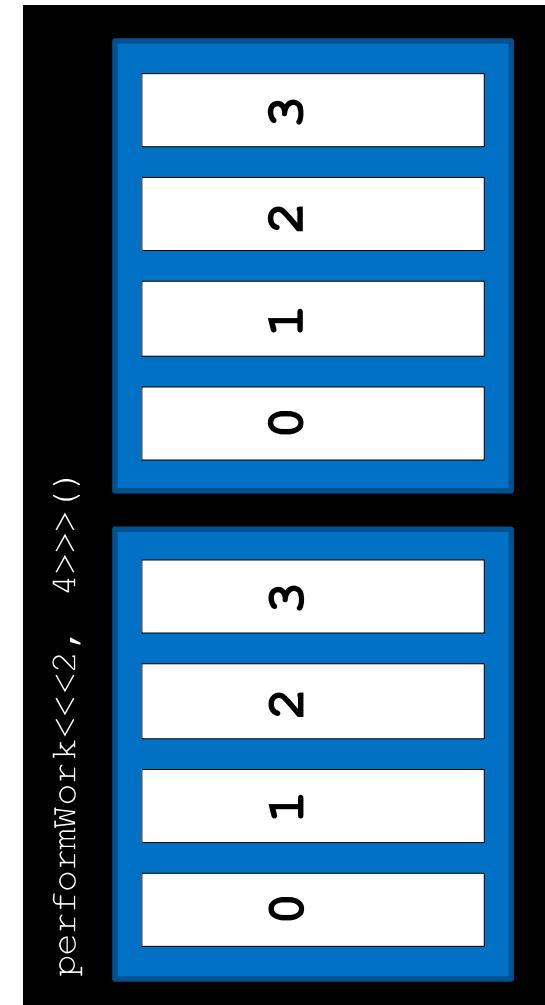
Grid-Stride Loops



Often there are more elements than threads in the grid

0	4	8	12	16	20	24	28
1	5	9	13	17	21	25	29
2	6	10	14	18	22	26	30
3	7	11	15	19	23	27	31

GPU
DATA

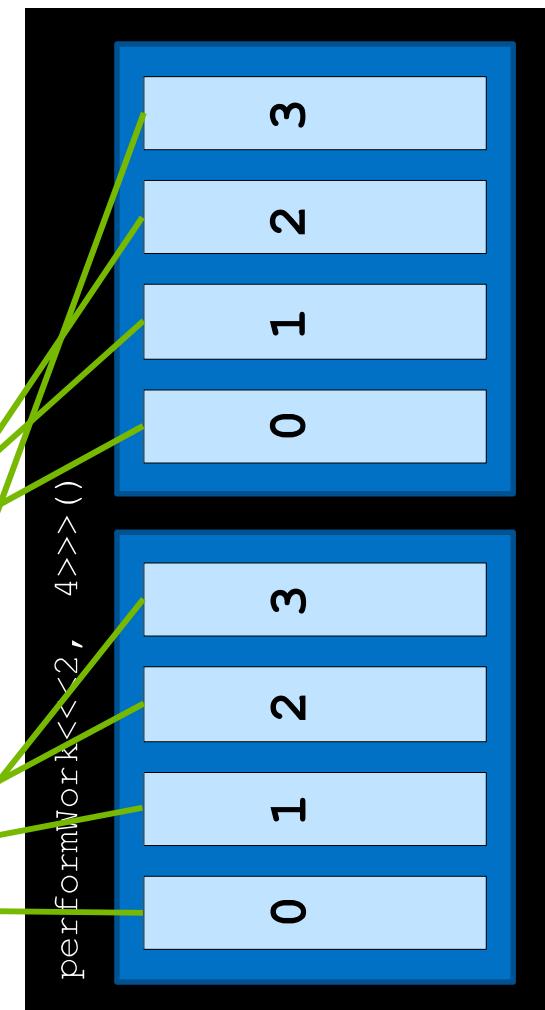


GPU

In such scenarios this cannot work on only one element

0	4	8	12	16	20	24	28
1	5	9	13	17	21	25	29
2	6	10	14	18	22	26	30
3	7	11	15	19	23	27	31

GPU DATA



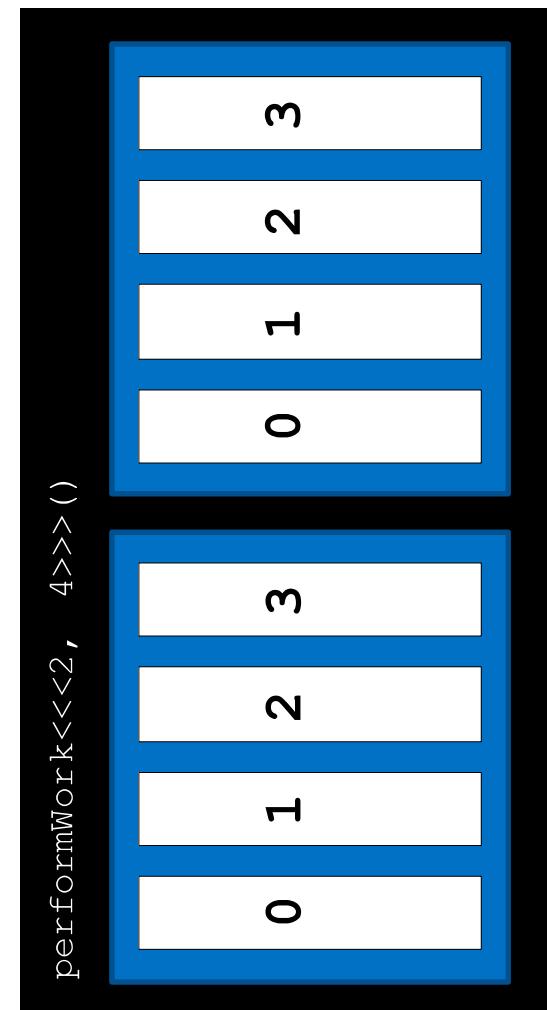
GPU

GPU DATA

0	4	8	12	16	20	24	28
1	5	9	13	17	21	25	29
2	6	10	14	18	22	26	30
3	7	11	15	19	23	27	31

... or else work is undone

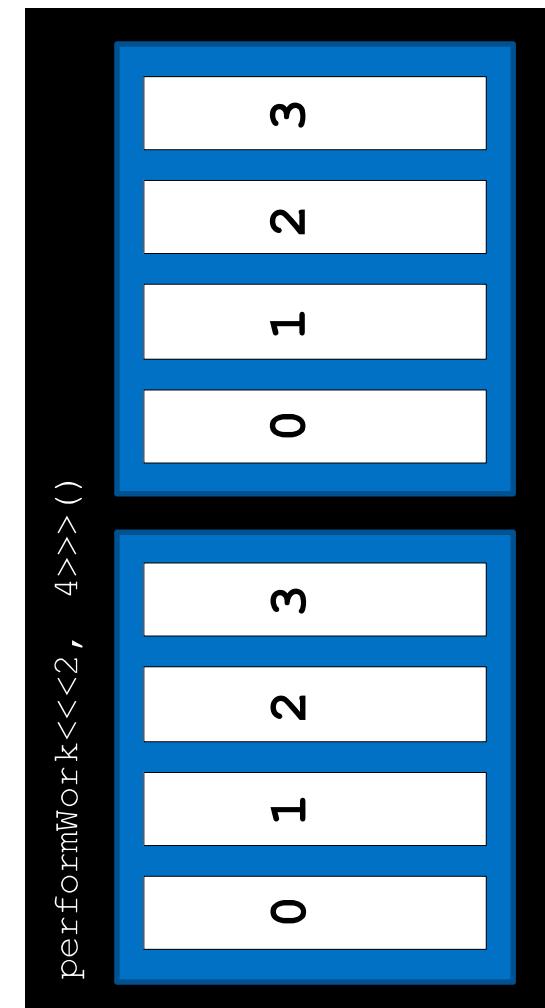
GPU



One way to address programmatically is via
grid-stride loops

0	4	8	12	16	20	24	28
1	5	9	13	17	21	25	29
2	6	10	14	18	22	26	30
3	7	11	15	19	23	27	31

GPU DATA

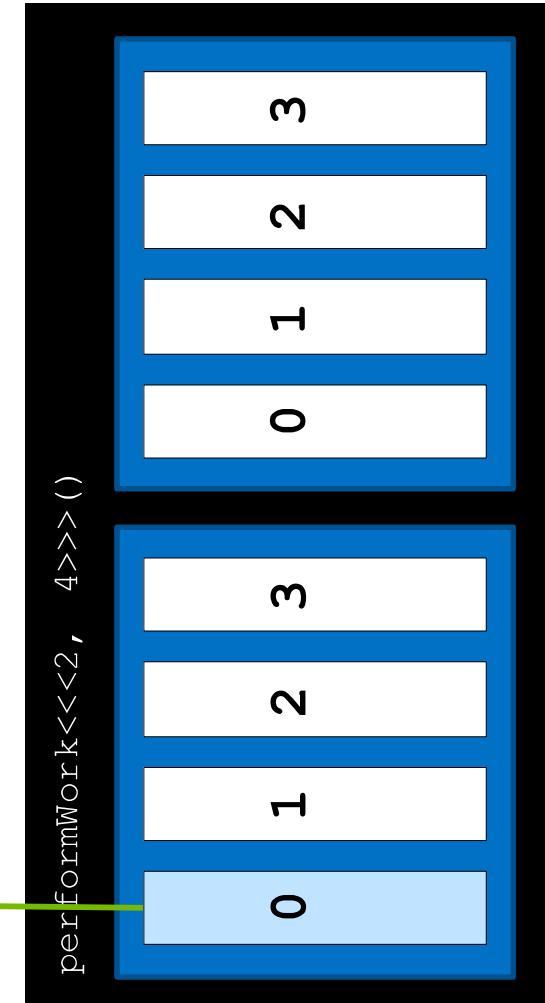


GPU

In a grid-stride loop
thread's first element
calculated as usual,
`threadIdx.x`
`blockIdx.x` *
`blockDim.x`

0	4	8	12	16	20	24	28
1	5	9	13	17	21	25	29
2	6	10	14	18	22	26	30
3	7	11	15	19	23	27	31

GPU
DATA

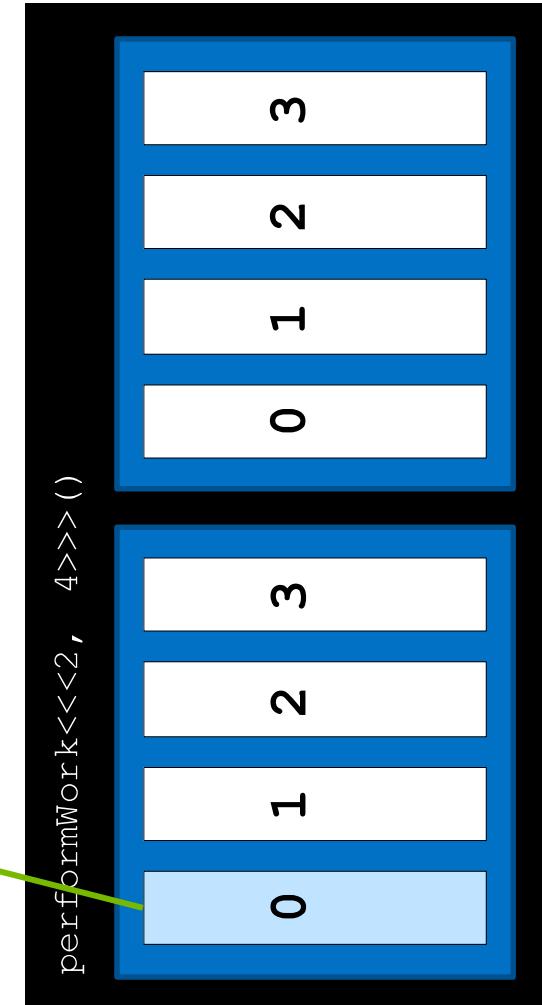


GPU

0	4	8	12	16	20	24	28
1	5	9	13	17	21	25	29
2	6	10	14	18	22	26	30
3	7	11	15	19	23	27	31

GPU DATA

The thread then str
forward by the numb
threads in the grid
(**blockDim.x**) , in this
gridDim.x) , in this
8

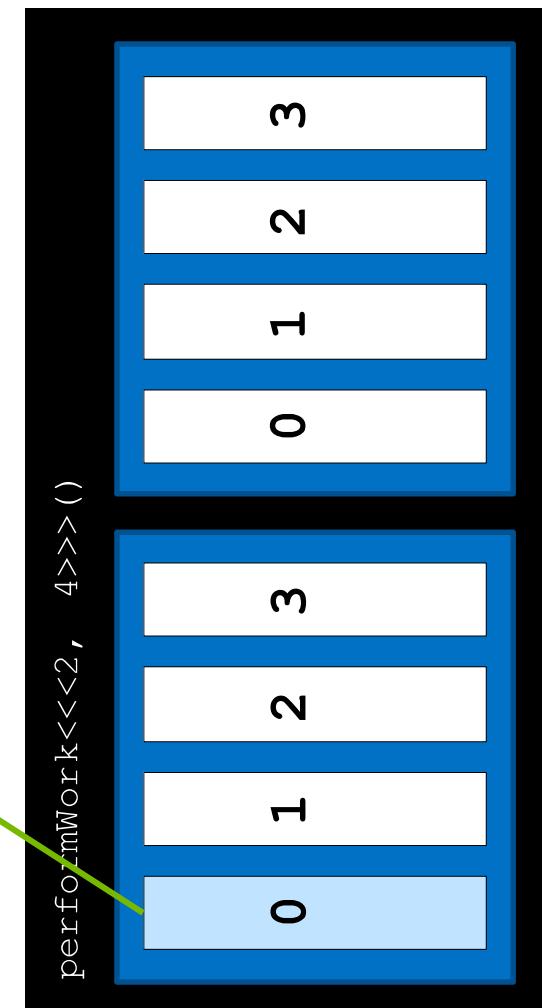


GPU

It continues in this way
its data index is greater
than the number of data
elements

0	4	8	16	20	24	28
1	5	9	13	17	21	25
2	6	10	14	18	22	26
3	7	11	15	19	23	27
						31

GPU
DATA

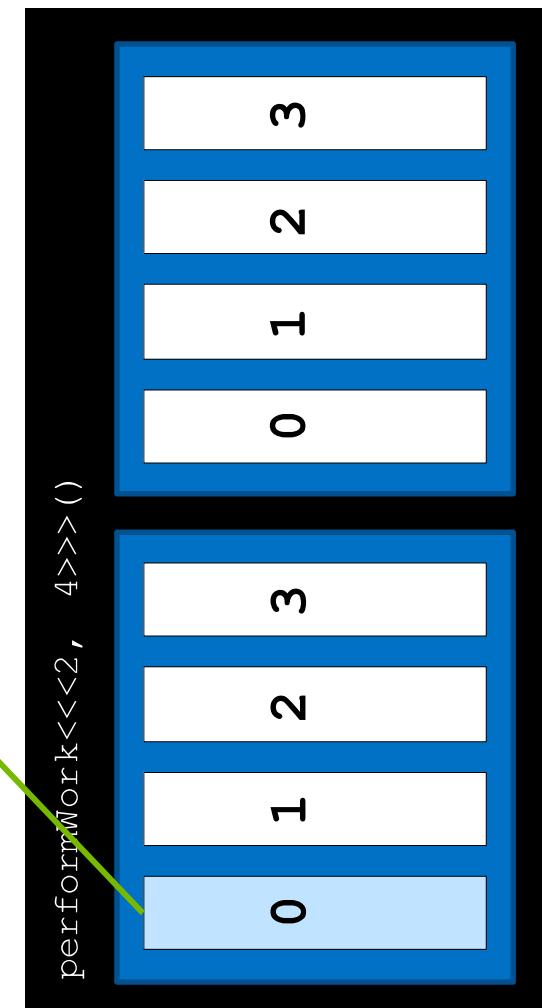


GPU

It continues in this way until its data index is greater than the number of data elements

0	4	8	12	16	20	24	28
1	5	9	13	17	21	25	29
2	6	10	14	18	22	26	30
3	7	11	15	19	23	27	31

GPU DATA

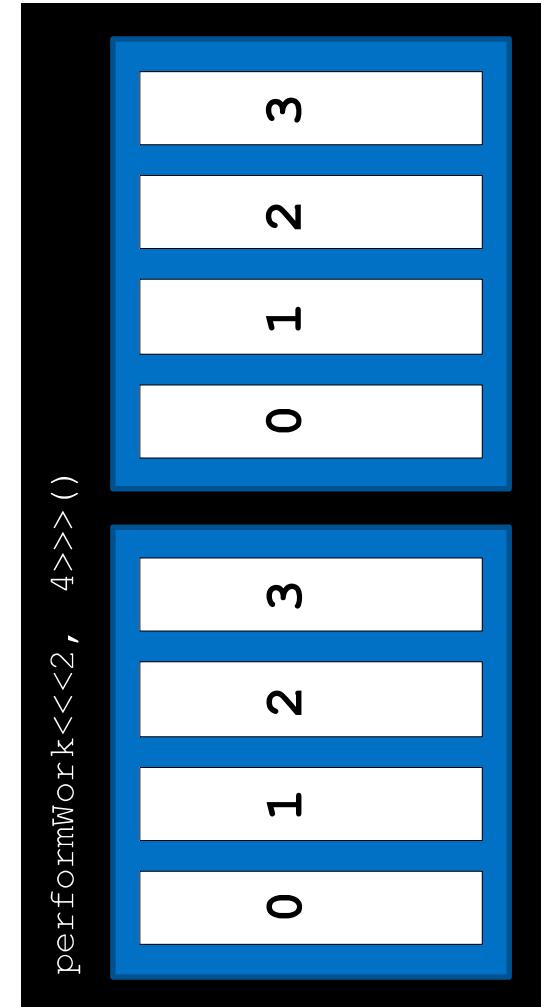


GPU

0	4	8	12	16	20	24	28
1	5	9	13	17	21	25	29
2	6	10	14	18	22	26	30
3	7	11	15	19	23	27	31

GPU
DATA

With all threads work
this way, all element
covered

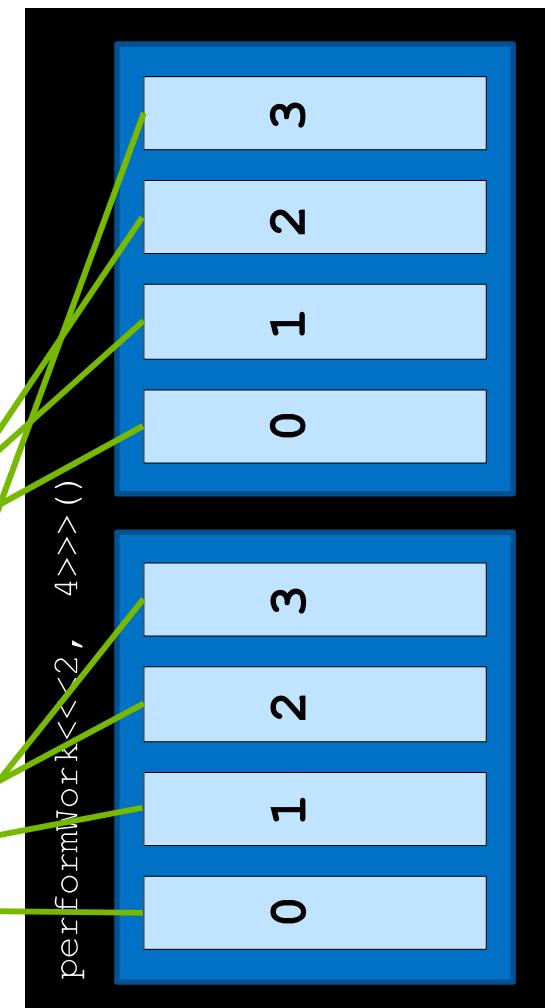


GPU

With all threads work
this way, all element
covered

0	4	8	12	16	20	24	28
1	5	9	13	17	21	25	29
2	6	10	14	18	22	26	30
3	7	11	15	19	23	27	31

GPU
DATA



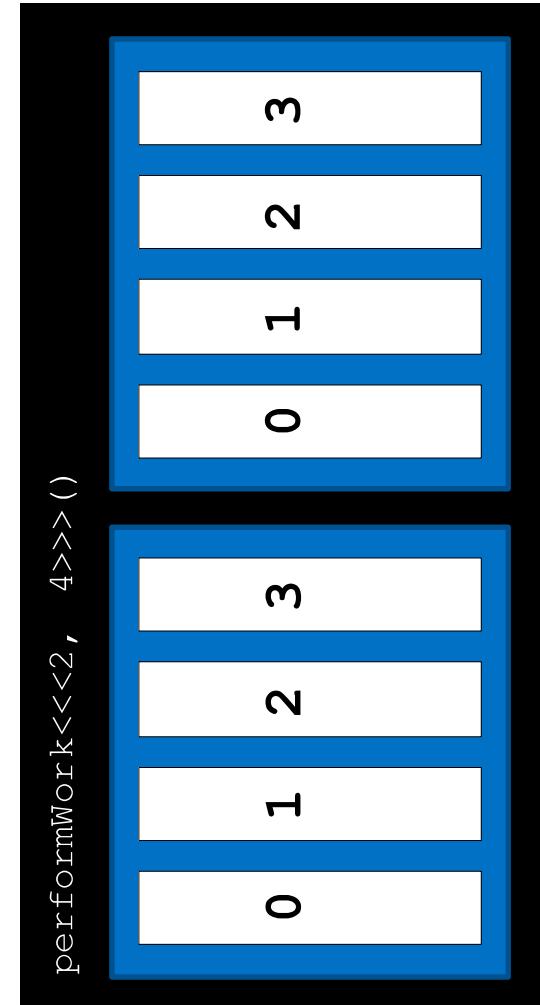
GPU

GPU DATA

0	4	8	12	16	20	24	28
1	5	9	13	17	21	25	29
2	6	10	14	18	22	26	30
3	7	11	15	19	23	27	31

With all threads work this way, all element covered

GPU

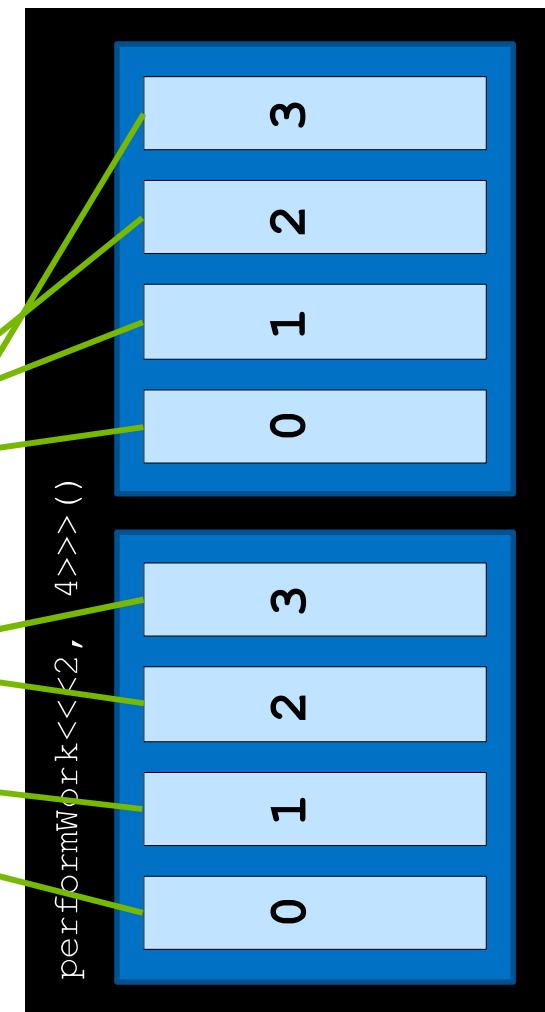


performWork<<<2, 4>>>()

With all threads work
this way, all element
covered

0	4	8	12	16	20	24	28
1	5	9	13	17	21	25	29
2	6	10	14	18	22	26	30
3	7	11	15	19	23	27	31

GPU
DATA



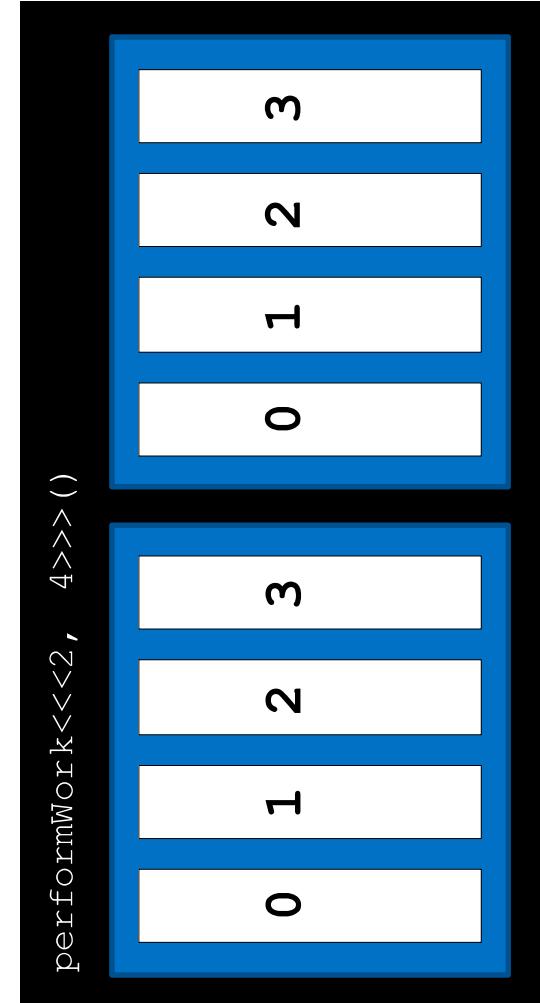
GPU

GPU DATA

0	4	8	12	16	20	24	28
1	5	9	13	17	21	25	29
2	6	10	14	18	22	26	30
3	7	11	15	19	23	27	31

With all threads work
this way, all element
covered

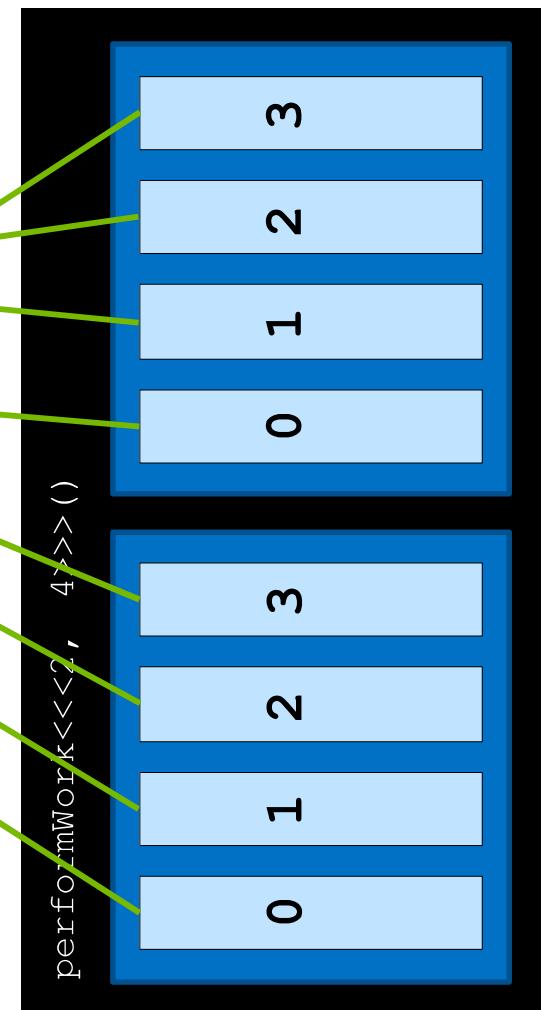
GPU



With all threads work
this way, all element
covered

0	4	8	12	16	20	24	28
1	5	9	13	17	21	25	29
2	6	10	14	18	22	26	30
3	7	11	15	19	23	27	31

GPU
DATA



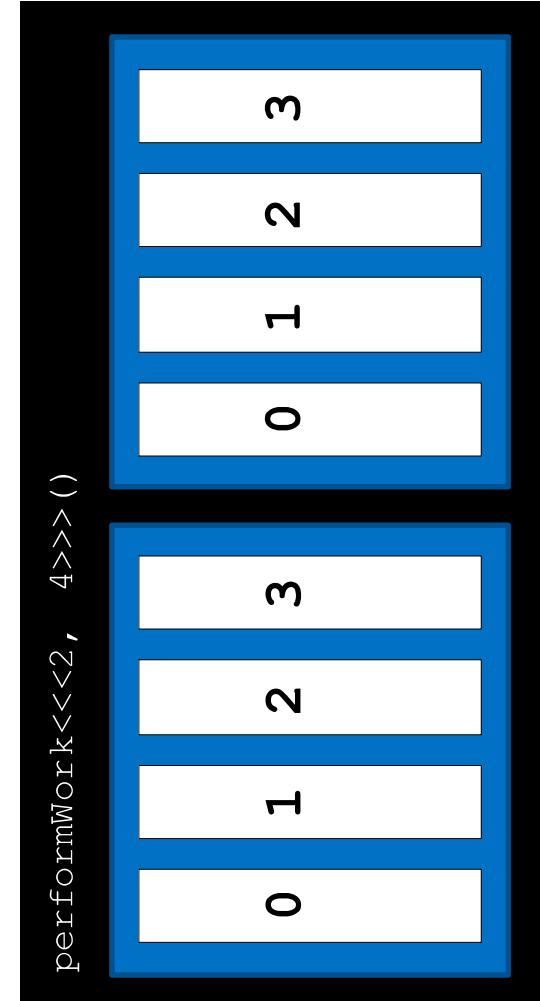
GPU

GPU DATA

0	4	8	12	16	20	24	28
1	5	9	13	17	21	25	29
2	6	10	14	18	22	26	30
3	7	11	15	19	23	27	31

With all threads work
this way, all element
covered

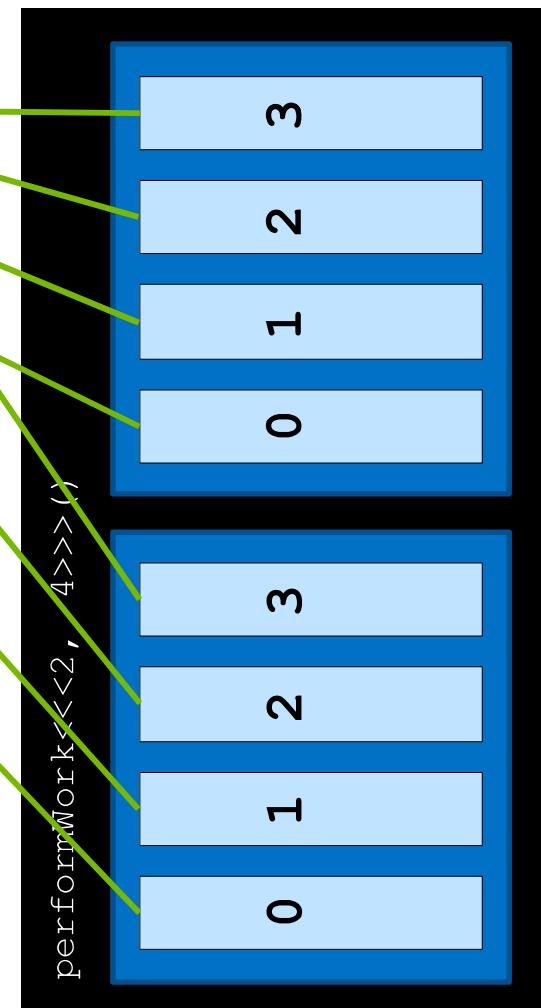
GPU



With all threads work
this way, all element
covered

0	4	8	12	16	20	24	28
1	5	9	13	17	21	25	29
2	6	10	14	18	22	26	30
3	7	11	15	19	23	27	31

GPU
DATA



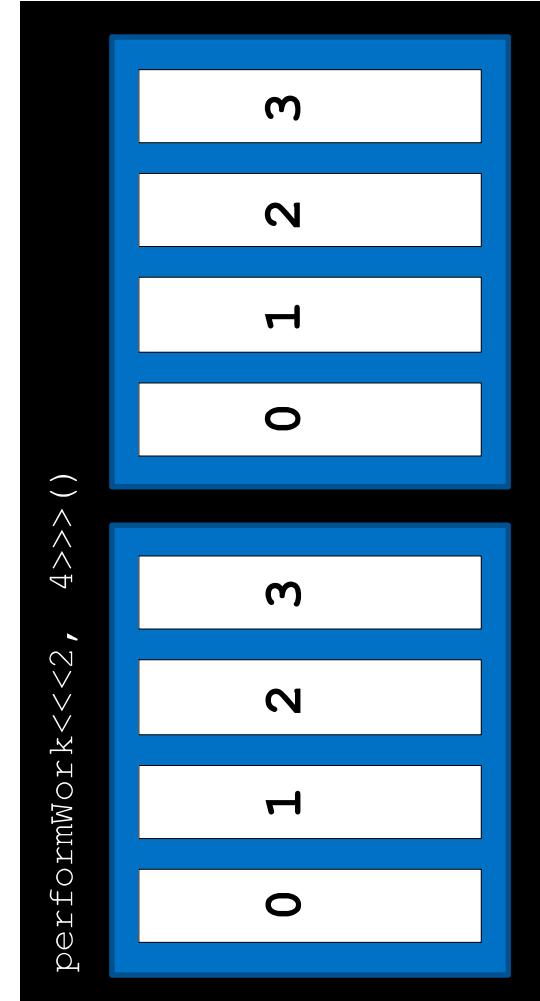
GPU

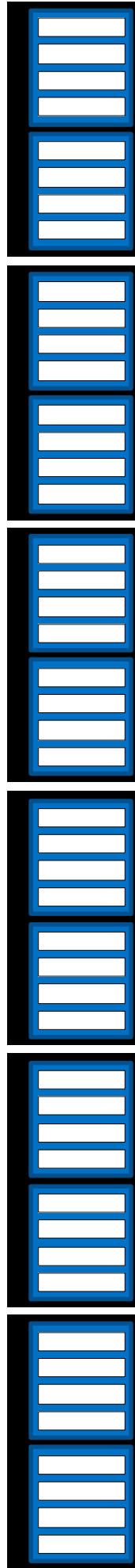
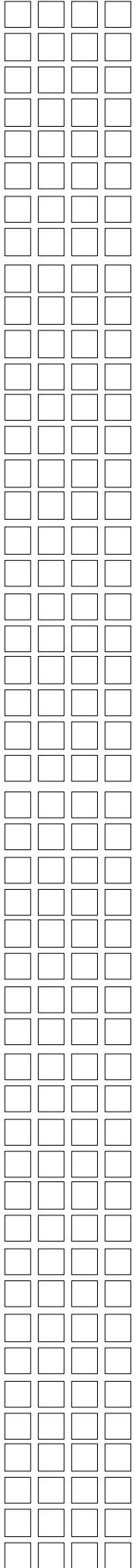
GPU DATA

0	4	8	12	16	20	24	28
1	5	9	13	17	21	25	29
2	6	10	14	18	22	26	30
3	7	11	15	19	23	27	31

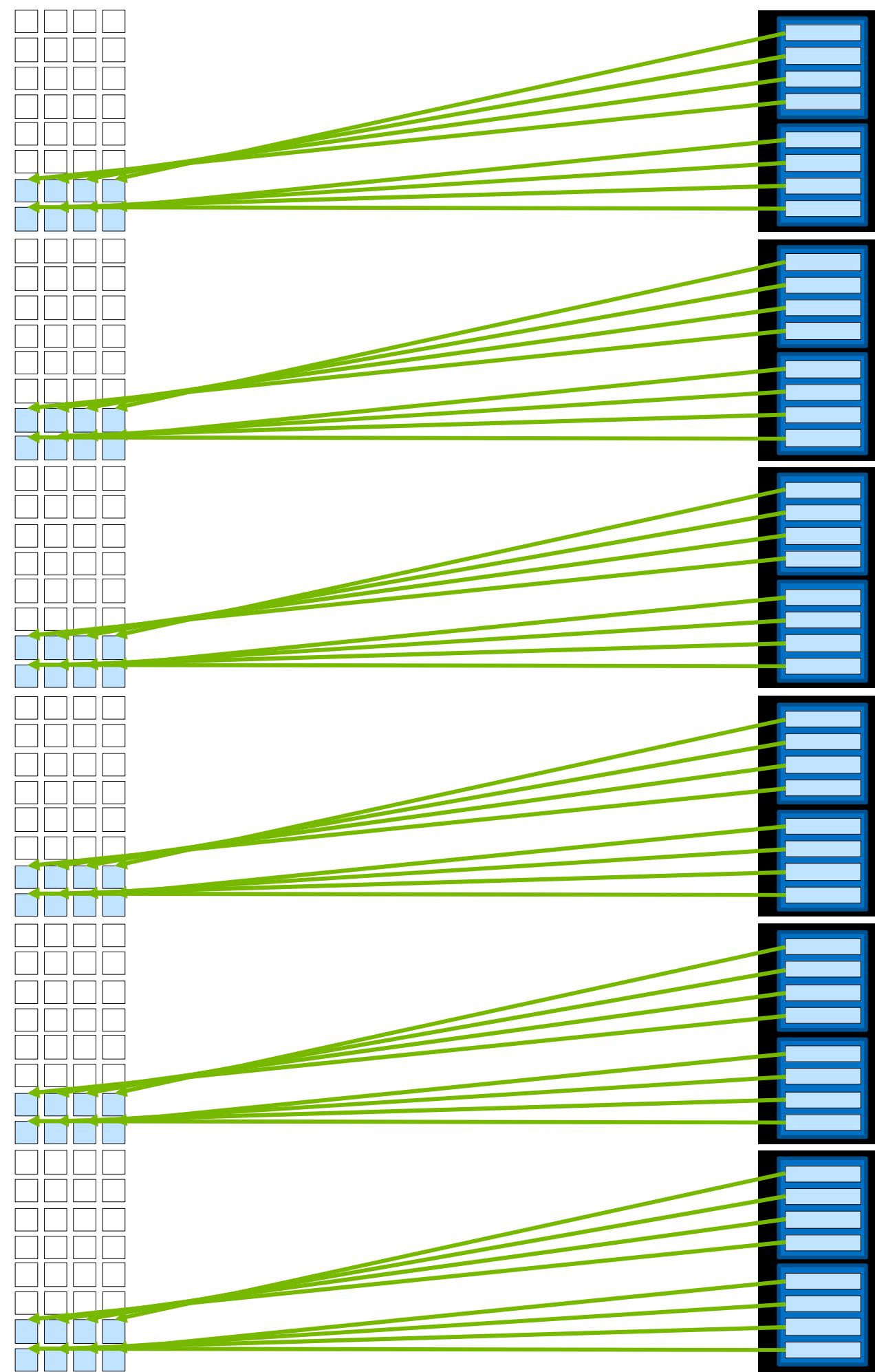
With all threads work
this way, all element
covered

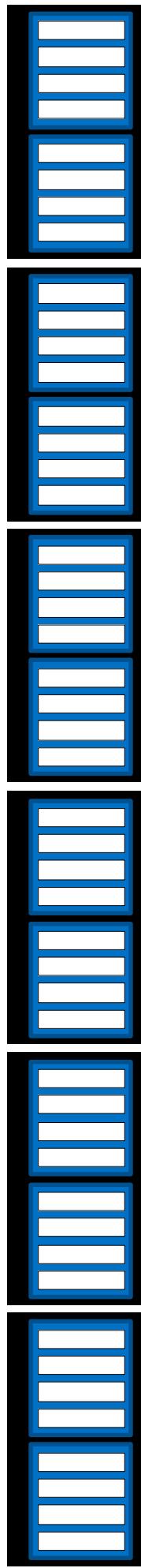
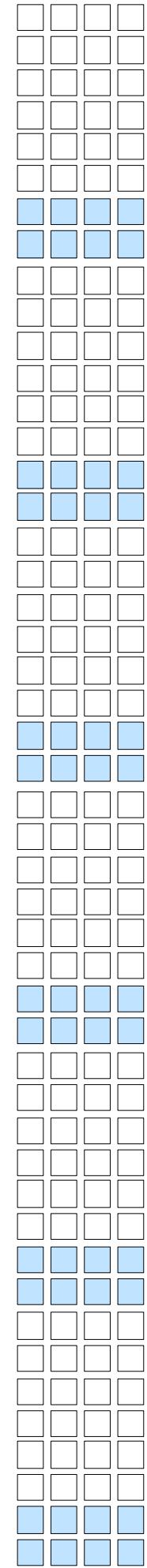
GPU

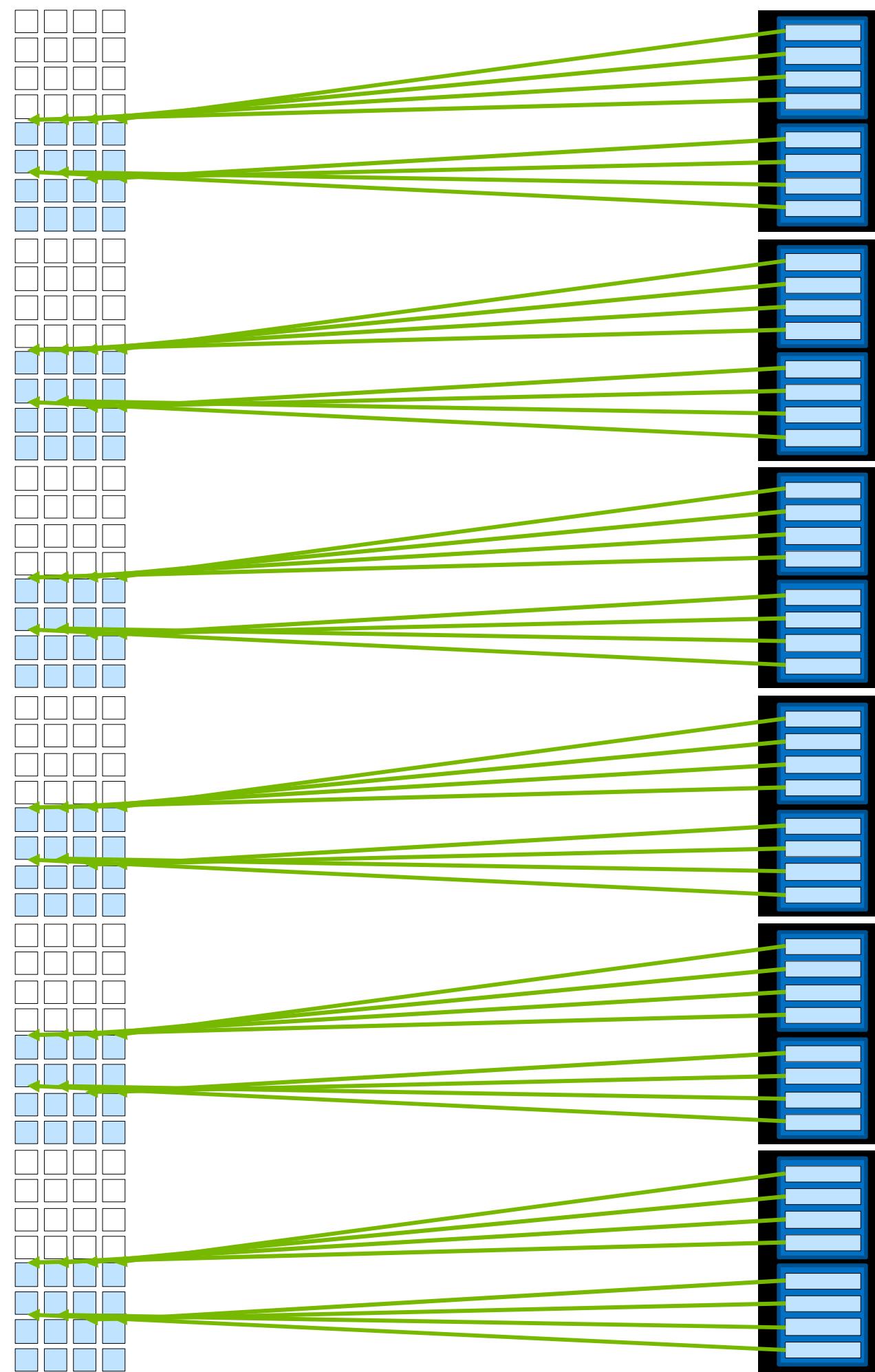


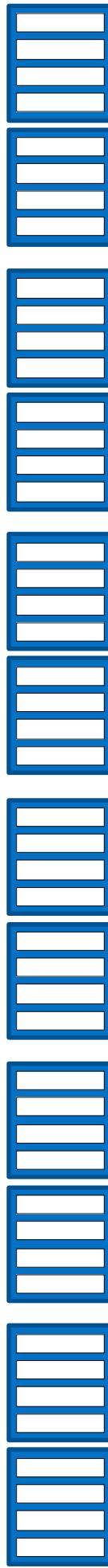
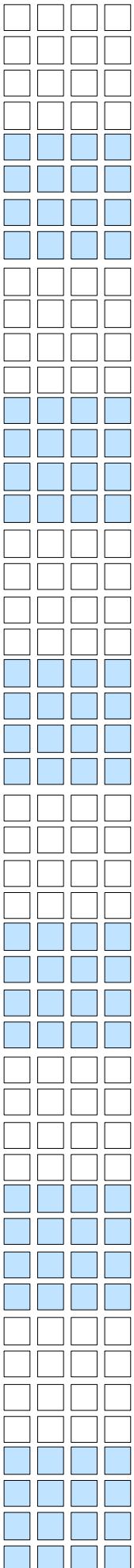


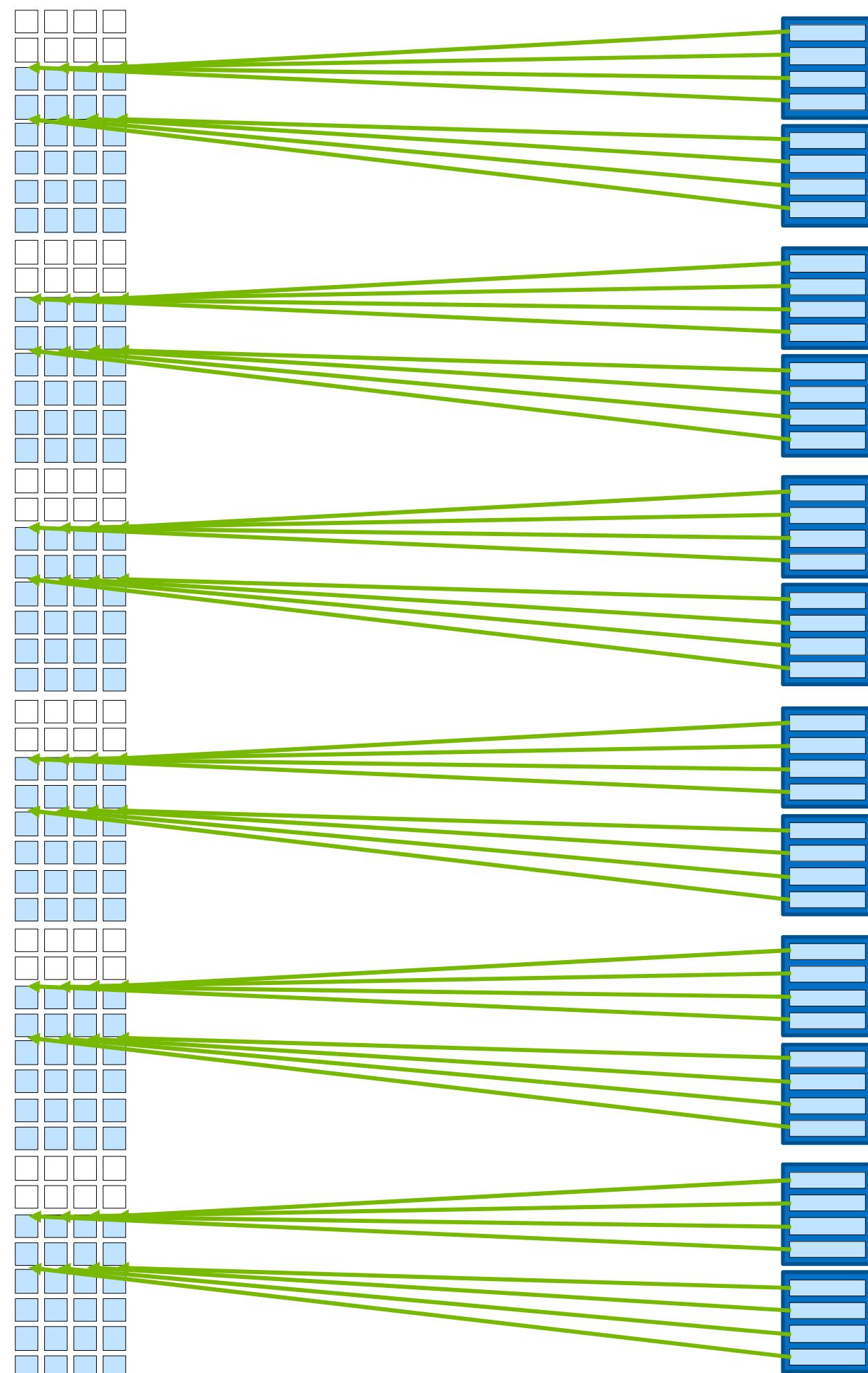
CUDA runs as many blocks in parallel at once as the GPU hardware supports, for massive parallelization

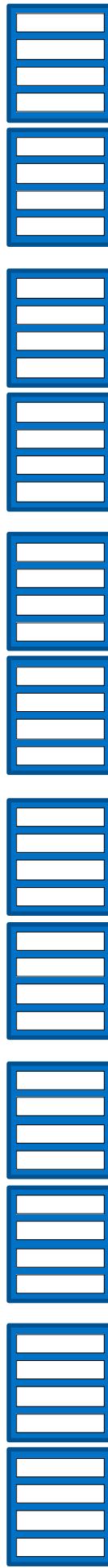
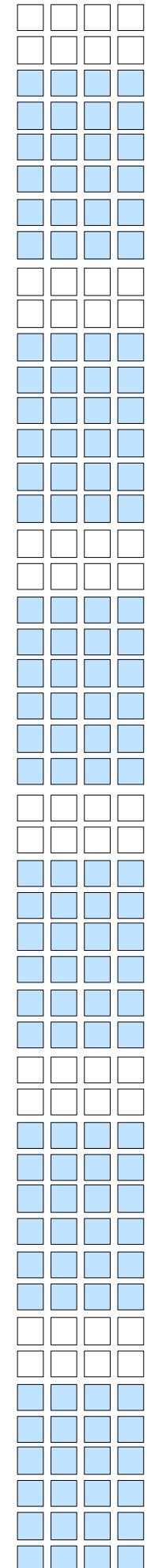


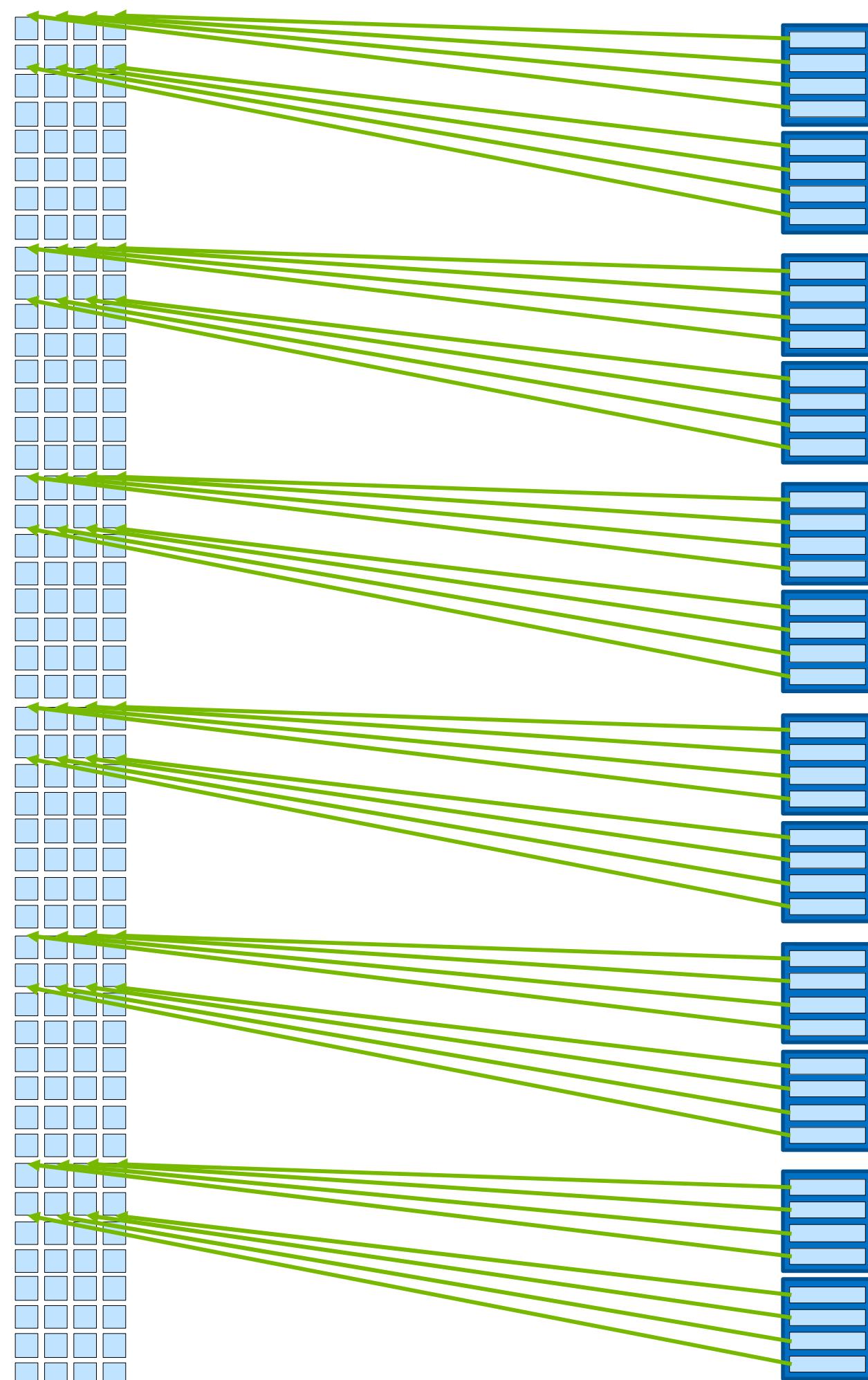


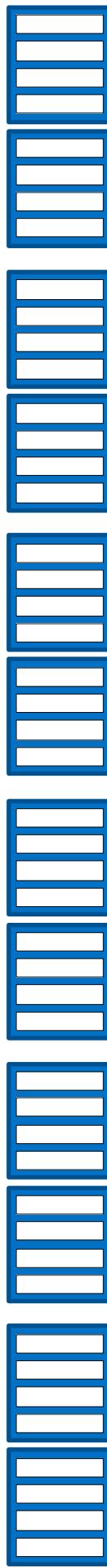
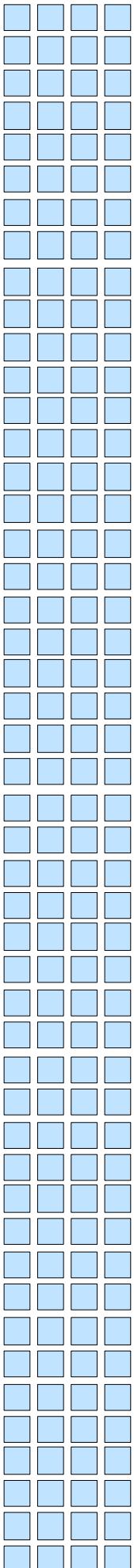












Glossary



Glossary

- **cudaMallocManaged()** : CUDA function to allocate memory accessible by both the CPU and GPUs. Memory allocated this way is called *unified memory* and is automatically migrated between the CPU and GPUs as needed.
- **cudaDeviceSynchronize()** : CUDA function that will cause the CPU to wait until the GPU is finished working.
- **Kernel**: A CUDA function executed on a GPU.
- **Thread**: The unit of execution for CUDA kernels.
- **Block**: A collection of threads.
- **Grid**: A collection of blocks.
- **Execution context**: Special arguments given to CUDA kernels when launched using the `<<...>>` syntax. It defines the number of blocks in the grid, as well as the number of threads in each block.
- **gridDim.x**: CUDA variable available inside executing kernel that gives the number of blocks in the grid
- **blockDim.x**: CUDA variable available inside executing kernel that gives the number of threads in the thread's block
- **blockIdx.x**: CUDA variable available inside executing kernel that gives the index the thread's block within the grid
- **threadIdx.x**: CUDA variable available inside executing kernel that gives the index the thread within the block
- **threadIdx.x + blockDim.x * blockIdx.x**: Common CUDA technique to map a thread to a data element
- **Grid-stride loop**: A technique for assigning a thread more than one data element to work on when there are more elements than the number of threads in the grid. The stride is calculated by `gridDim.x * blockDim.x`, which is the number of threads in the grid.



DEEP
LEARNING
INSTITUTE



www.nvidia.com/dli