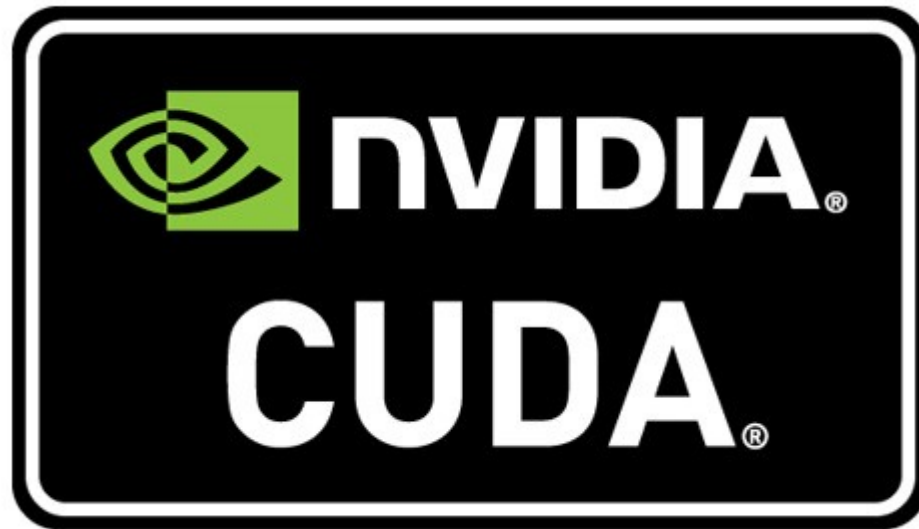
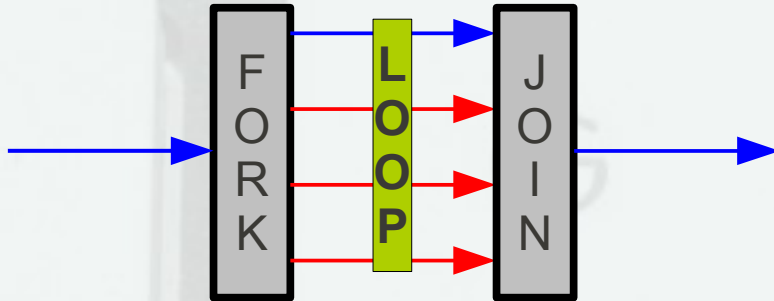


A Brief Introduction to CUDA



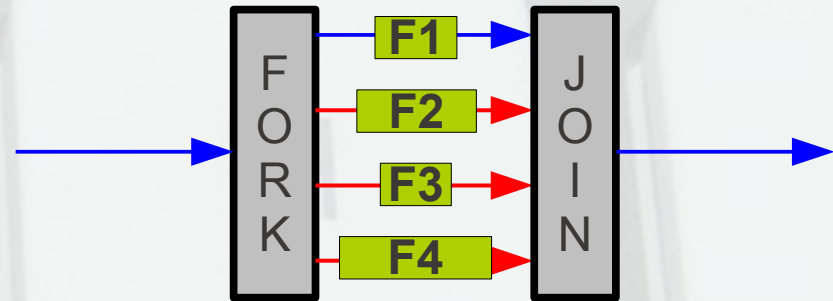
2

Types of Parallelism

DATA
PARALLELISM

Cores execute *same instructions*

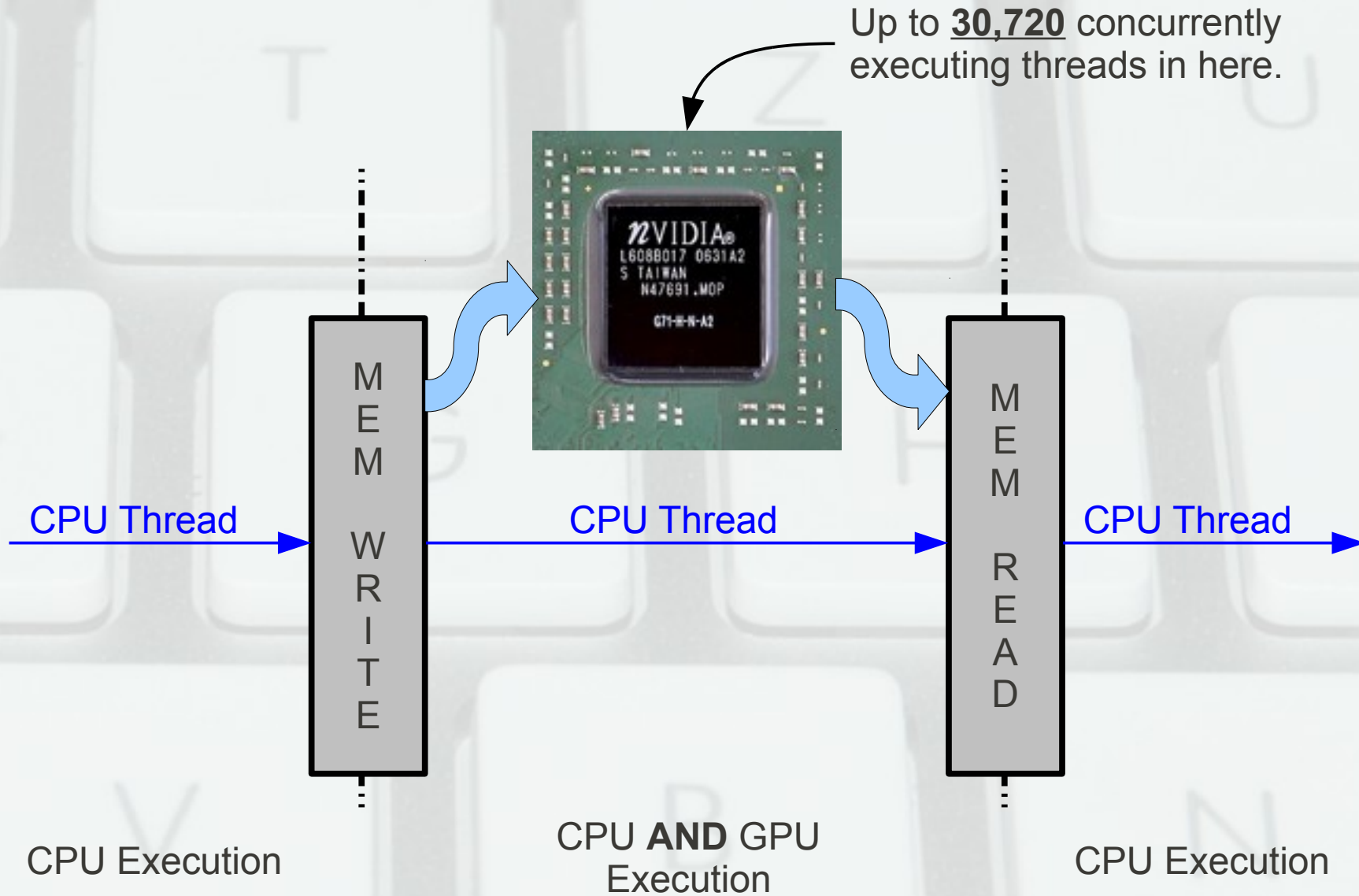
...but operate on *different data*.

FUNCTIONAL
PARALLELISM

Cores execute *different instructions*

...and can read same data &
should write different data.

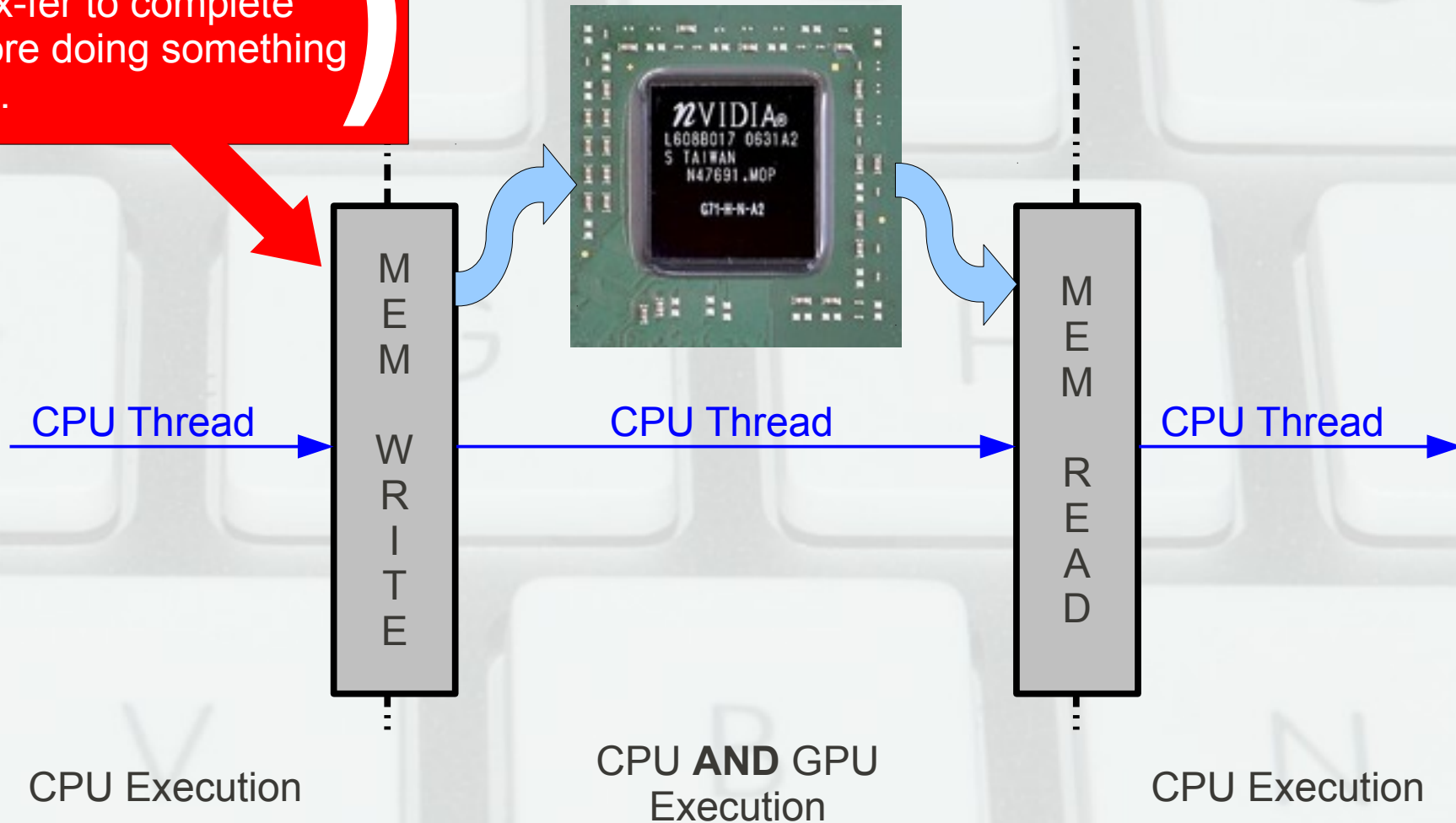
BASIC IDEA



BASIC IDEA

CPU → GPU
Memory transfers are
blocking.

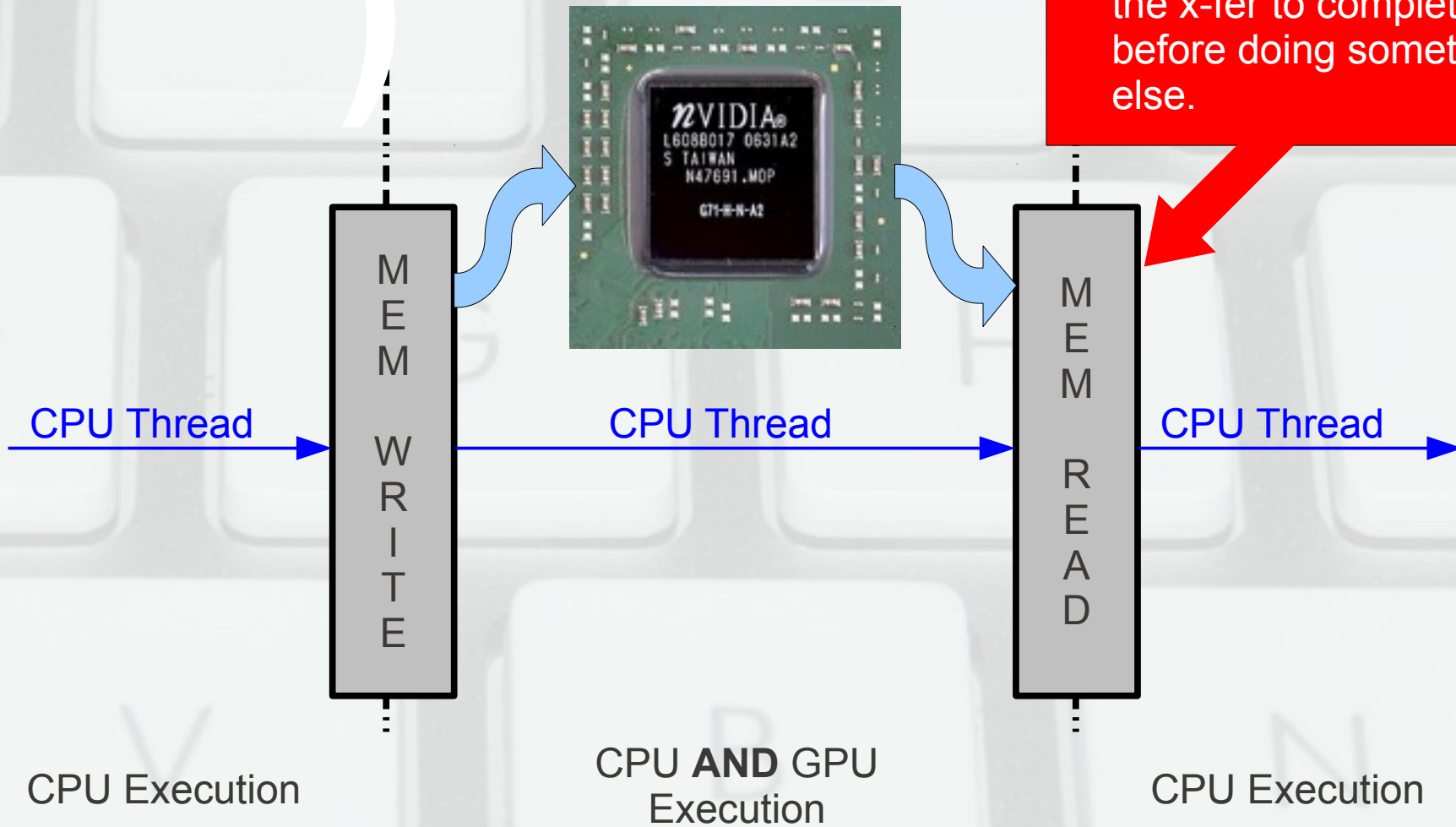
(The CPU must wait for
the x-fer to complete
before doing something
else.)



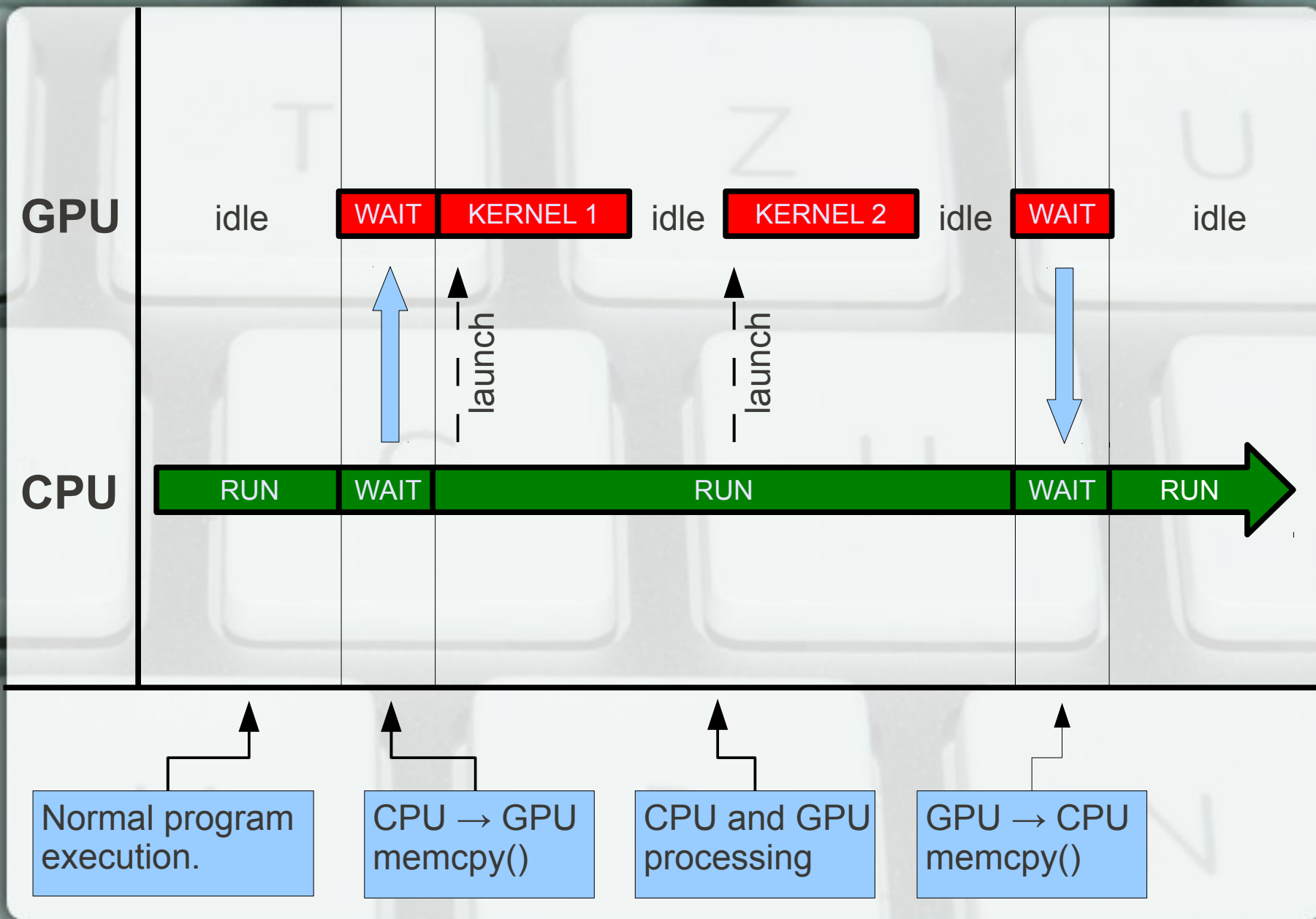
BASIC IDEA

GPU → CPU
Memory transfers are
also blocking.

The CPU must wait for
the x-fer to complete
before doing something
else.



BASIC EXECUTION FLOW



BASIC EXECUTION FLOW

KERNEL 1

A diagram illustrating the basic execution flow. A red rectangular box labeled "KERNEL 1" is positioned at the top. A thick black arrow points from a larger black rectangular box below it up towards the "KERNEL 1" box. The black box contains text explaining that kernels consist of many parallel threads and that thread IDs are used for indexing into arrays, with a note that most array accesses must be coalesced.

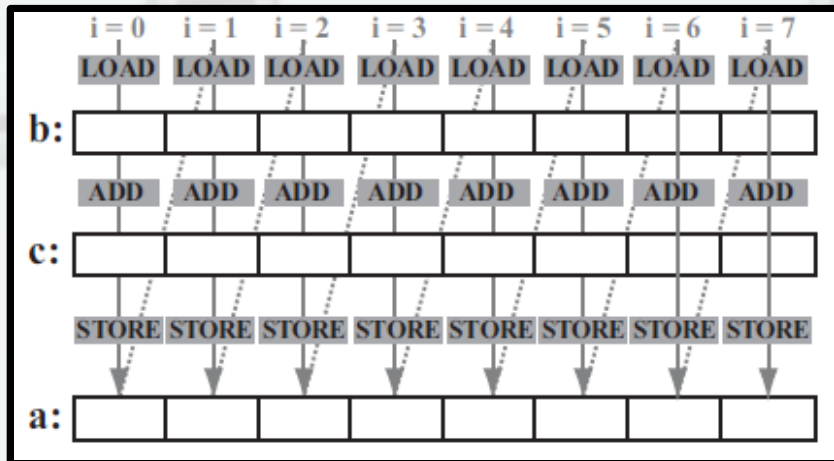
Kernels consist of MANY parallel threads

- Thread IDs are used to index into arrays
- (Most) Array accesses must be coalesced.

The Most Basic Example

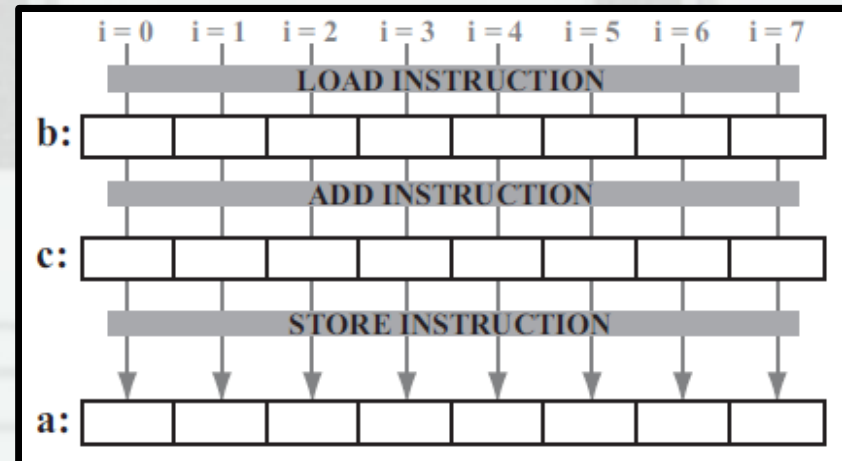
CPU

```
void
func (float* a, float* b, float* c)
{
    int i;
    for (i=0; i<8; i++) {
        a[i] = b[i] + c[i];
    }
}
```



GPU Kernel

```
__global__ void
kernel (float* a, float* b, float* c)
{
    int i = threadIdx.x;
    a[i] = b[i] + c[i];
}
```



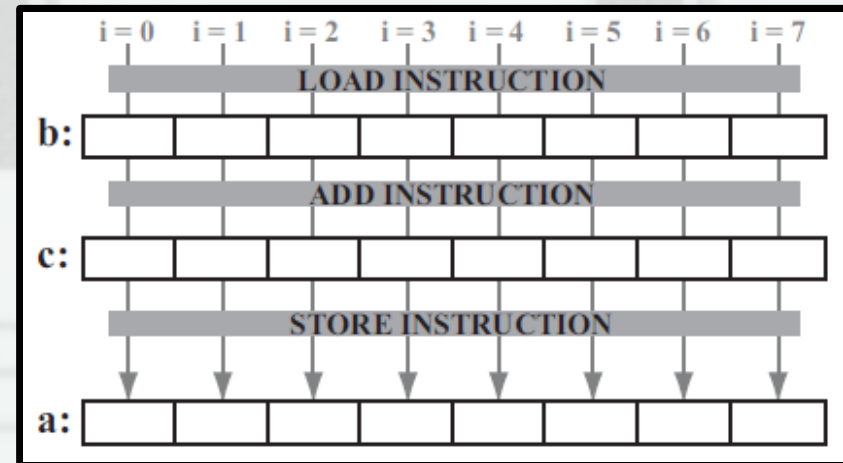
The Most Basic Example

GPU Kernel

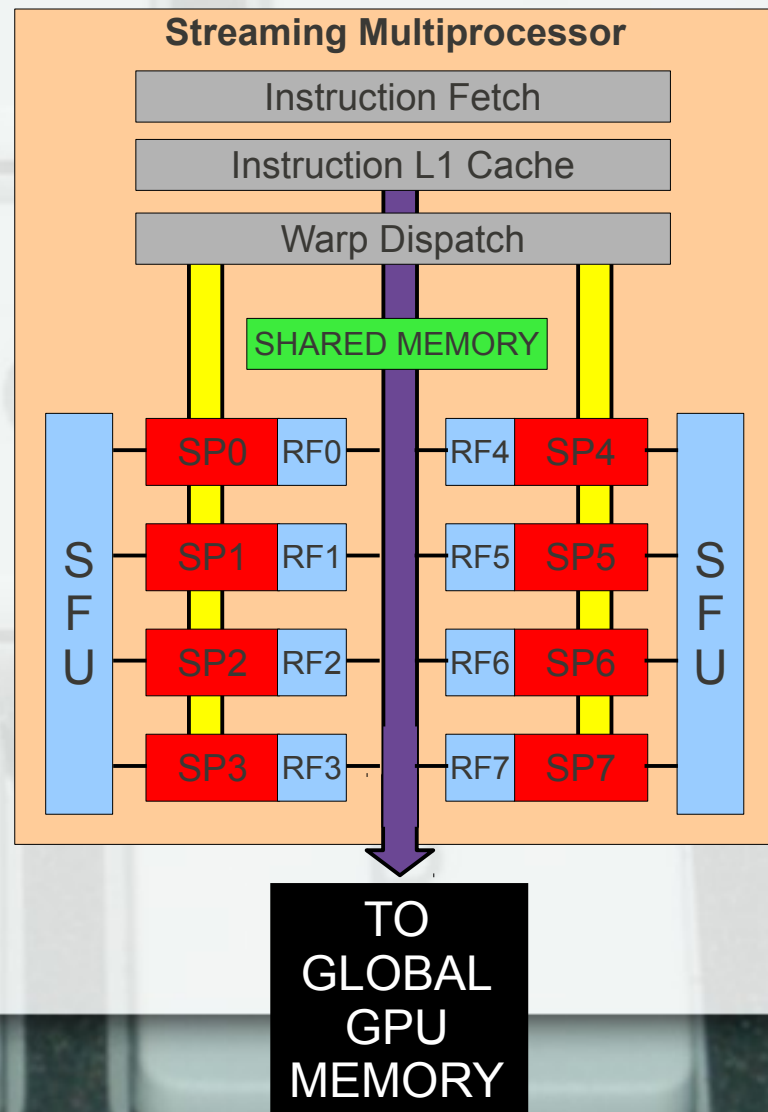
```
__global__ void  
kernel (float* a, float* b, float* c)  
{  
    int i = threadIdx.x;  
    a[i] = b[i] + c[i];  
}
```

Thread IDs used to index
into arrays.

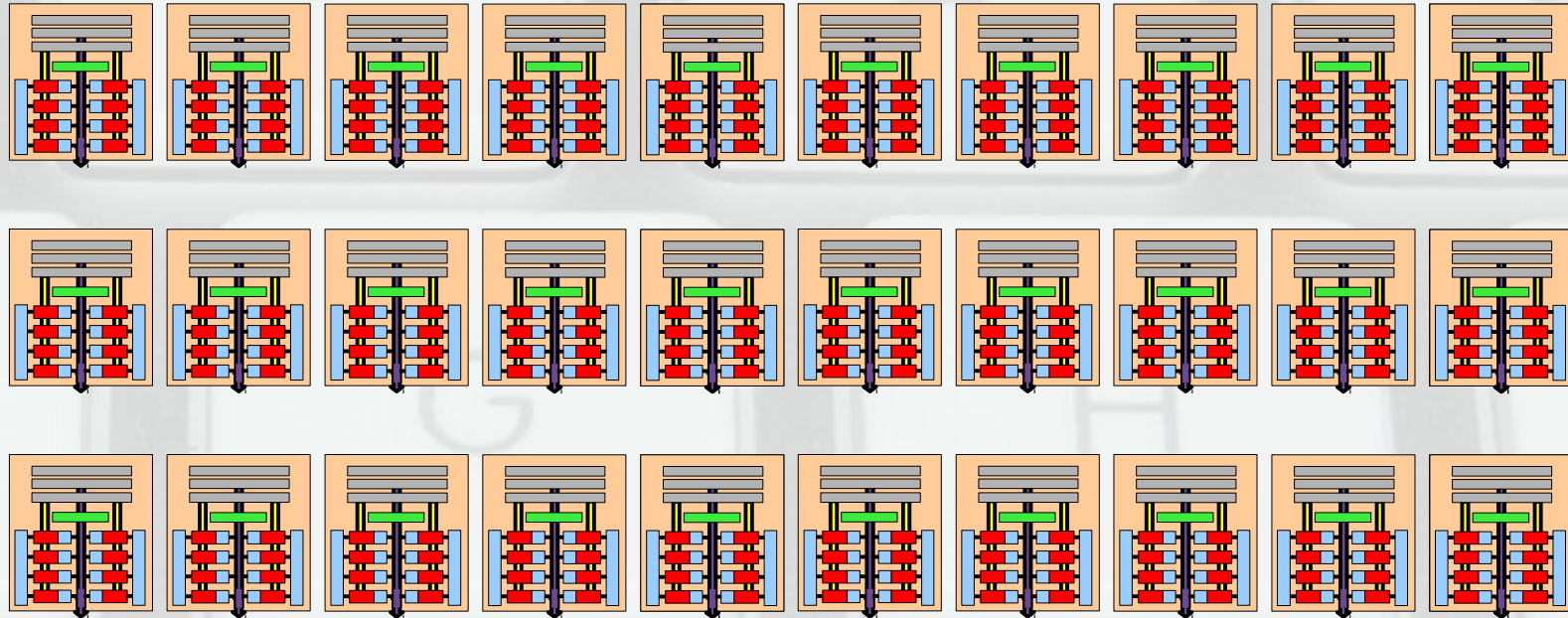
Data reads & writes are
coalesced.



Threading Strategies

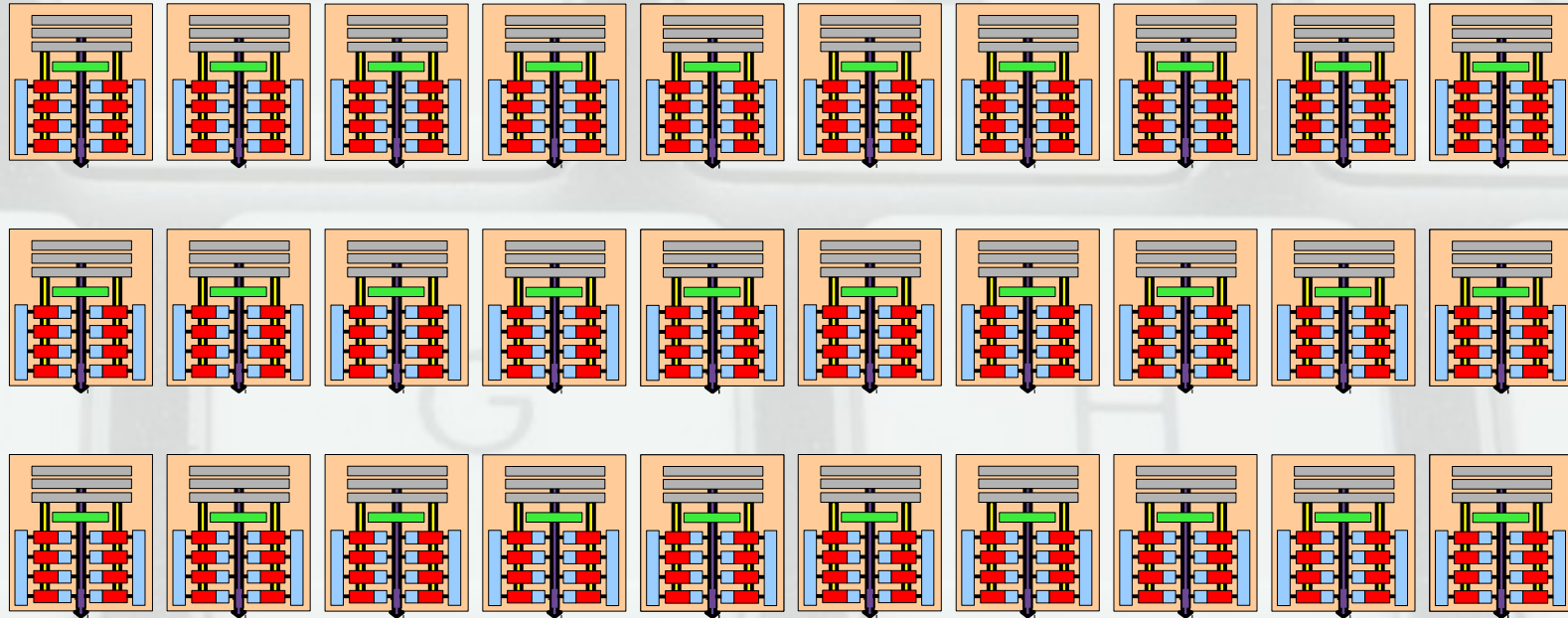


Threading Strategies



30
SMs

Threading Strategies

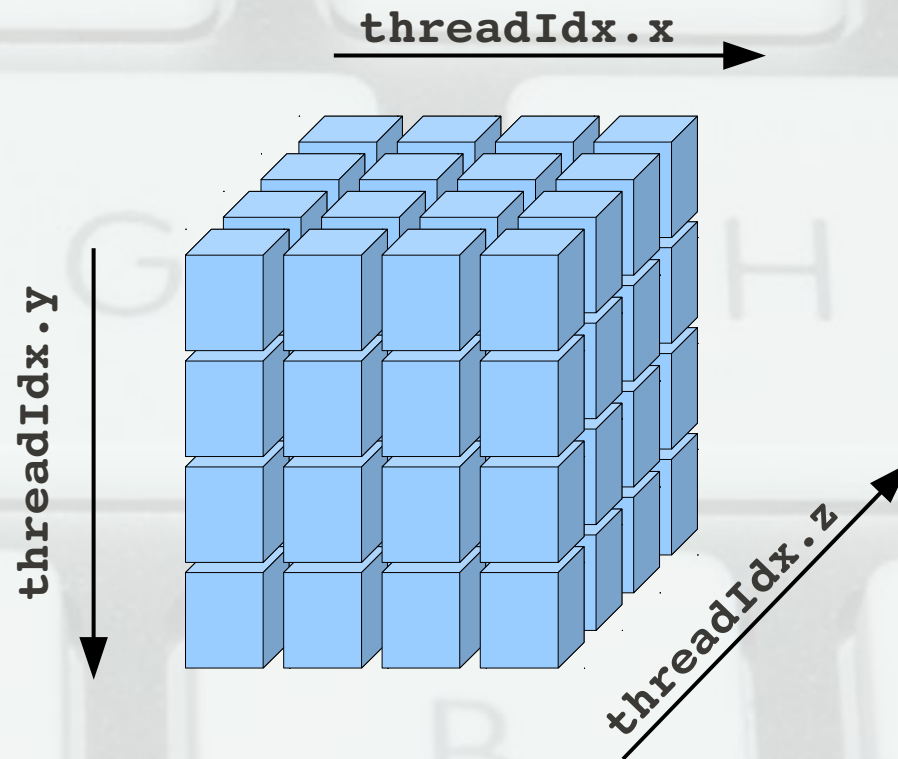


Up to
8 thread blocks
per SM*

*1024 thread capacity per SM. Block assignment subject to SM resource availability.

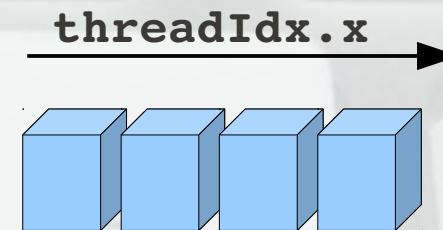
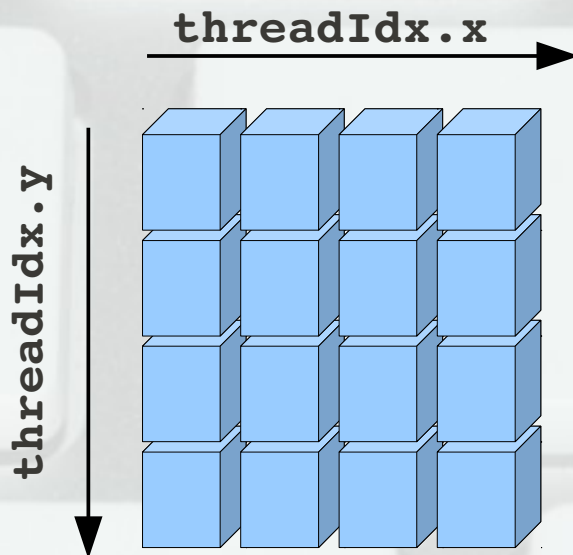
Threading Strategies

Threads are organized into 3D blocks



Threading Strategies

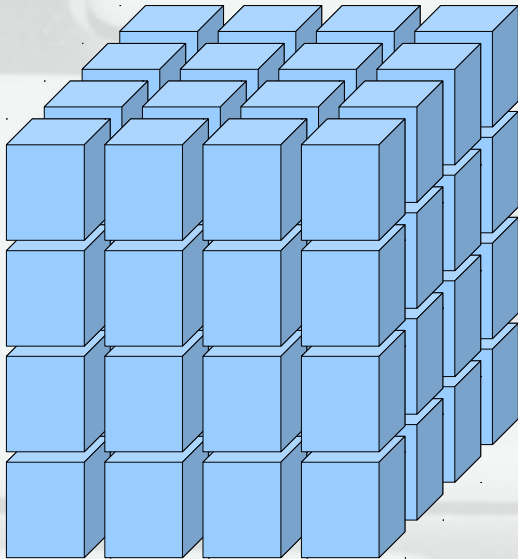
2D and 1D configurations are also possible



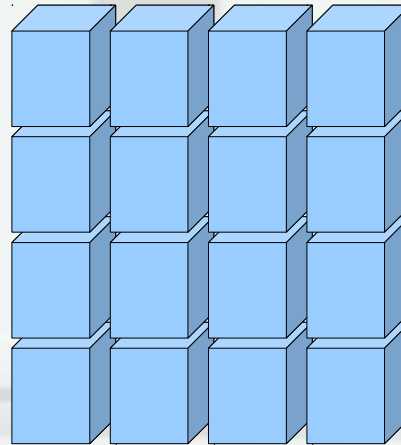
This makes indexing into arrays easier.

Threading Strategies

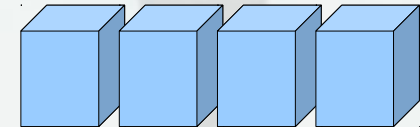
2D and 1D configurations are also possible



Good for 3D data



Good for 2D data



Good for 1D data

Max threads per block is always
512

Threading Strategies

GPU Kernel – 1D Data

```
__global__ void  
kernel (float* a, float* b, float* c)  
{  
    int i = threadIdx.x; // thread x coord  
  
    a[i] = b[i] + c[i];  
}
```

Threading Strategies

GPU Kernel – 2D Data

```
__global__ void  
kernel (float* a, float* b, float* c)  
{  
    int i = threadIdx.x; // thread x coord  
    int j = threadIdx.y; // thread y coord  
  
    int d = blockDim.x; // block x dim  
  
    int index = d.x*j + i;  
  
    a[index] = b[index] + c[index];  
}
```

Threading Strategies

GPU Kernel – 3D Data

```
__global__ void
kernel (float* a, float* b, float* c)
{
    int i = threadIdx.x; // thread x coord
    int j = threadIdx.y; // thread y coord
    int k = threadIdx.z; // thread z coord

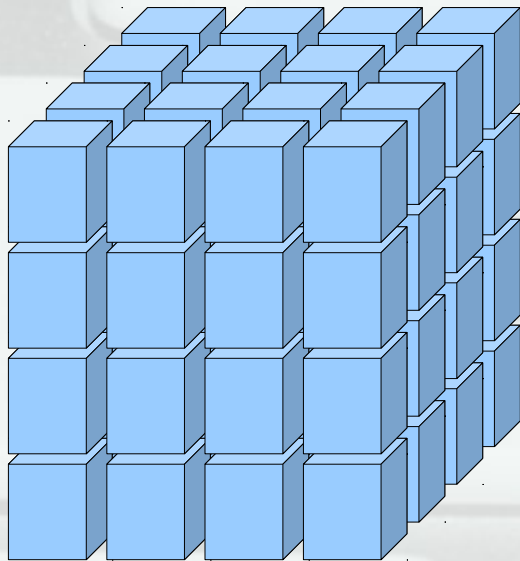
    int2 d;
    d.x = blockDim.x; // block x dim
    d.y = blockDim.y; // block y dim

    int index = d.x*d.y*k + d.x*j + i;

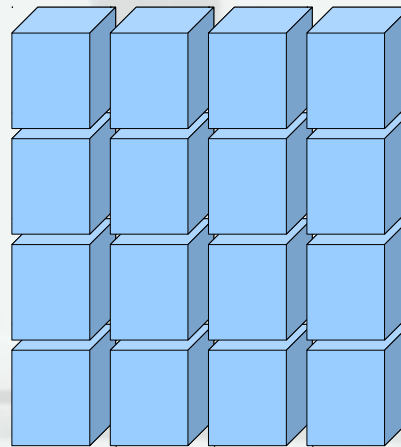
    a[index] = b[index] + c[index];
}
```


Threading Strategies

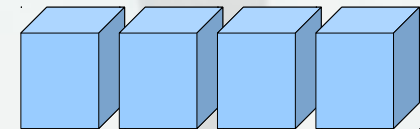
2D and 1D configurations are also possible



Good for 3D data



Good for 2D data



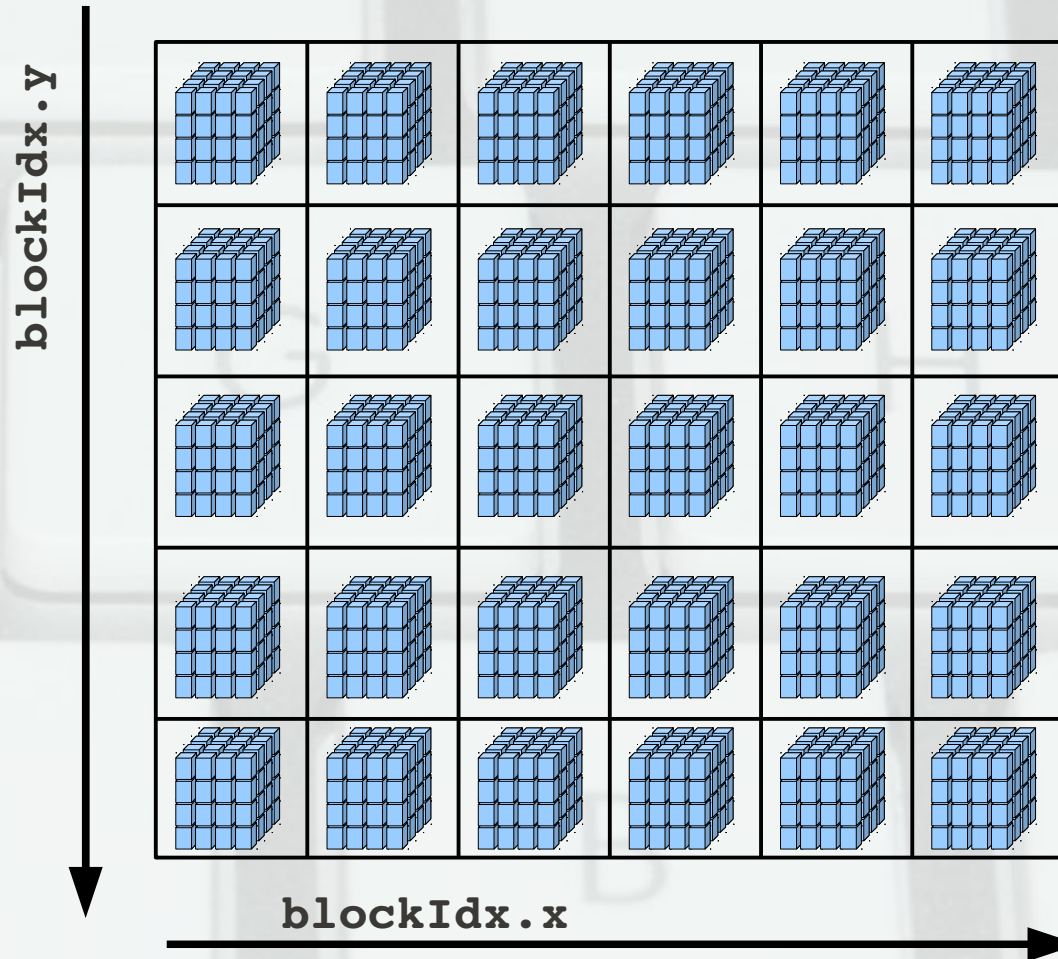
Good for 1D data

What if my input data has more than 512 elements?

Max threads per block is always
512

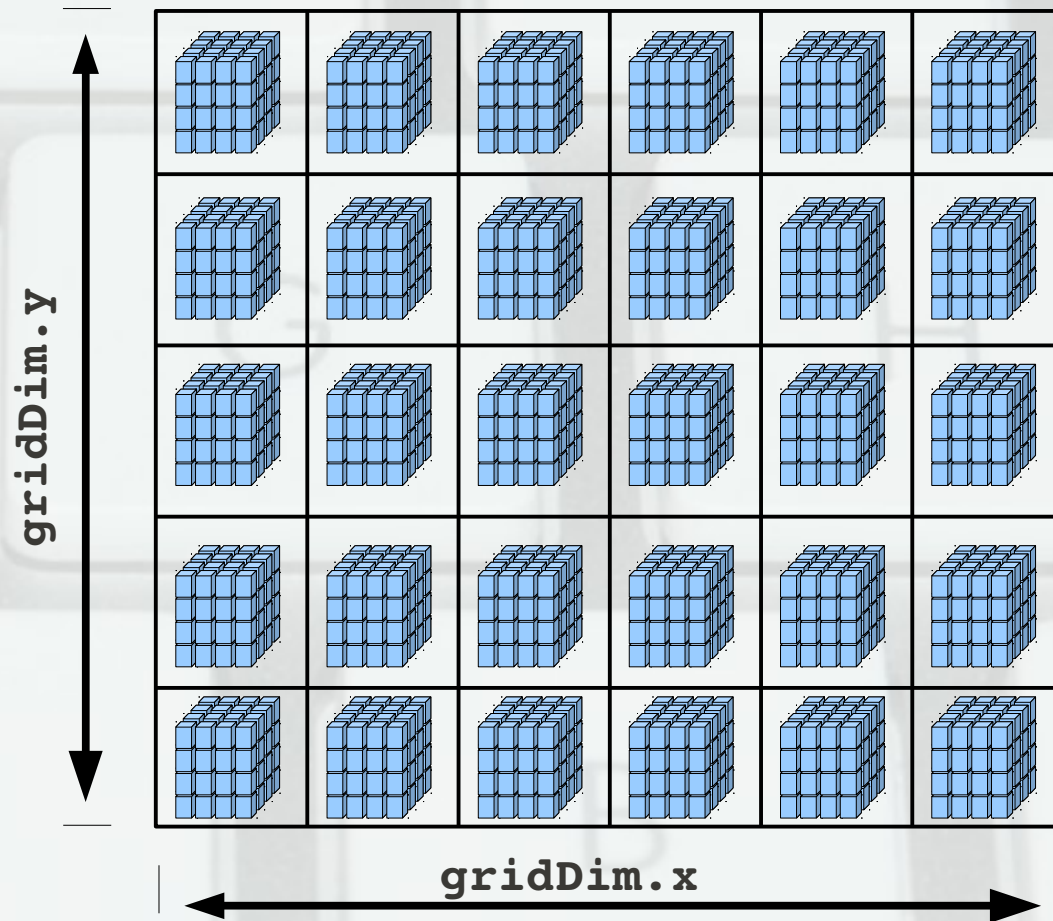
Threading Strategies

Grid of Thread Blocks:



Threading Strategies

Grid of Thread Blocks:



Threading Strategies

GPU Kernel – 1D Data

```
__global__ void  
kernel (float* a, float* b, float* c)  
{  
    int i = threadIdx.x;    // thread x coord  
  
    int block = blockIdx.x; // block x coord  
  
    int dim = blockDim.x;   // block x dim  
  
    int index = block*dim + i;  
  
    a[index] = b[index] + c[index];  
}
```

Threading Strategies

GPU Kernel – 2D Data

```
__global__ void
kernel (float* a, float* b, float* c)
{
    int i = threadIdx.x;  // thread x coord
    int j = threadIdx.y;  // thread y coord

    int block_size = blockDim.x * blockDim.y;

    int block_idx =
        blockIdx.y*gridDim.x + blockIdx.x;

    int offset = block_idx*block_size;

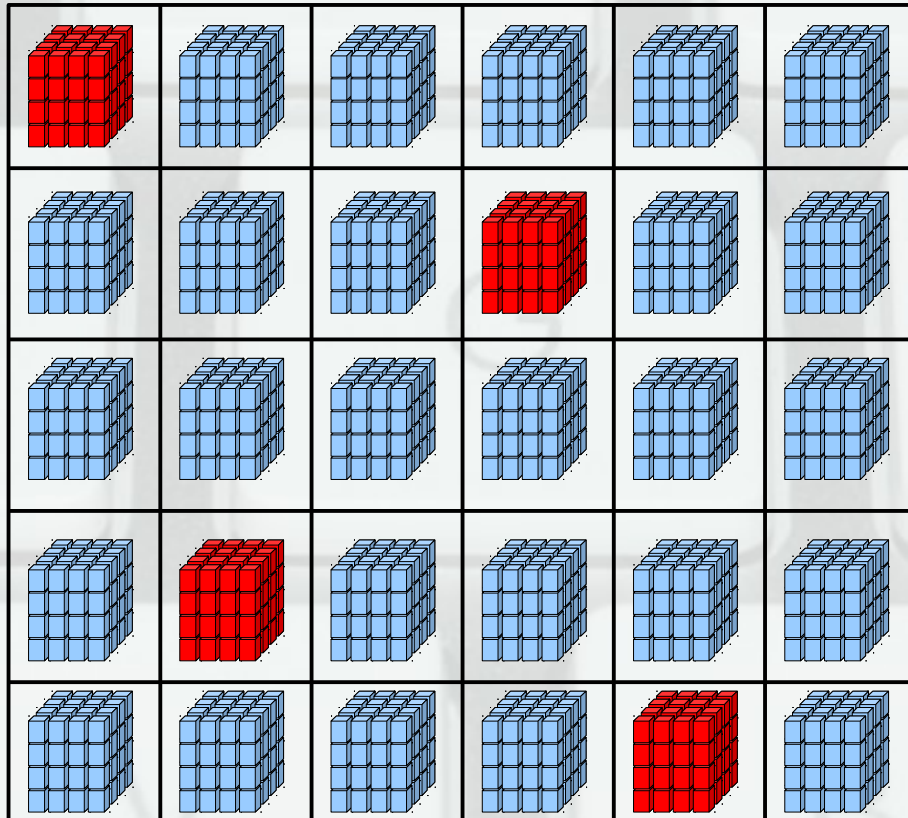
    int local_index = blockDim.x*j + i;

    int index = offset + local_index;

    a[index] = b[index] + c[index];
}
```


DATA PARALLELISM

Blocks may be assigned to SMs **out of order**. This is okay.



As a result, there can be no data dependencies between blocks.

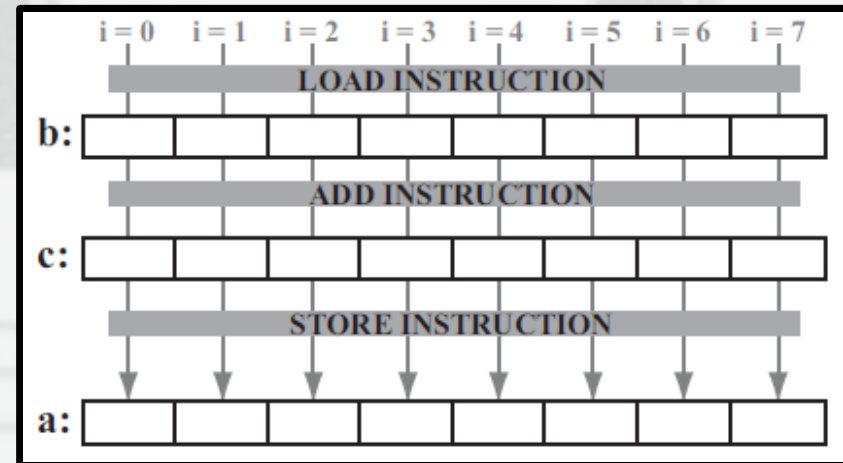
The Most Basic Example

GPU Kernel

```
__global__ void  
kernel (float* a, float* b, float* c)  
{  
    int i = threadIdx.x;  
    a[i] = b[i] + c[i];  
}
```

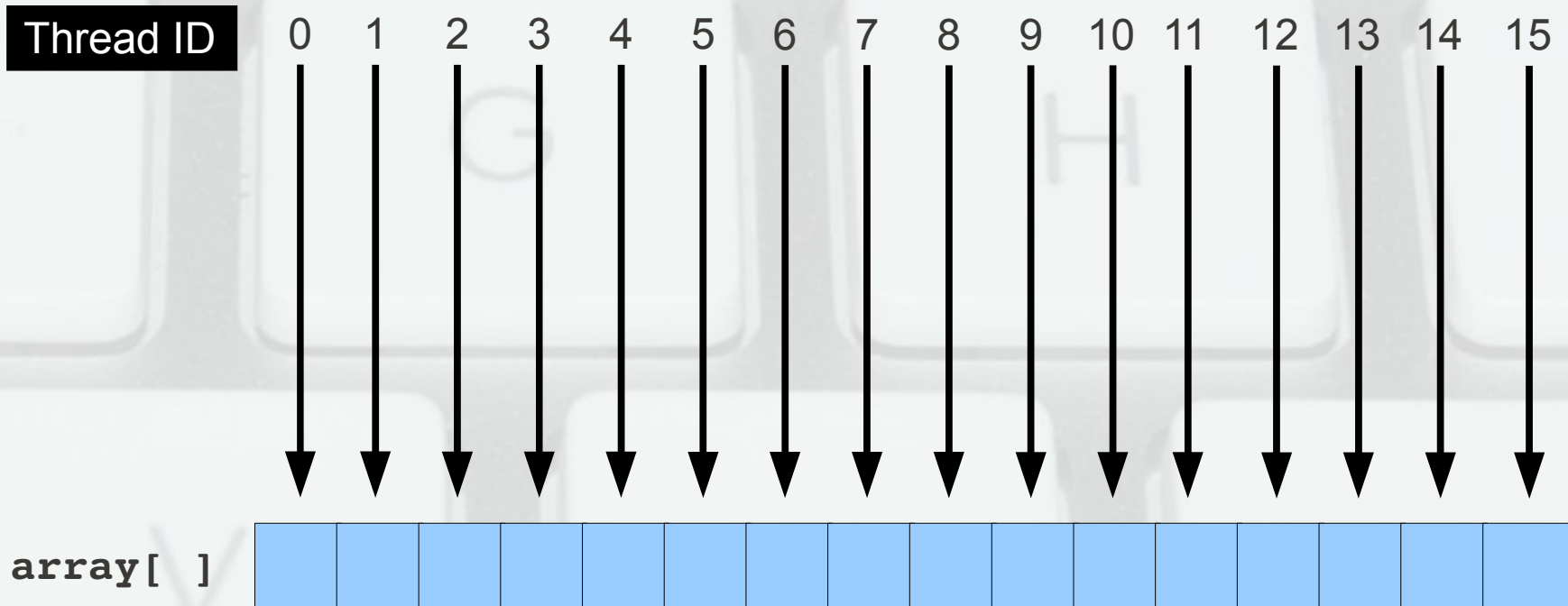
✓ Thread IDs used to index into arrays.

Data reads & writes are *coalesced*.



Data Coalescence

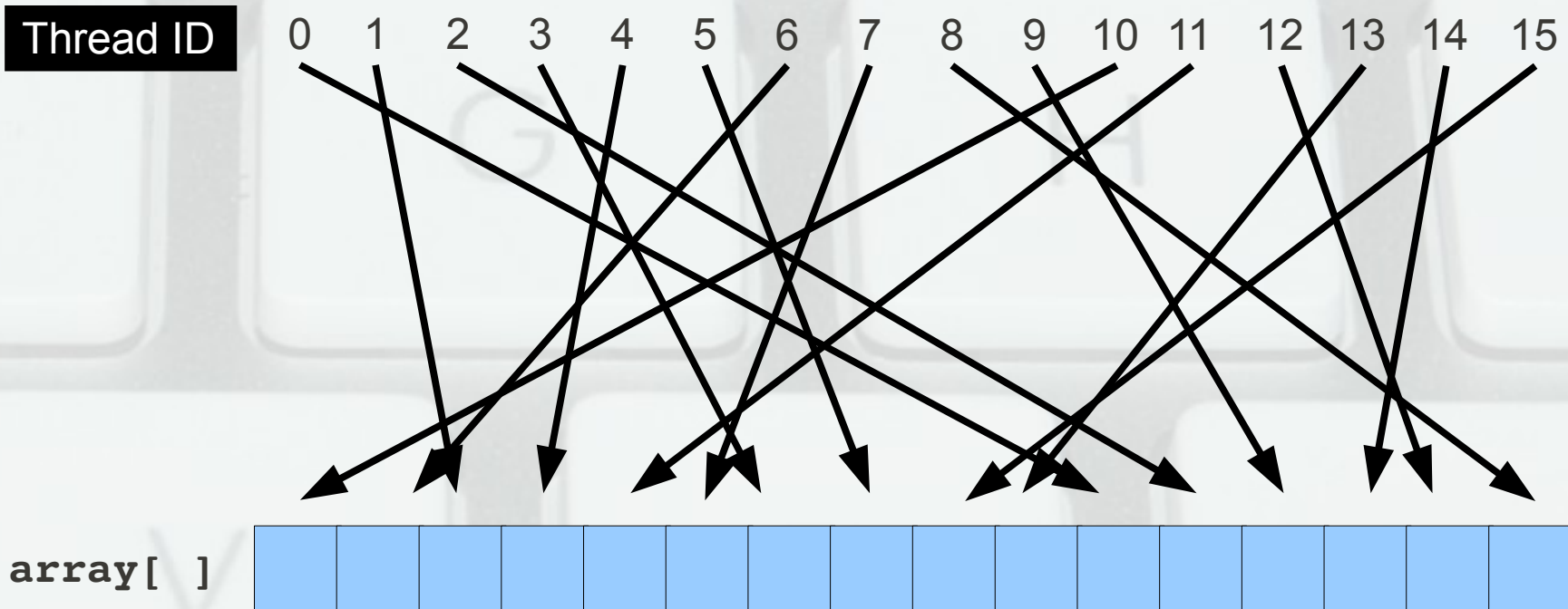
Good ✓



Data Coalescence

Bad X

(These writes get serialized)



Data Coalescence

Bad X

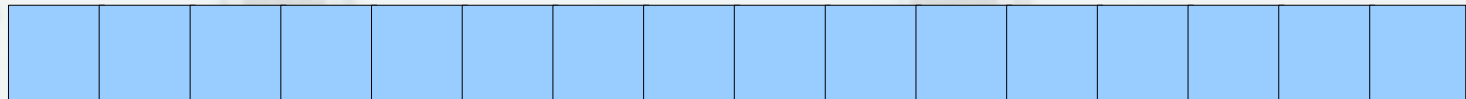
(These writes get serialized)

Thread ID

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

**Reads follow the
same rules!**

array[]



Essence of Structure

Thread IDs are used to index into data.

&

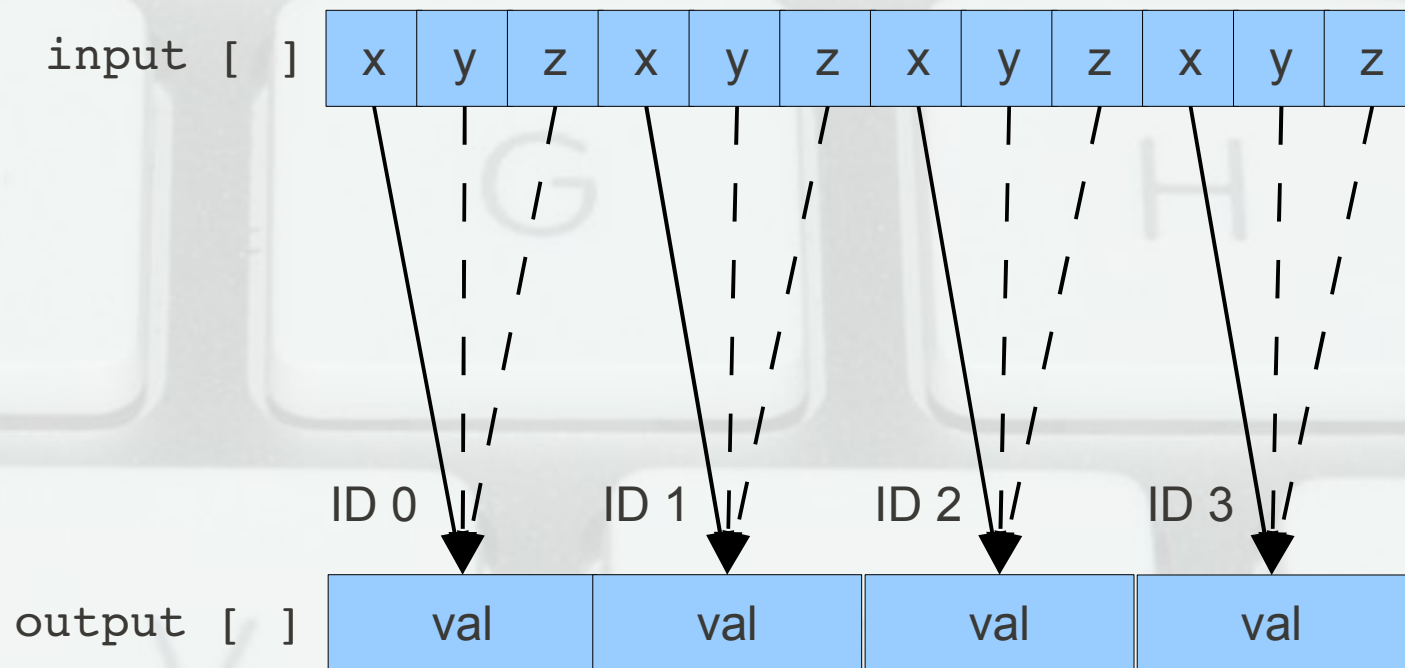
**Consecutive thread IDs must
access consecutive data.**

- **Threading scheme is heavily**
- **dependent on data structure.**

Essence of Structure

Tip 1: *Interleaving is probably a bad idea.*

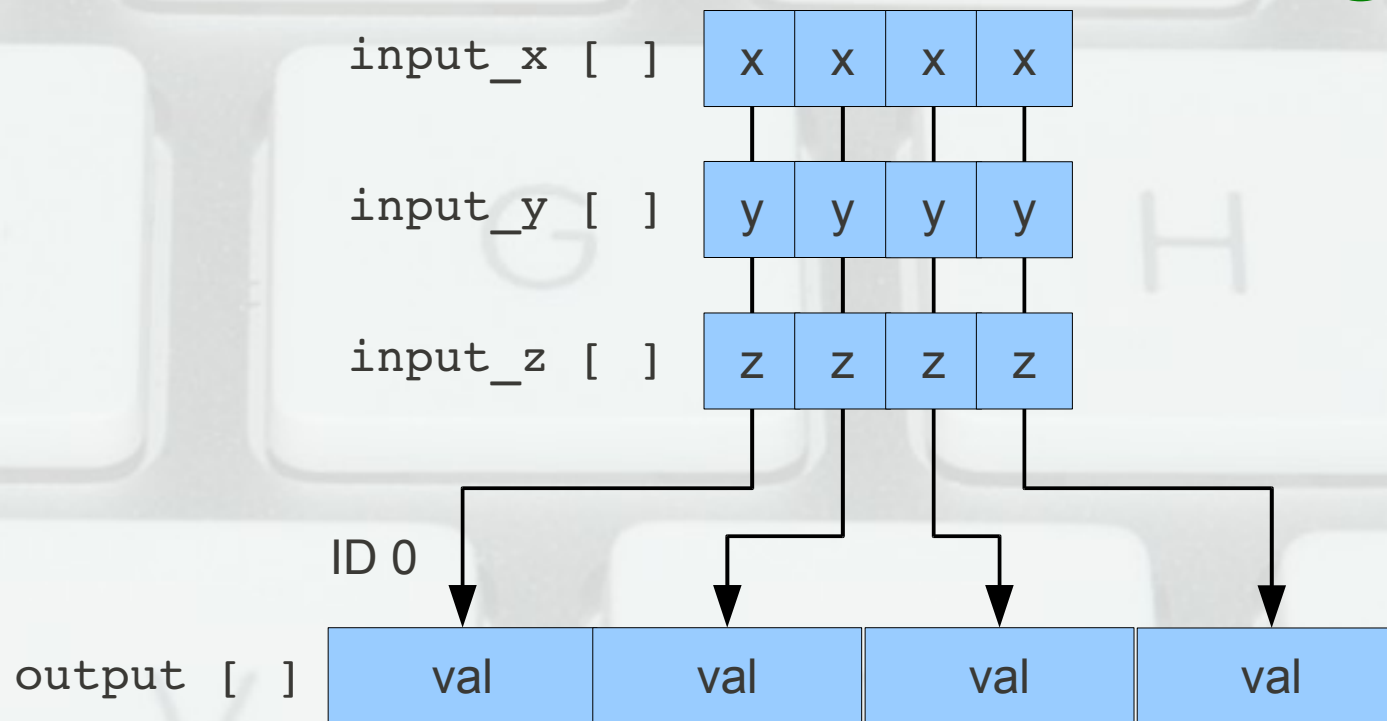
Bad 



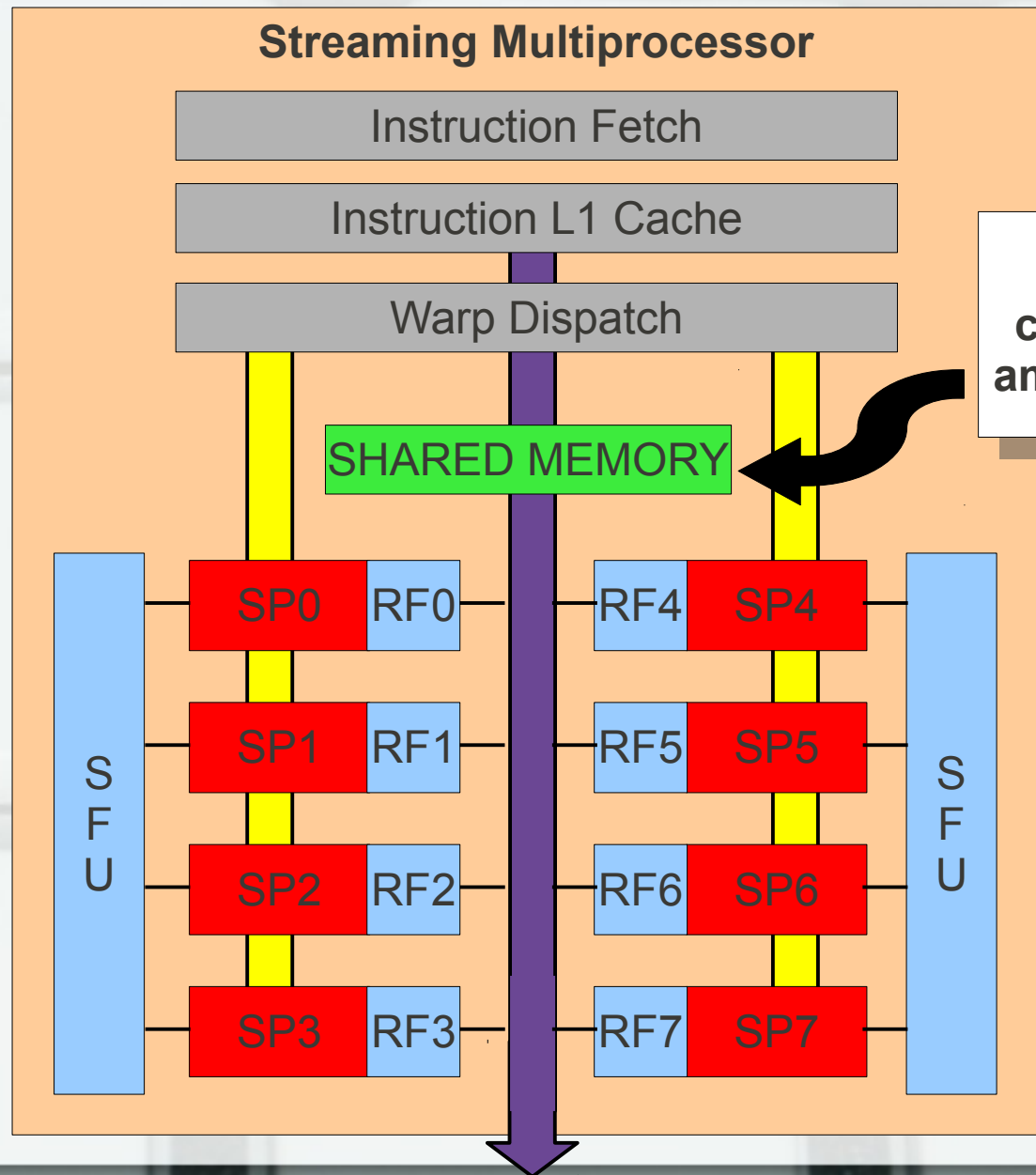
Essence of Structure

Tip 1: *Interleaving is probably a bad idea.*

Good ✓



SHARED MEMORY



Inter-block
communication
amongst threads.

TO GLOBAL
GPU MEMORY

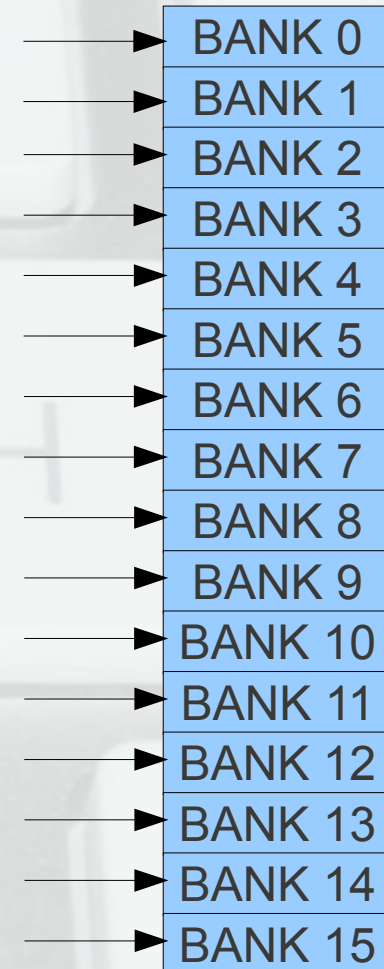
SHARED MEMORY

Each bank is 1KB
(16KB Shared Memory **TOTAL**)

Accessed by the **HALF WARP**
(16 threads at a time)

Each thread in a half warp should
access a different bank!
(Avoid **BANK CONFLICTS**)

MUCH FASTER THAN GLOBAL MEMORY

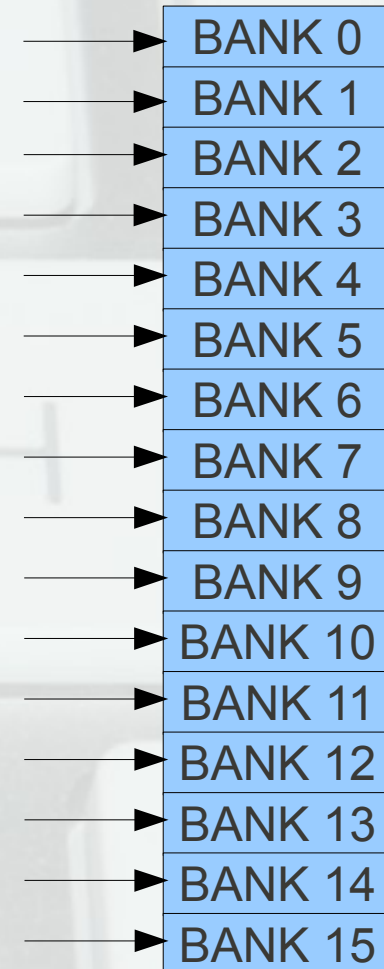


SHARED MEMORY

Can only be populated **DURING** kernel execution.

Has the thread block level **SCOPE**.

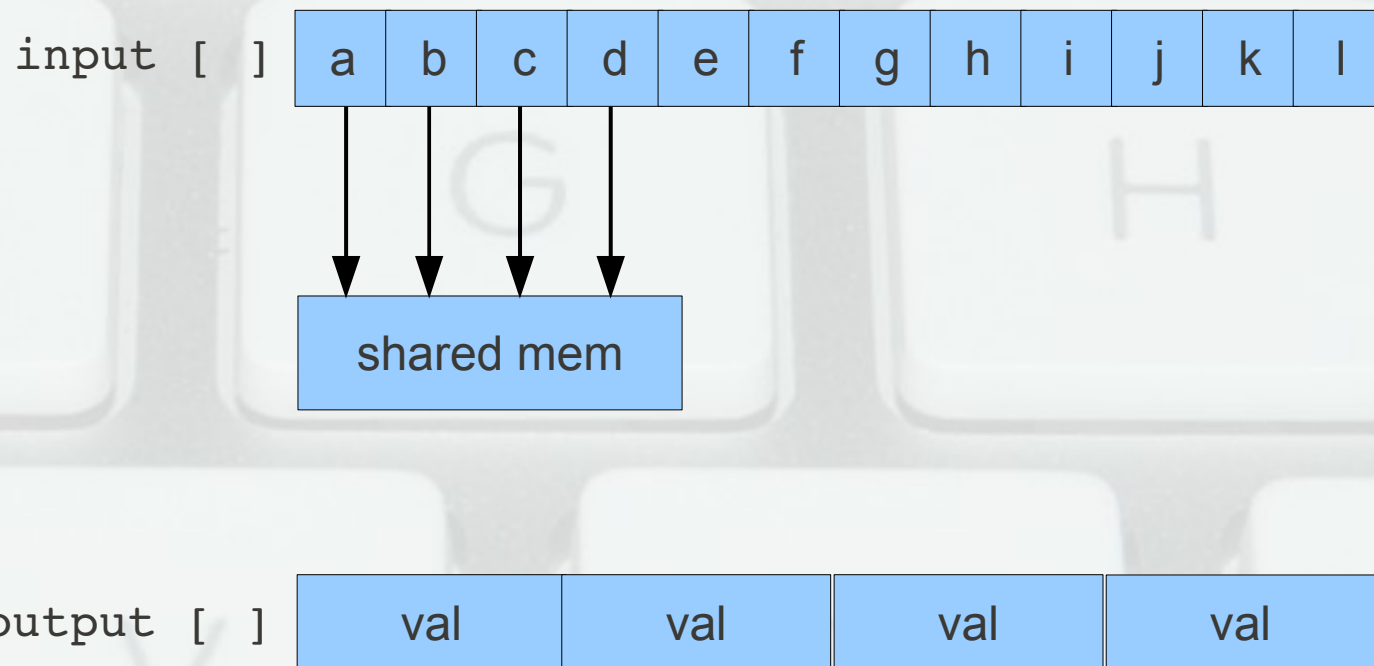
Amount of shared memory for each thread block is specified **AT LAUNCH**



Essence of Structure

Tip 2: Many to fewer: assign thread blocks, not threads.

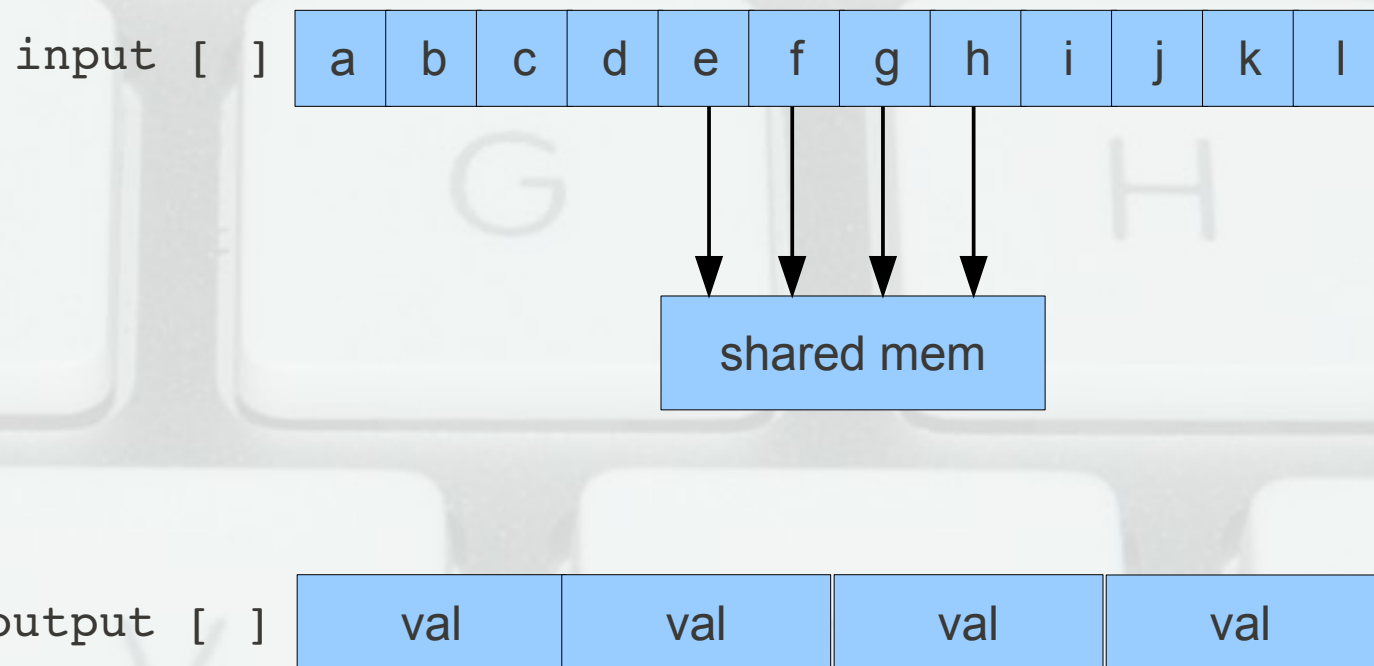
(example – 4 threads per block)



Essence of Structure

Tip 2: Many to fewer: assign thread blocks, not threads.

(example – 4 threads per block)



Essence of Structure

Tip 2: Many to fewer: assign thread blocks, not threads.

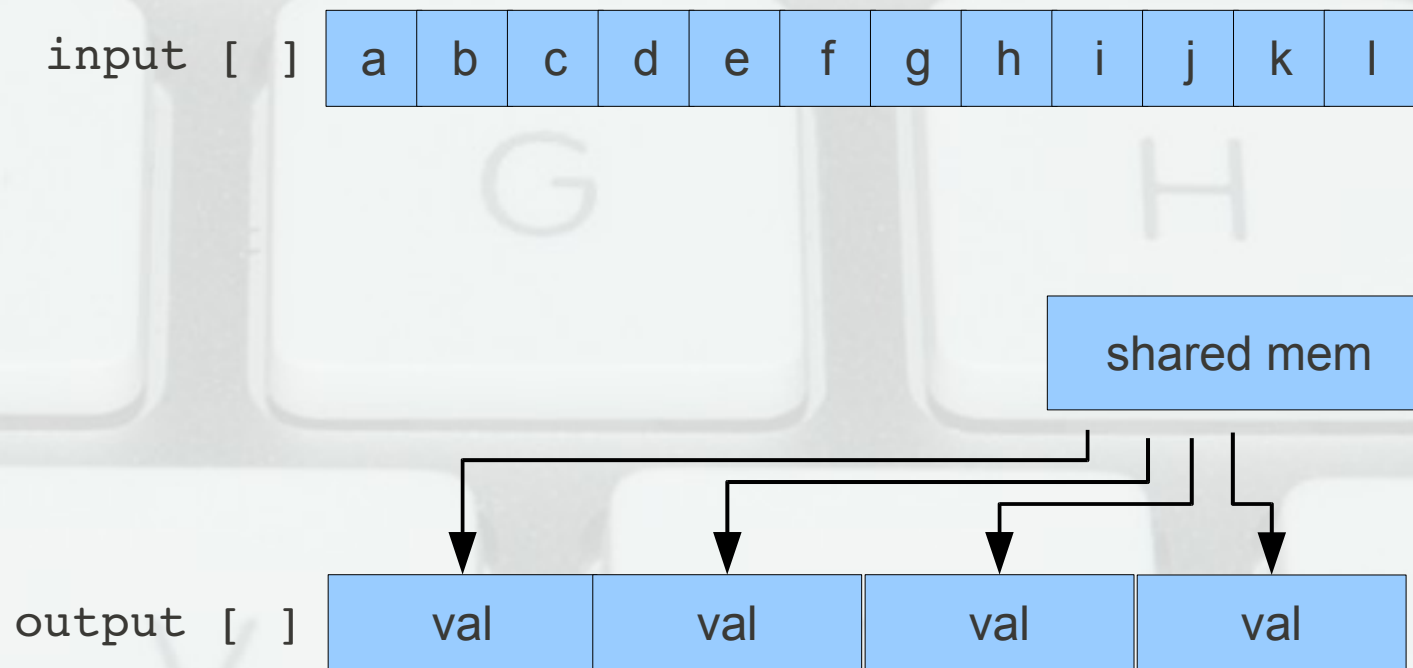
(example – 4 threads per block)



Essence of Structure

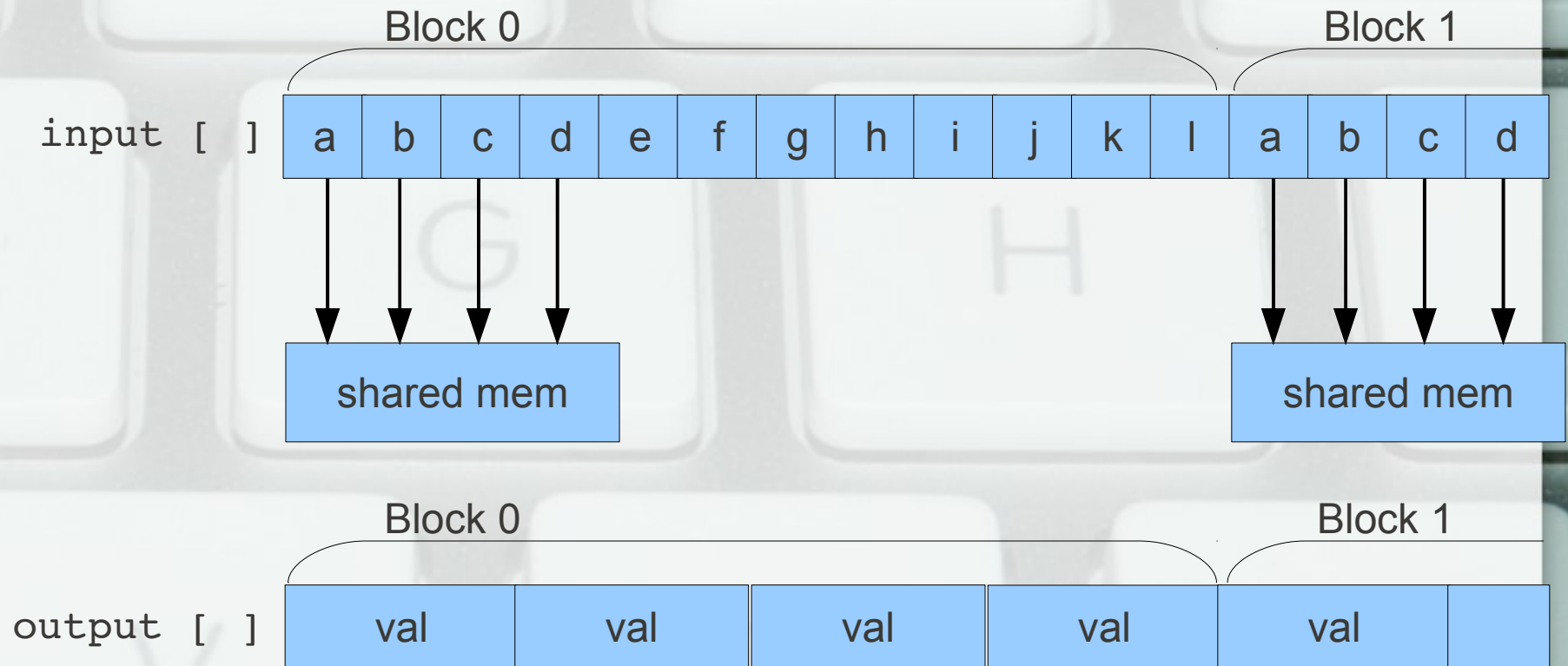
Tip 2: Many to fewer: assign thread blocks, not threads.

(example – 4 threads per block)



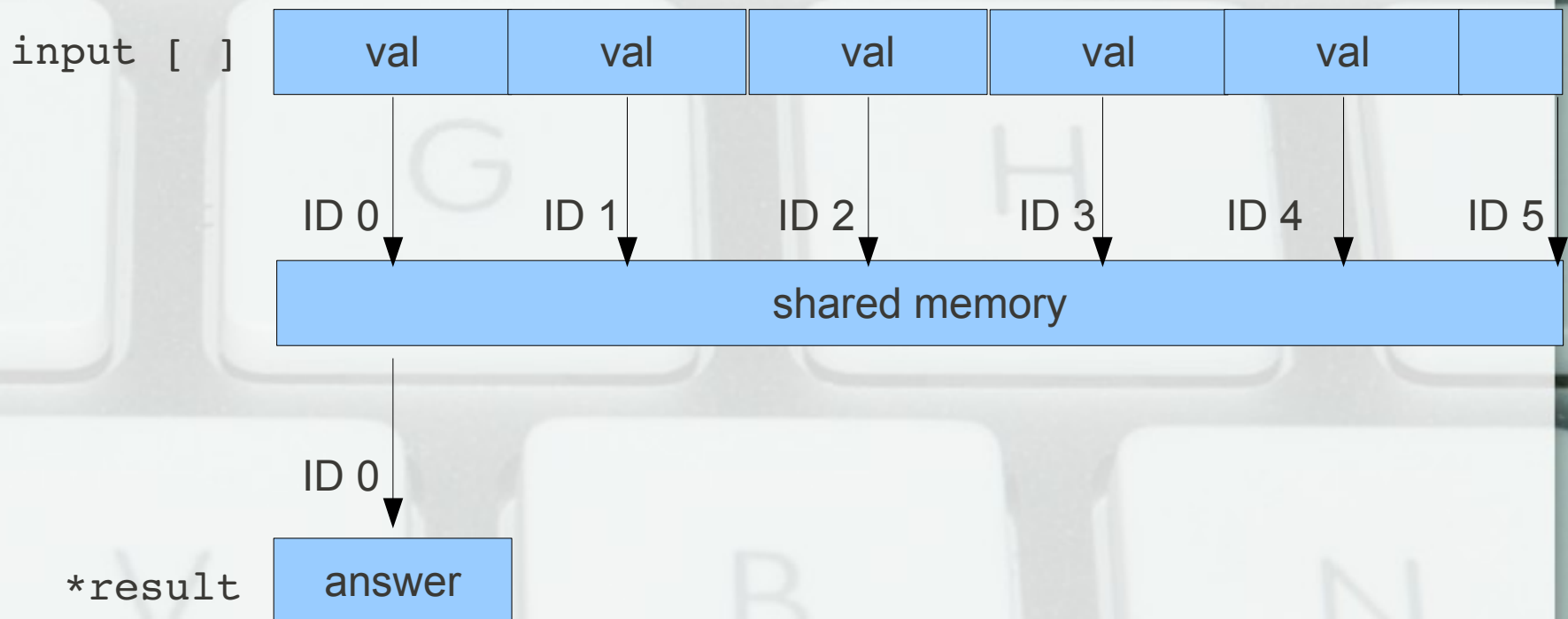
Essence of Structure

Tip 2: Many to fewer: assign thread blocks, not threads.



Essence of Structure

Tip 3: *Many to one: parallel sum reductions*



cuda-gdb command reference

`thread <<<(b.x, b.y), (t.x, t.y, t.z)>>>`

switch to specified CUDA thread

`break function OR file:line#`

set breakpoint. pauses execution at specified point.

`info cuda threads`

reports which line # each active thread is currently on

`info break`

list all breakpoints and show their numbers

`info cuda threads all`

reports line #s for all CUDA threads

`disable breakpoint#`

temporarily disable a breakpoint

`info cuda state`

reports memory allocation and symbols on GPU

`enable breakpoint#`

enable a previously disabled breakpoint

cuda-gdb command reference

`delete breakpoint#`

deletes specified breakpoint

`run`

begin program execution

`print variable`

displays the current contents of the specified variable

`step`

execute next line then break

`print array[start]@end`

displays a slice of an array from element *start* thru element *end*

`list`

show code listing around current breakpoint

`watch variable`

set watchpoint. program will break when *variable* is modified

`kill`

terminal program execution

`quit`

exit debugger