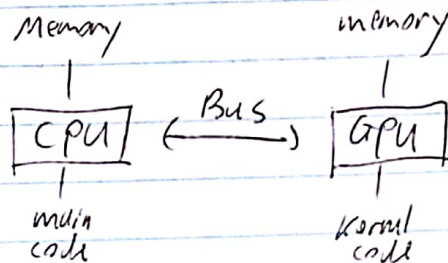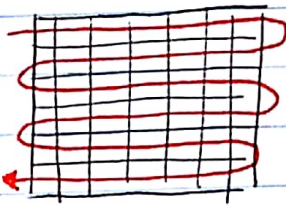# Learn CUDA in an Afternoon

## Introduction

- Why faster?
  many more cores, different kind of memory.
- But, GPUs cannot be used on their own. Not good for an OS, for example.
- GPUs accelerate computationally demanding sections of code (which we call _kernels_). kernels are decomposed to run in parallel on multiple cores.
- Separate memory spaces

```
   Memory              memory
     |                   |
   [CPU]  <-Bus->      [GPU]
     |                   |
   main               kernel
   code               code
```
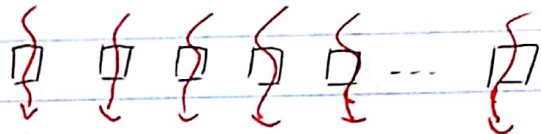
- CUDA is an extension to C/C++ that allows programming GPUs.
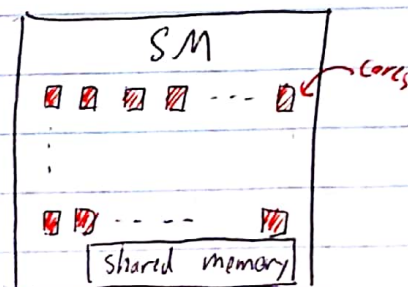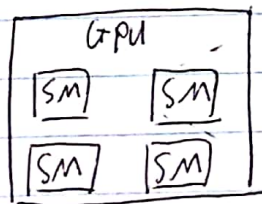
## Stream computing

serial or sequential                    _Parallel_



## Hardware



SMS - streaming multiprocessors. The number of SMs, and cores per SM _varries_ across generation.

- ⟹ We need an abstraction that will perform well across different generations of GPUs.
- This is abstracted as a grid of thread blocks.
- Each block in a grid contains multiple threads, mapping onto the cores in an SM.
- So, the mapping is

$$block \longrightarrow SM$$
$$thread \longrightarrow core\ in\ an\ SM$$

- Key thing is that we don't need to know the details of the hardware.
- Instead, we oversubscribe, and the system will perform the scheduling automatically. Use more blocks than SMs, and more threads than cores. This way everything is kept as busy as possible, giving the best possible performance.

## CUDA
- dim3 type:

$$dim3\ \ my\_xyz\_values\ (Xvalue,\ yvalue,\ zvalue);$$

- Access with .x (as usual).

## Hotel Example

- Serial solution

```
for( i=0; i<N; i++) {
    result [i] = 2*i;
}
```

- parallel ①

```
__global__   void mykernel (int *result)
α
    int i = threadIdx.x;  ②
    result [i] = 2*i;
}
```

① specify that this is a kernel

② internal var unique to each thread in a block. It's a dim3 type but since our problem is 1D, we don't use .y and .z.

- Launching the kernel from the CPU:

```
dim3    blocks Per Grid (1, 1, 1) ;      // 1 block
dim3    threads Per Block (N, 1, 1) ;    // N threads

my Kernel <<<< blocks Per Grid, threads PerBlock >>>> (result) ;
```

- The above example won't be fast, since it only uses 1 block (which maps to __1__ SM). To use multiple blocks:

```
__global__    void my Kernel ( int  x result)
{
    int  i = blockIdx.x * blockDim.x + thread Idx.x ;
    result [i] = 2 * x ;
}

---
dim3 blocks Per Grid (N/256, 1, 1) ;    // assuming N % 256 == 0
dim3 threads Per Block (256, 1, 1) ;

my Kernel <<<< blocks Per Grid, threads PerBlock >>>> (result) ;
...
```

- We have chosen to use 256 threads per block, which is typically a good number (why?)

- Vector addition example:

```
__global__    add Vectors ( float * a,  float * b,  float * c)  {
    int  i = block Idx.x * block Dim.x + thread Idx.x ;
    c[i] = a[i] + b[i] ;
}
---
dim3 bpg (N/256, 1, 1) ;
dim3 tpb (256, 1, 1) ;
add Vectors <<<< bpg, tpb >>>> (a, b, c) ;
```

## 2D Example

```
__global__ void matrixAdd( float a[N][N], float b[N][N],
                           float c[N][N] ) {

    int j = blockIdx.x * blockDim.x + threadIdx.x;
    int i = blockIdx.y * blockDim.y + threadIdx.y;

    c[i][j] = a[i][j] + b[i][j];
}

int main() {
    dim3 bpg (N/16, N/16, 1)
    dim3 tpb (16, 16, 1);
    matrixAdd <<<< bpg, tpb >>>> (a, b, c);
}
```
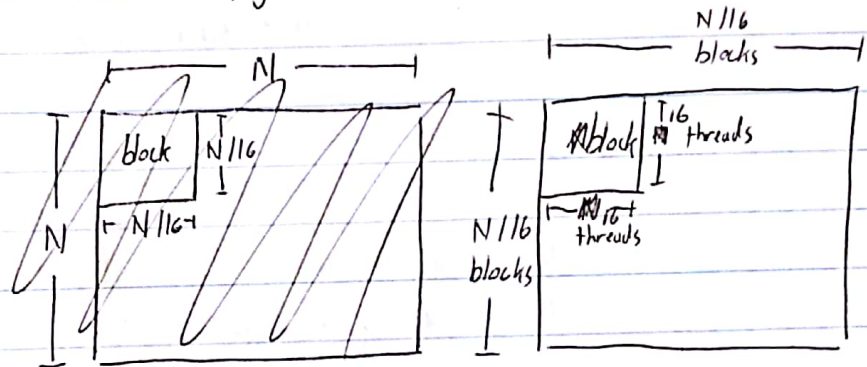


## Memory Management

- Stuff in the kernel has to point at GPU memory.
- Alloc and free memory on the GPU with:

```
float *a;
cudaMalloc (&a, N*N sizeof (float));

...

cudaFree (a);
```

- Copying memory

destination     Source

```
cudaMemcpy (array_device, array_host, N*sizeof (float),
            cudaMemcpy Host To Device );

cudaMemcpy (array_device, array_host, N*sizeof (float),
            cudaMemcpy Device To Host );
```

* transfers are relatively slow!

$$\begin{pmatrix} device - GPU \\ host - CPU \end{pmatrix}$$

## Sync Between Host and Device

- kernel calls are non-blocking (i.e., host program continues after it calls the kernel).
- Use cuda Thread Synchronize() to wait for kernel to finish
- Standard cuda Memcpy Calls are blocking (non-blocking variants exist)
- It's good practice to just put the cuda Thread Synchronize().

## Sync Between CUDA Threads

- To sync between threads in a block, use sync threads()
- Example: Communicate a value between threads. (Assume x is local and array is shared.)

```
if (thread Idx.x == 0)  array[0] = x;
Sync threads();
if (thread Idx.x == 1)  x = array[0];
```

- It is not possible to communicate between different blocks. Must instead exit kernel and start a new one.

## GPU Performance Inhibitors.

- Copying data to/from device
- Device under-utilisation / GPU memory latency
- GPU memory bandwidth
- Code branching

## Host - Device Data Copy

- Copying from host to device is expensive. We want to minimize these copies.
- keep data resident on device. May require importing routines to device, even if they are not computationally expensive.
- Might be quicker to calculate something from scratch on device instead of copying.

# Data Copy Optimization Example

```
                    Loop over timesteps
                        inexpensive_routine_on_host (data_on_host)
                        copy_ data from host to device
                        expensive_routine_on_device (data on_device)
                        copy data from device to host
                    end
```

- plis slor!
  port inexpensive routine to device and move data copies
  outside of loop.

```
                    Copy data from host to device
                    Loop over timesteps
                        inexpensive_routine_on_device (data_on_device)
                        expensive_routine_on_device (data_on_device)
                    end
                    copy data from device to host
```

## Exposing parallelism

- GPU performance relies on the usage of many threads.
- If a lot of code remains serial, effectiveness of GPU will
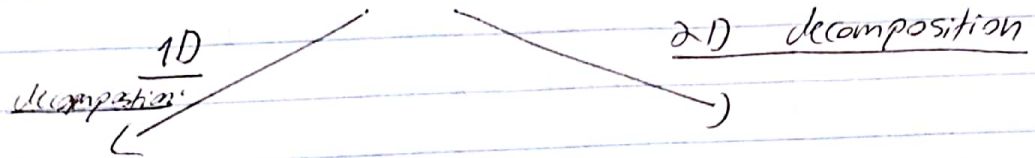  be limited (Amdahl's law)

## Occupancy and Memory Latency Hiding

- Decompose loops to threads
- For best performance, # threads >> # cores
- Accesses to GPU memory have several hundred cycles latency.
  When a thread stalls waiting for data, if another thread can
  switch in this latency can be hidden.
- NVIDIA GPUs have very fast thread switching.

## Example

original code

- Loop over i from 1 to 512
  - Loop over J from 1 to 512
    - independent iteration

**1D decomposition**

- Calc i from thread/block ID
  - Loop over J from 1 to 512
    - independent iteration

X  512 threads

**2D decomposition**

- Calc i & J from thread/block ID
- independent iteration

V  262,144 threads

## Memory coalescing

- GPUs have high peak memory bandwidth. But this is achieved only when data is accessed for multiple threads in a single transaction: memory coalescing.
- To achieve this, <u>ensure that consecutive threads access consecutive memory location</u> ❗
- otherwise, memory accesses are serialised, significantly degrading performance.

## Example

- Consecutive threads are those with consecutive thread Idx. X values.
- Do consecutive threads access consecutive memory locations?

  ✓  index = blockIdx.X * blockDim.X + threadIdx.X;
      output[index] = 2 * input[index];

## Example

In C, outermost index runs fastest: J here

|   |   |
|---|---|
| X | `i = blockIdx.x * blockDim.x + threadIdx.x;` |
| not coalesced. | `for ( J=0; J<N; J++)` |
| Consecutive | `    output[i][J] = 2 * input[i][J];` |
| threadIdx.x | |

Corresponds to Consecutive
i values!

```
J = blockIdx.x * blockDim.x + threadIdx.x;
for (i=0; i<N; i++)
    output[i][J] = 2 * input[i][J];
```