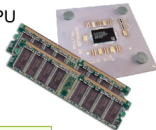# Data Migration (cont.)

Parallel
↳ Kernel Programming.

Host = CPU + RAM

Host Programming.

Serial
↳ CPU
↳ activity for GPU
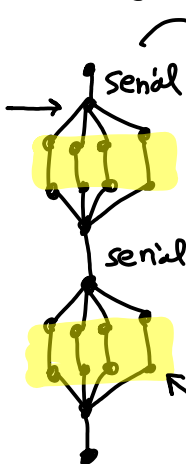↳ Copy data
↳ execute gpu

Host = GPU + RAM
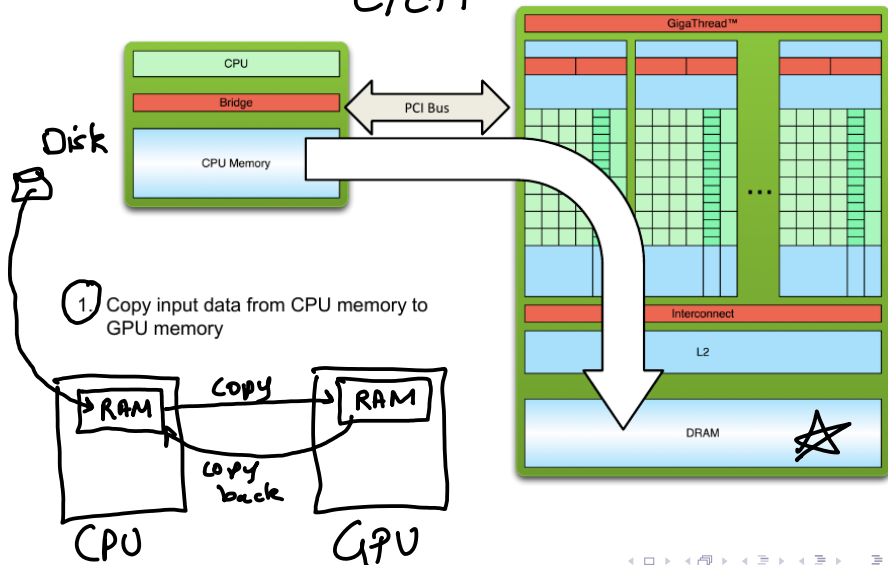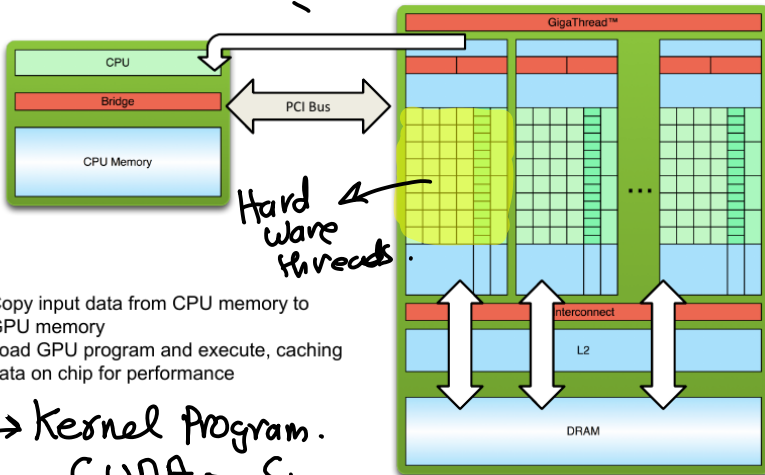
serial

parallel code

serial code

parallel code
serial code

serial
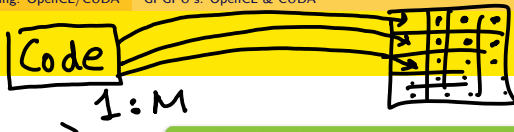


Figure 4: Porting portions of your code to GPU

# Data Migration (cont.)

Serial Programming.

C/C++



Disk

1 Copy input data from CPU memory to GPU memory

CPU

RAM — copy → RAM

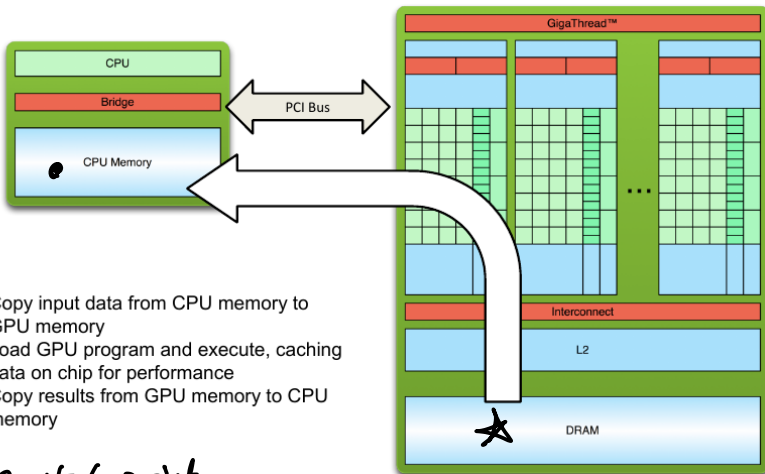copy back

GPU

# Data Migration (cont.)



1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance

→ Kernel Program.
CUDA — C

# Data Migration (cont.)



1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance
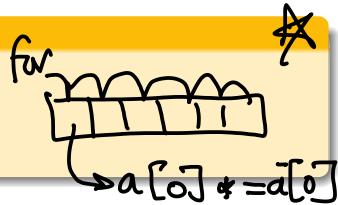3. Copy results from GPU memory to CPU memory

cout / print.

## Squaring an Array

*kernel.* (handwritten)

### C Code

```
void squareCPU(float *a, int size) {
    for (int x = 0; x < size; x++)
        a[x] *= a[x];
}
```

*for* (handwritten annotation)

*→ a[0] *= a[0]* (handwritten)

### OpenCL Kernel Code

```
__kernel void squareCPU(__global float *a) {
    int x = get_global_id(0);
    a[x] *= a[x];
}
```

### CUDA Kernel Code

*flag.* (handwritten annotation)

```
__global void hello(float *a) {
    int x = threadIdx.x;
    a[x] *= a[x];
}
```

*for* ⓧ (handwritten)

*x = thread id.* (handwritten)

*Object . member;* (handwritten annotation)

*a[0] *= a[0]* (handwritten)

code

C → Serial → CPU

flag
Parallel. → GPU

Compiler. (Help needed)

# Building up towards Hello World

Code

CPU

```
int main()
{
    printf("hello world!\n");
    return 0;
}
```

GPU

## Compilation

- nvcc hello_world.cu
- ./a.out

## Building up towards Hello World: Inserting GPU Code

*serial*

*function call.* ←

```
__global__ void myKernel(void)
{
}
```

*retrieve.*

```
int main()
{
    myKernel<<<1,1>>>();
```

80 % ← ⟨⟨⟨1 , 1 ⟩⟩⟩

```
    printf("hello world!\n");
    return 0;
}
```

// indicates function runs on *Parallel*
// device and is called from
// host code

*ftn name();*

// like function call, but
// with more information
// 1st digit number of blocks
// 2nd digit threads per block

- nvcc will separate source code into host and device components
  - Device functions processed by NVIDIA compiler
  - Host functions processed by standard host compiler

# Building up towards Hello World: Inserting GPU Code

*(handwritten annotations at top right: crossed-out "copy", circled "2" with "exec", crossed-out "copy")*

```
__global__ void myKernel(void)      // indicates function runs on
{                                    // device and is called from
}   empty.                          // host code


int main()
{              serial
    myKernel<<<1,1>>>();    ②        // like function call, but
    1 thread                         // with more information
         serial.                     // 1st digit number of blocks
                                     // 2nd digit threads per block
    printf("hello world!\n");
    return 0;
}
```
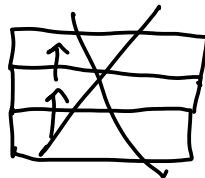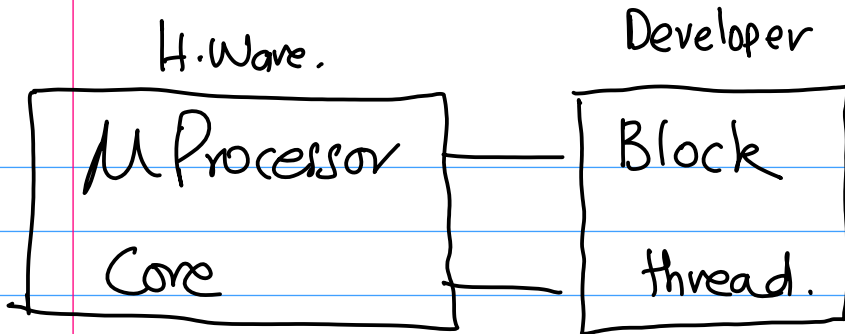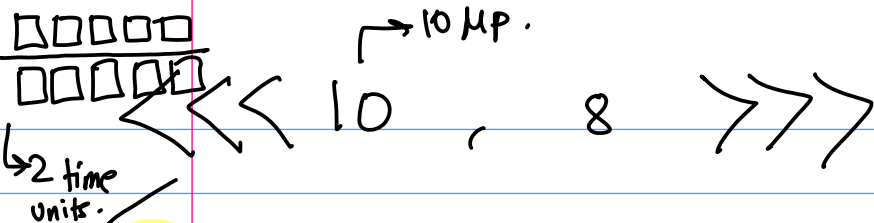
- nvcc will separate source code into host and device components
  - Device functions processed by NVIDIA compiler
  - Host functions processed by standard host compiler

H.Ware.                    Developer

| MProcessor | — | Block |
| Core | — | thread. |

1 block + 1 thread.

↳ 1 MP. ├ 1 core.

5 MP $\longrightarrow$ 8 Cores.   $\Rightarrow$ $8 \times 5 = 40$ cores.

□ □ □ □ □
□ □ □ □ □  $\ll\!\!\langle$ 10 , 8 $\rangle\!\!\gg$

$\rightarrow$ 10 MP.

$\rightarrow$ 2 time units.

$MP_0$   $MP_1$   $MP_2$   $MP_3$   $MP_4$

| $C_0$ | $C_4$ |
|-------|-------|
| $C_1$ | $C_5$ |
| $C_2$ | $C_6$ |
| $C_3$ | $C_7$ |

parallel → executed somewhat sequentially (32)
1 MP.

[ ∫∫∫∫∫∫∫∫∫∫
_____

[ ∫∫∫∫∫∫∫∫∫∫
_____

4 time units
32 threads.

[ ∫∫∫∫∫∫∫∫∫∫
_____

[ ∫∫∫∫∫∫∫∫∫

int main( )

• ∫

    int a = 5

    int b = 6

①   int c = a + b;
②
③



Create
memory locations
on gpu.

① copy (a,b)

② exec

③ copy

CPU           GPU

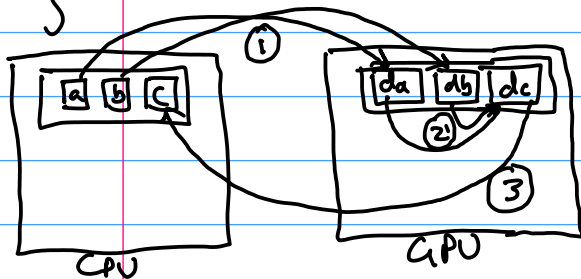a b c      da db dc

```
/* a, b, c are pointers to memory location on the device */
__global__ void add(int *a, int *b, int *c)
{ *c = *a + *b;
}
int main()
{ int a=2, b=7, c, *da, *db, *dc;

    cudaMalloc((void **)&da, sizeof(int)); // allocate
    cudaMalloc((void **)&db, sizeof(int)); // on device
    cudaMalloc((void **)&dc, sizeof(int));

    // Copying in blocking mode
    cudaMemcpy(da, &a, sizeof(int), cudaMemcpyHostToDevice);
    cudaMemcpy(db, &b, sizeof(int), cudaMemcpyHostToDevice);
    add<<<1,1>>>(da, db, dc);
    cudaMemcpy(&c, dc, sizeof(int), cudaMemcpyDeviceToHost);

    printf("a + b = %d\n", c);
    cudaFree(da); cudaFree(db); cudaFree(dc);
    return 0;
}
```

Handwritten annotations:

CPU

int *da = &a;

a  b  c
da  db  dc

GPU

destination → source

*c = *a + *b

```
/* a, b, c are pointers to memory location on the device */
__global__ void add(int *a, int *b, int *c)
{  *c = *a + *b;
}
int main()
{  int a=2, b=7, c, *da, *db, *dc;

   cudaMalloc((void **)&da, sizeof(int)); // allocate
   cudaMalloc((void **)&db, sizeof(int)); // on device
   cudaMalloc((void **)&dc, sizeof(int));

   // Copying in blocking-mode
   cudaMemcpy(da, &a, sizeof(int), cudaMemcpyHostToDevice);
   cudaMemcpy(da, &a, sizeof(int), cudaMemcpyHostToDevice);
   add<<<1,1>>>(da, db, dc);
   cudaMemcpy(&c, dc, sizeof(int), cudaMemcpyDeviceToHost);

   printf("a + b = %d\n", c);
   cudaFree(da); cudaFree(db); cudaFree(dc);
   return 0;
}
```

```
/* a, b, c are pointers to memory location on the device */
__global__ void add(int *a, int *b, int *c)
{   *c = *a + *b;
}
int main()
{   int a=2, b=7, c, *da, *db, *dc;

    cudaMalloc((void **)&da, sizeof(int)); // allocate
    cudaMalloc((void **)&db, sizeof(int)); // on device
    cudaMalloc((void **)&dc, sizeof(int));

    // Copying in blocking-mode
    cudaMemcpy(da, &a, sizeof(int), cudaMemcpyHostToDevice);
    cudaMemcpy(da, &a, sizeof(int), cudaMemcpyHostToDevice);
    add<<<1,1>>>(da, db, dc);
    cudaMemcpy(&c, dc, sizeof(int), cudaMemcpyDeviceToHost);

    printf("a + b = %d\n", c);
    cudaFree(da); cudaFree(db); cudaFree(dc);
    return 0;
}
```
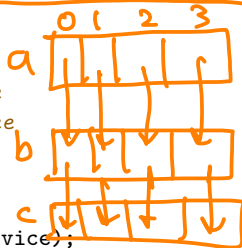
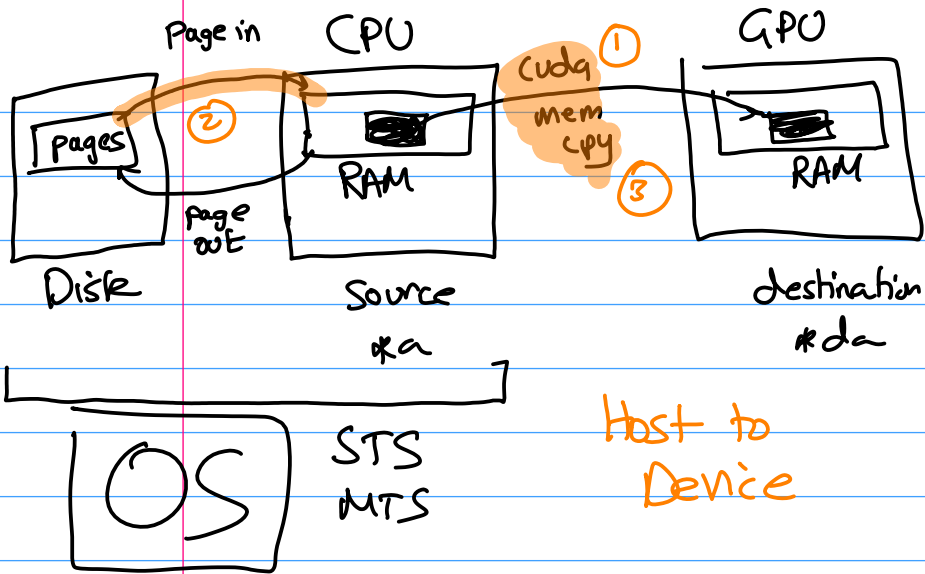$c[x] = a[x] + b[x]$

4 threads

<<< 1, 4 >>>

Page in     CPU                                    GPU



① cuda mem cpy

② 

③

| pages |

Disk          RAM          RAM

page out

Source          destination
*a              *da

STS
MTS

OS

Host to Device

```
/* a, b, c are pointers to memory location on the device */
__global__ void add(int *a, int *b, int *c)
{  *c = *a + *b;
}
int main()
{  int a=2, b=7, c, *da, *db, *dc;

    cudaMalloc((void **)&da, sizeof(int)); // allocate
    cudaMalloc((void **)&db, sizeof(int)); // on device
    cudaMalloc((void **)&dc, sizeof(int));

    // Copying in blocking-mode
    cudaMemcpy(da, &a, sizeof(int), cudaMemcpyHostToDevice);
    cudaMemcpy(da, &a, sizeof(int), cudaMemcpyHostToDevice);
    add<<1,1>>>(da, db, dc);
    cudaMemcpy(&c, dc, sizeof(int), cudaMemcpyDeviceToHost);

    printf("a + b = %d\n", c);
    cudaFree(da); cudaFree(db); cudaFree(dc);
    return 0;
}
```

# Running in Parallel

- so far; pointing parameters to GPU, and running single thread on GPU.
- Time to look at how to run things in parallel.

### Code Change

```
// add<<<1, 1>>>(da, db, dc);
   add<<<N, 1>>>(da, db, dc);
```

- Instead of running add() once, execute it N times in parallel
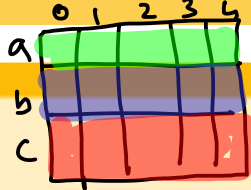
### Changes in Kernel Code        <<< N, 1 >>>

```
__global__ void add(int *a, int *b, int *c)
{
    c[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x];
}
```

## Changes in Host Code

```
int main()
{   int *a, *b, *c, *da, *db, *dc, N=5, i;

    a = (int*)malloc(sizeof(int)*N);  // allocate host mem
    b = (int*)malloc(sizeof(int)*N);  // and assign random
    c = (int*)malloc(sizeof(int)*N);  // memory

    cudaMalloc((void **)&da, sizeof(int)*N);
    cudaMalloc((void **)&db, sizeof(int)*N);
    cudaMalloc((void **)&dc, sizeof(int)*N);

    cudaMemcpy(da, &a, sizeof(int)*N, cudaMemcpyHostToDevice);
    cudaMemcpy(da, &a, sizeof(int)*N, cudaMemcpyHostToDevice);
    add<<<N,1>>>(da, db, dc);
    cudaMemcpy(&c, dc, sizeof(int)*N, cudaMemcpyDeviceToHost);
                 ⟶ N MP − 1 core.
    for (i = 0; i < N; i++)
        printf("a[%d] + b[%d] = %d\n", i, i, c[i]);
}
```
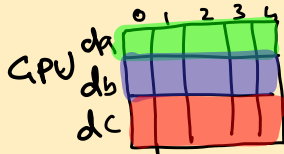
CPU

GPU

1 MP
<<< 1, N >>>
N cores

## That was Blocks in Parallel. What about Threads in Parallel

$$\langle\langle\langle N, 1 \rangle\rangle\rangle \longrightarrow c[blockIdx.x]$$

### Function Call Change

```
// add<<<1, 1>>>(da, db, dc); // single thread GPU
// add<<<N, 1>>>(da, db, dc); // N blocks on GPU
   add<<<1, N>>>(da, db, dc); // N threads on GPU
```

### Changes in Kernel Code

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

- Rest of Host code would be the same

$$c[x] = a[x] + b[x]$$

$$\langle\langle\langle 1, N \rangle\rangle\rangle \longrightarrow c[threadIdx.x]$$