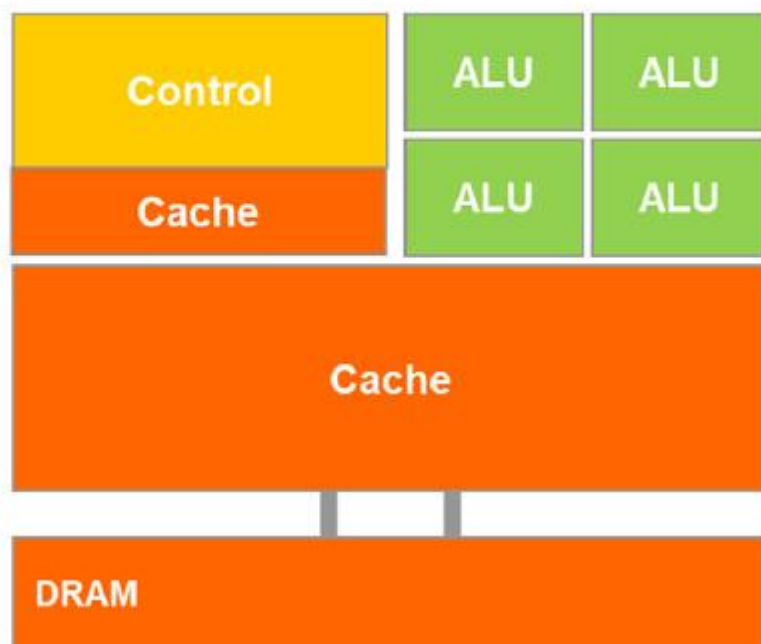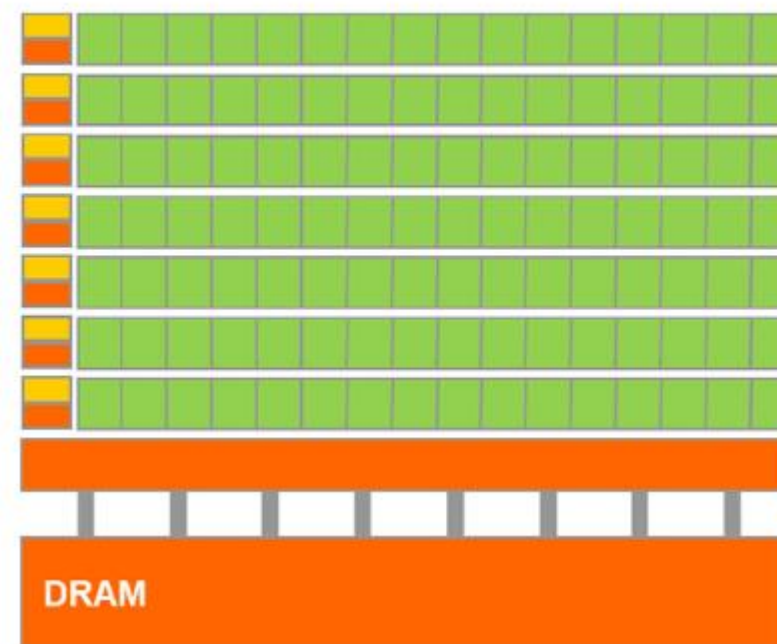# Parallel Programming

- Parallel programming is like getting help from friends to clean up a big mess. Instead of one person doing all the work, everyone works together at the same time. Each person cleans a small part, and when everyone finishes, the whole place is clean faster. This way, the job is done quicker because many people are helping at once.
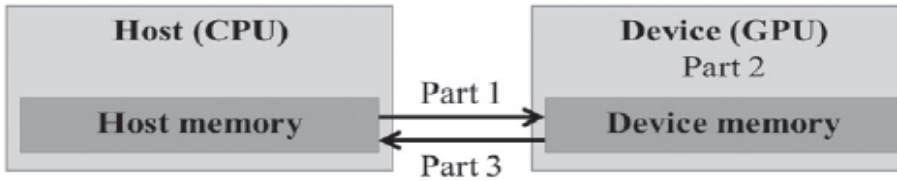
**FIGURE 1.1**

CPUs and GPUs have fundamentally different design philosophies: (A) CPU design is latency oriented; (B) GPU design is throughput-oriented.

```
01     void vecAdd(float* A, float* B, float* C, int n) {
02          int  size = n* sizeof(float);
03          float  *d_A *d_B, *d_C;
04
05          // Part 1: Allocate device memory for A, B, and C
06          // Copy A and B to device memory
07          ...
08
09          // Part 2: Call kernel – to launch a grid of threads
10          // to perform the actual vector addition
11          ...
12
13          // Part 3: Copy C from the device memory
14          // Free device vectors
15          ...
16     }
```

**FIGURE 2.5**

Outline of a revised vecAdd function that moves the work to a device.

```
01    void vecAdd(float* A_h, float* B_h, float* C_h, int n) {
02        int size = n * sizeof(float);
03        float *A_d, *B_d, *C_d;
04
05        cudaMalloc((void **) &A_d, size);
06        cudaMalloc((void **) &B_d, size);
07        cudaMalloc((void **) &C_d, size);
08
09        cudaMemcpy(A_d, A_h, size, cudaMemcpyHostToDevice);
10        cudaMemcpy(B_d, B_h, size, cudaMemcpyHostToDevice);
11
12        // Kernel invocation code - to be shown later
13        ...
14
15        cudaMemcpy(C_h, C_d, size, cudaMemcpyDeviceToHost);
16
17        cudaFree(A_d);
18        cudaFree(B_d);
19        cudaFree(C_d);
20    }
```

**FIGURE 2.8**

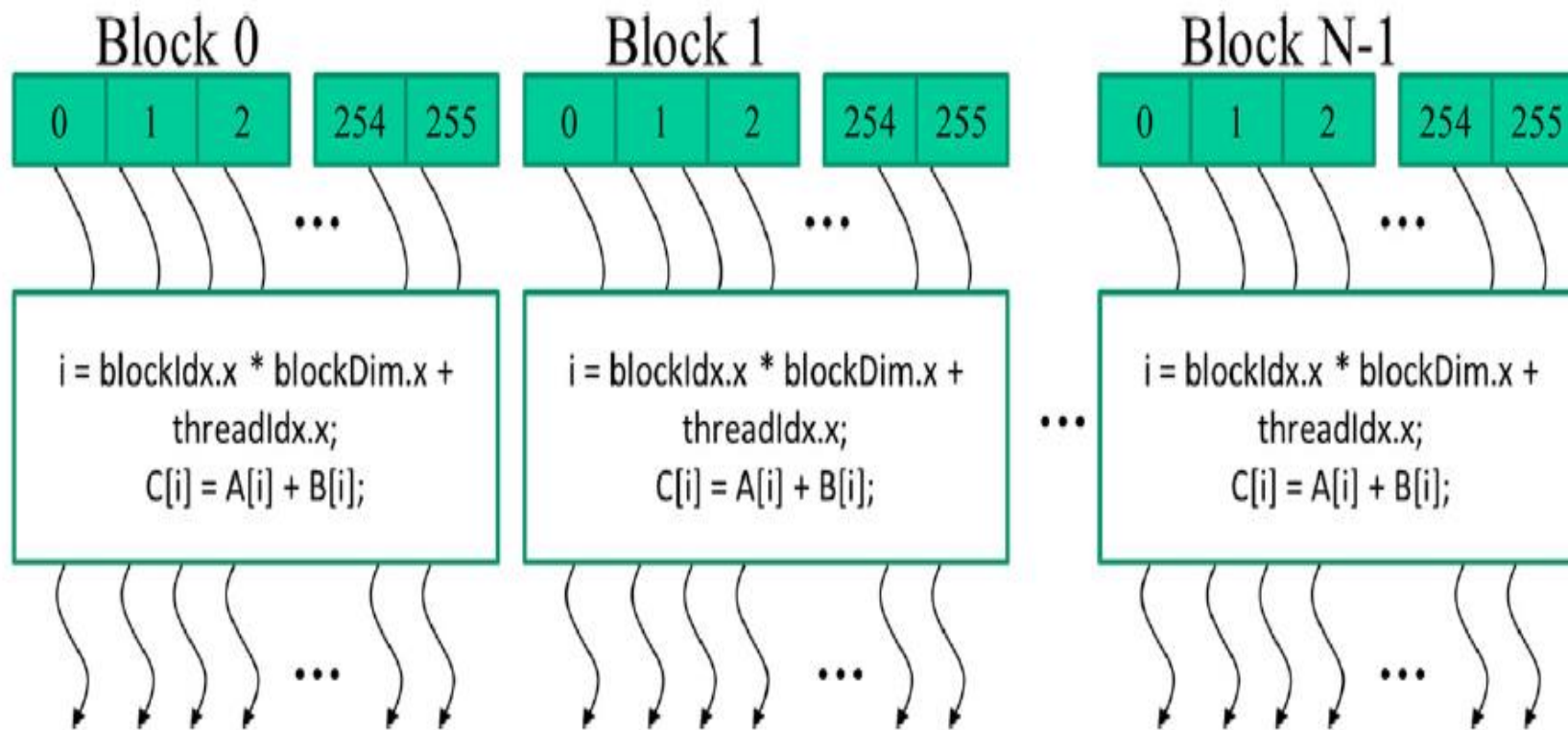A more complete version of vecAdd().

**FIGURE 2.9**

All threads in a grid execute the same kernel code.

```
01    // Compute vector sum C = A + B
02    // Each thread performs one pair-wise addition
03    __global__
04    void vecAddKernel(float* A, float* B, float* C, int n) {
05        int i = threadIdx.x + blockDim.x * blockIdx.x;
06        if (i < n) {
07            C[i] = A[i] + B[i];
08        }
09    }
```

**FIGURE 2.10**

A vector addition kernel function.

```
01    void vecAdd(float* A, float* B, float* C, int n) {
02        float *A_d, *B_d, *C_d;
03        int size = n * sizeof(float);
04
05        cudaMalloc((void **) &A_d, size);
06        cudaMalloc((void **) &B_d, size);
07        cudaMalloc((void **) &C_d, size);
08
09        cudaMemcpy(A_d, A, size, cudaMemcpyHostToDevice);
10        cudaMemcpy(B_d, B, size, cudaMemcpyHostToDevice);
11
12        vecAddKernel<<<ceil(n/256.0), 256>>>(A_d, B_d, C_d, n);
13
14        cudaMemcpy(C, C_d, size, cudaMemcpyDeviceToHost);
15
16        cudaFree(A_d);
17        cudaFree(B_d);
18        cudaFree(C_d);
19    }
```