# CHAPTER 2

# Heterogeneous data parallel computing

**FIGURE 2.1**

Conversion of a color image to a grayscale image.
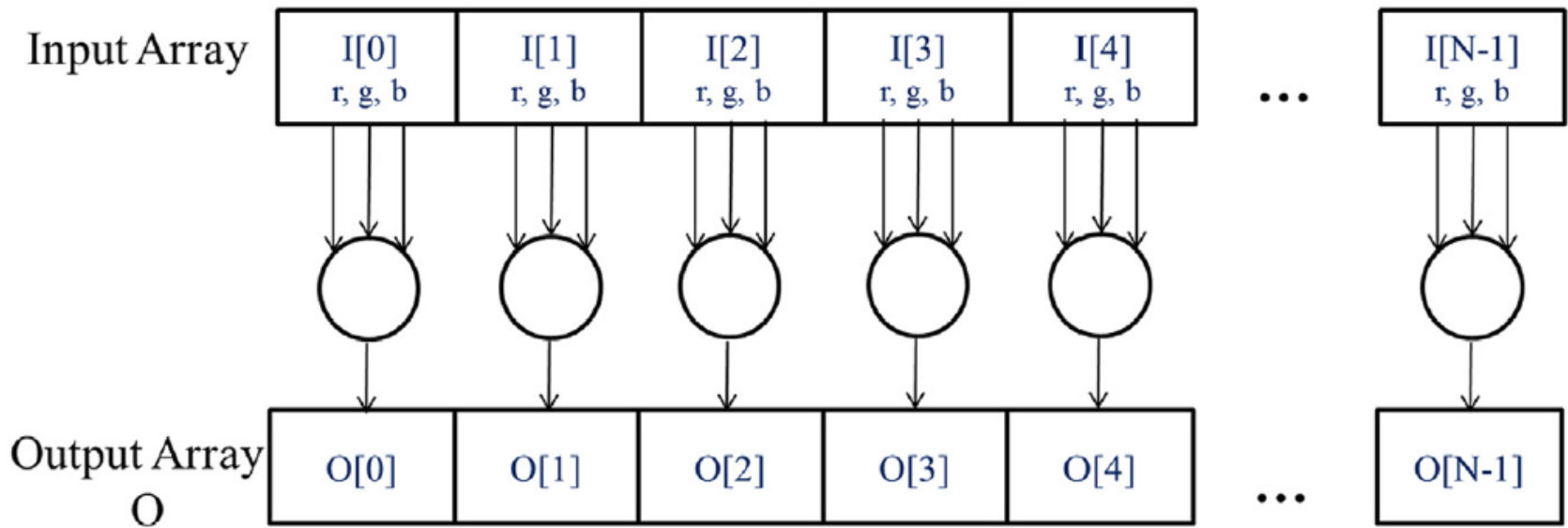
**FIGURE 2.2**

Data parallelism in image-to-grayscale conversion. Pixels can be calculated independently of each other.

**CPU Serial Code**

**Device Parallel Kernel**

`KernelA<<< nBlk, nTid >>>(args);`

**CPU Serial Code**

**Device Parallel Kernel**

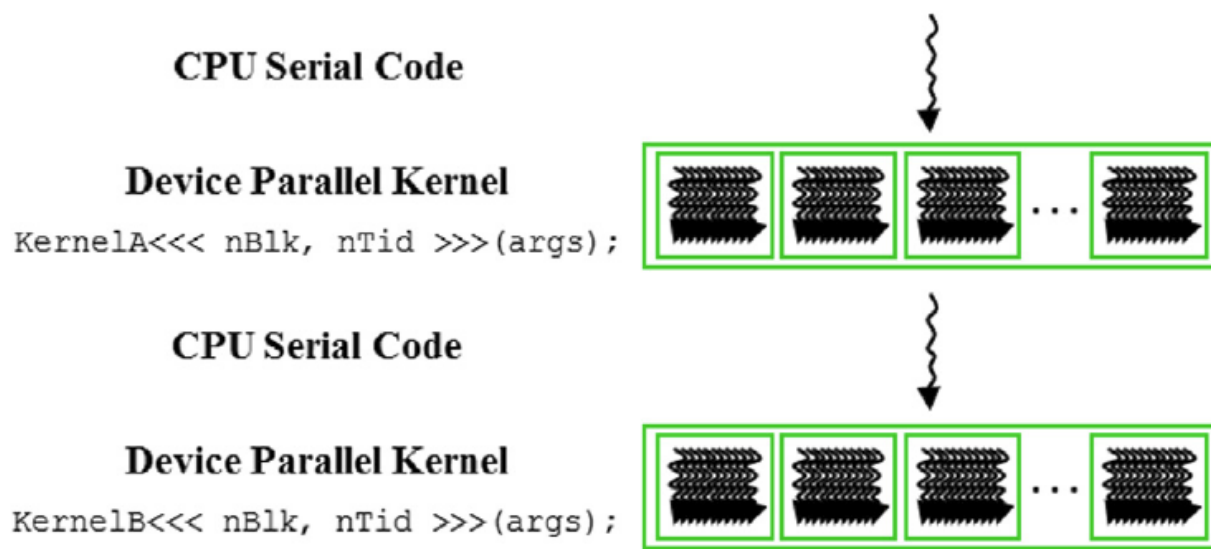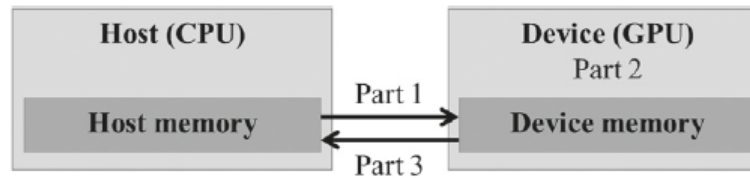`KernelB<<< nBlk, nTid >>>(args);`

**FIGURE 2.3**

Execution of a CUDA program.

```
01    // Compute vector sum C_h = A_h + B_h
02    void vecAdd(float* A_h, float* B_h, float* C_h, int n) {
03        for (int i = 0; i < n; ++i) {
04            C_h[i] = A_h[i] + B_h[i];
05        }
06    }
07    int main() {
08        // Memory allocation for arrays A, B, and C
09        // I/O to read A and B, N elements each
10        ...
11        vecAdd(A, B, C, N);
12    }
```

**FIGURE 2.4**

A simple traditional vector addition C code example.

```
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.

cudaMalloc()
- Allocates object in the device global memory
- Two parameters
  o **Address of a pointer** to the allocated object
  o **Size** of allocated object in terms of bytes

cudaFree()
- Frees object from device global memory
  o **Pointer** to freed object

**FIGURE 2.6**

CUDA API functions for managing device global memory.

cudaMemcpy()
- memory data transfer
- Requires four parameters
  - Pointer to destination
  - Pointer to source
  - Number of bytes copied
  - Type/Direction of transfer

**FIGURE 2.7**

CUDA API function for data transfer between host and 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**

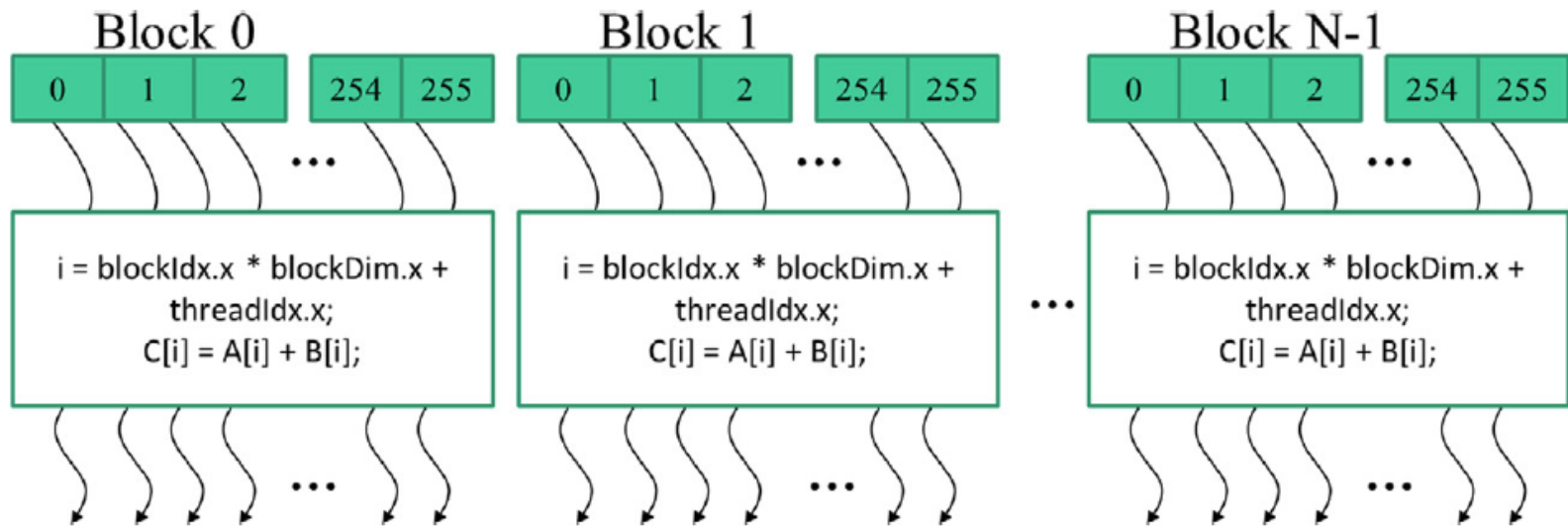A more complete version of vecAdd().

**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.

| Qualifier Keyword | Callable From | Executed On | Executed By |
|---|---|---|---|
| __host__ (default) | Host | Host | Caller host thread |
| __global__ | Host (or Device) | Device | New grid of device threads |
| __device__ | Device | Device | Caller device thread |

**FIGURE 2.11**

CUDA C keywords for function declaration.

```
01    int vectAdd(float* A, float* B, float* C, int n) {
02        // A_d, B_d, C_d allocations and copies omitted
03        ...
04        // Launch ceil(n/256) blocks of 256 threads each
05        vecAddKernel<<<ceil(n/256.0), 256>>>(A_d, B_d, C_d, n);
06    }
```

**FIGURE 2.12**

A vector addition kernel call statement.

```
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    }
```

**FIGURE 2.13**

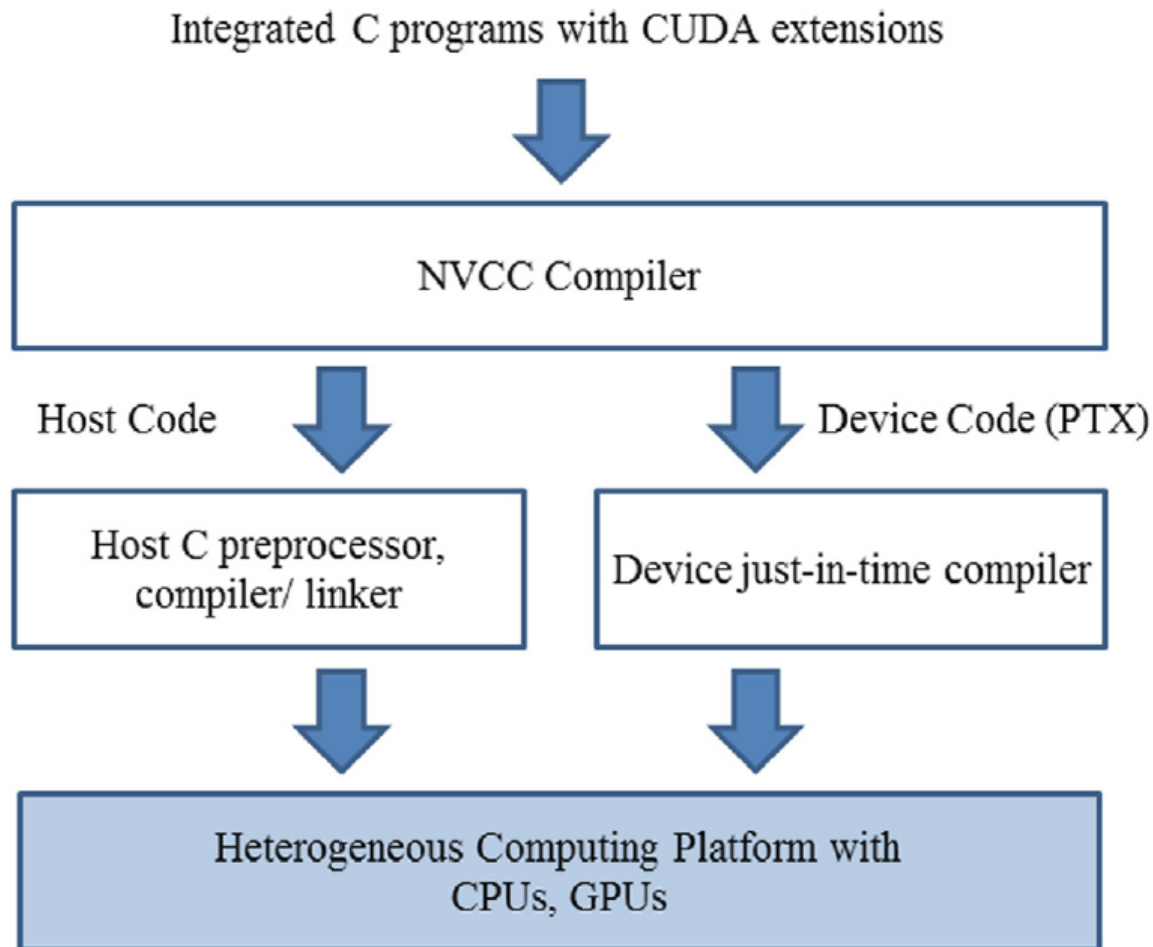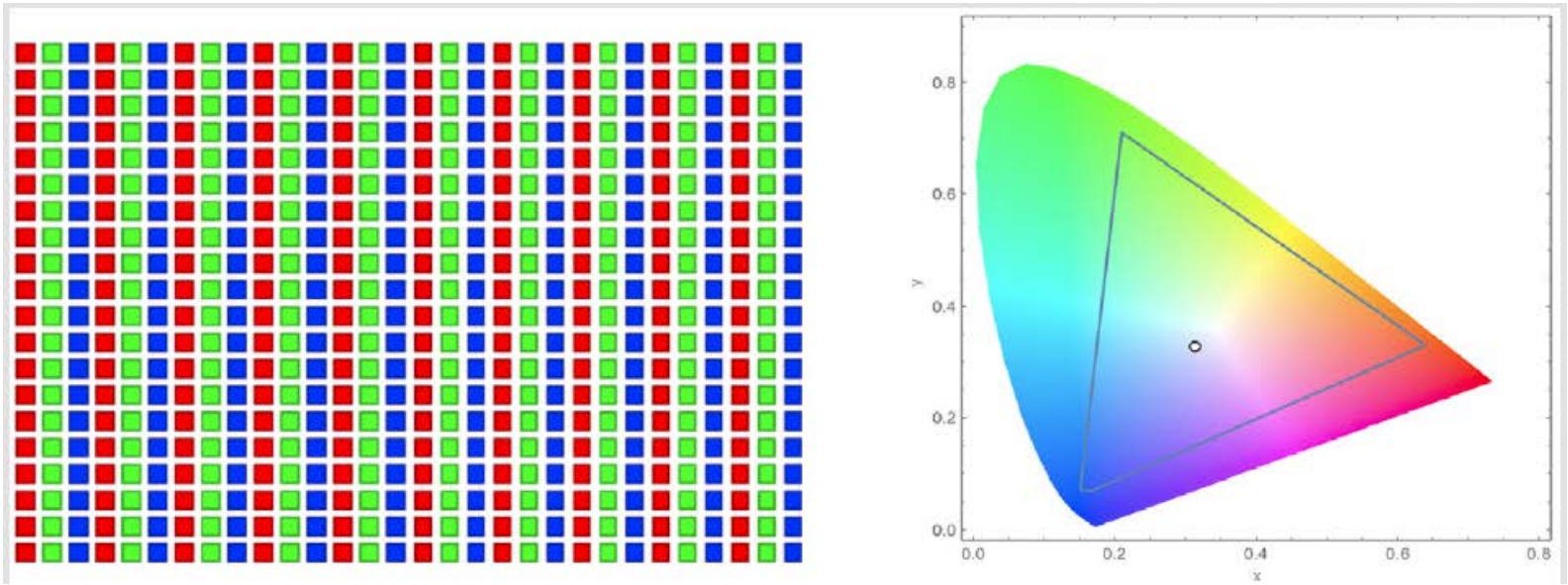A complete version of the host code in the vecAdd function.

**FIGURE 2.14**

Overview of the compilation process of a CUDA C program.

In-text figure 1