

CHAPTER 3

Multidimensional grids and data

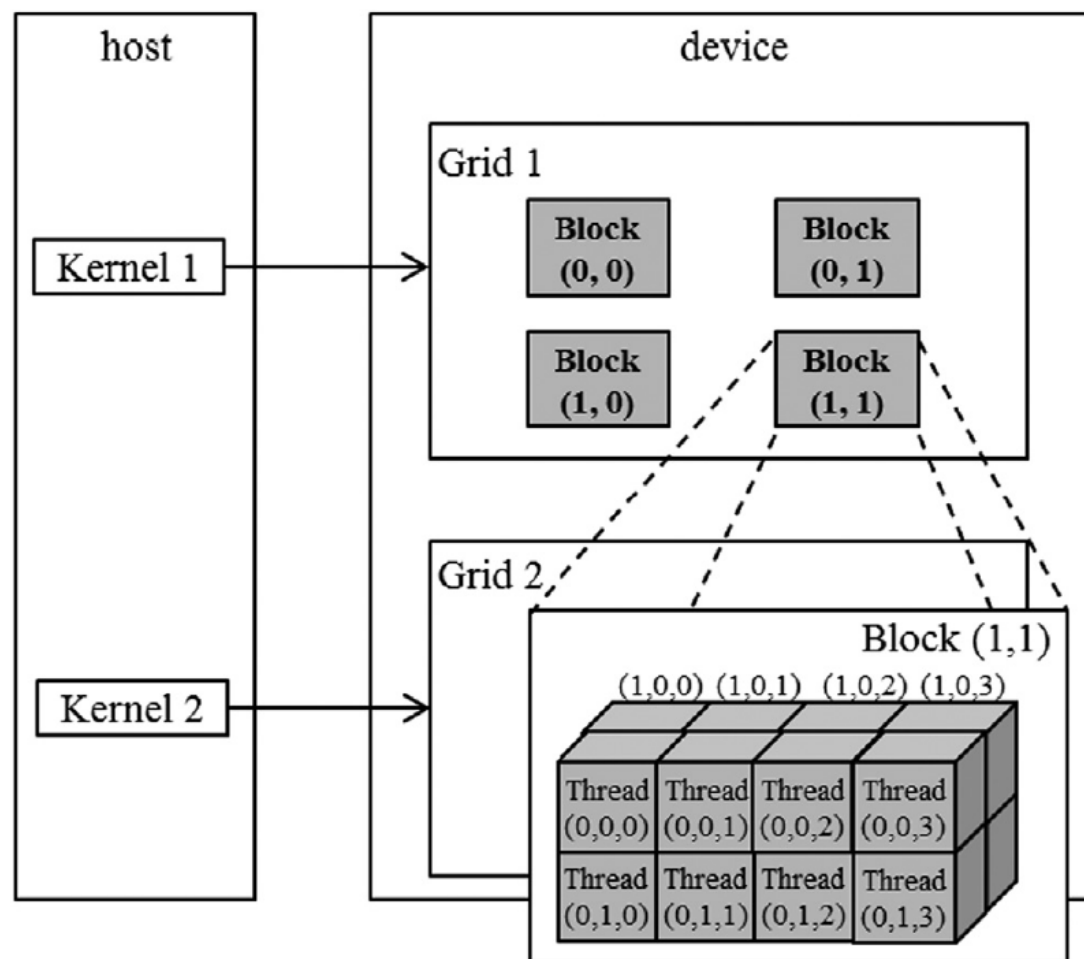
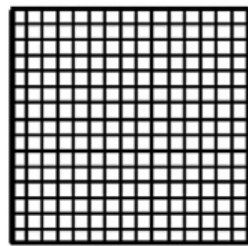


FIGURE 3.1

A multidimensional example of CUDA grid organization.



16×16 blocks

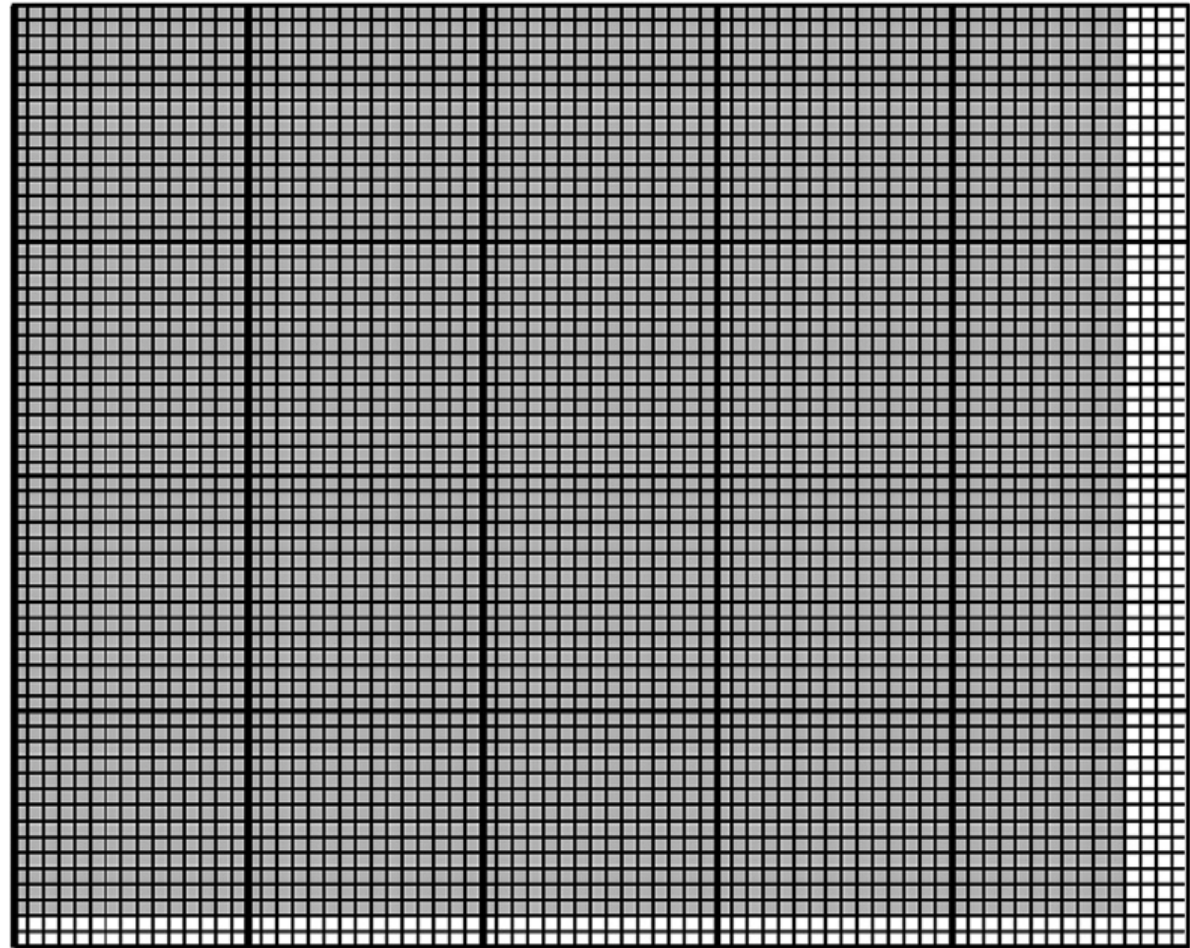


FIGURE 3.2

Using a 2D thread grid to process a 62×76 picture P.

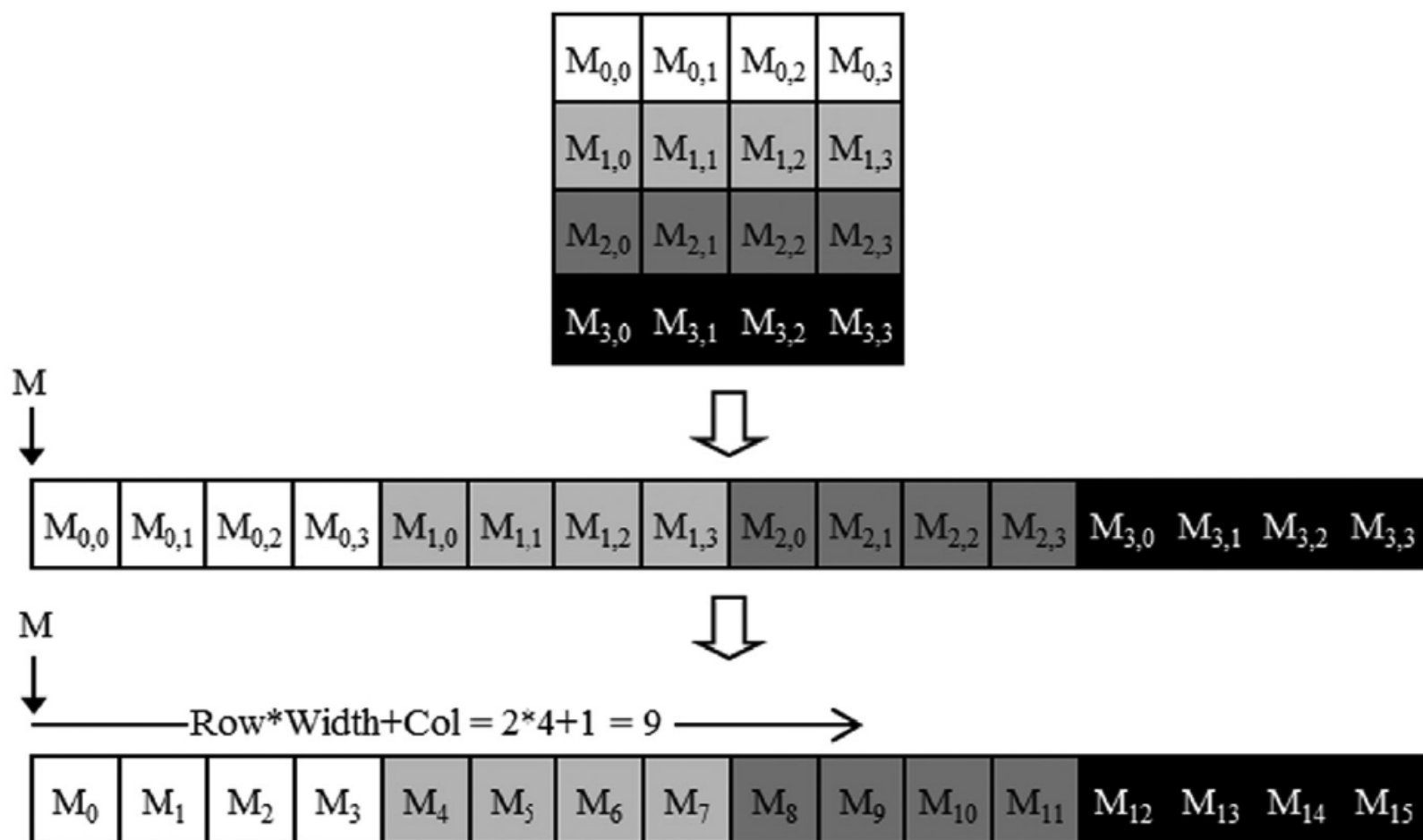


FIGURE 3.3

Row-major layout for a 2D C array. The result is an equivalent 1D array accessed by an index expression $j * \text{Width} + i$ for an element that is in the j th row and i th column of an array of Width elements in each row.

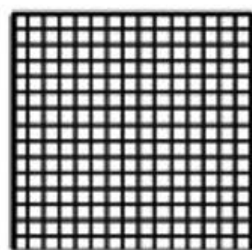
```

01 // The input image is encoded as unsigned chars [0, 255]
02 // Each pixel is 3 consecutive chars for the 3 channels (RGB)
03 __global__
04 void colorToGrayscaleConversion(unsigned char * Pout,
05                                unsigned char * Pin, int width, int height) {
06     int col = blockIdx.x*blockDim.x + threadIdx.x;
07     int row = blockIdx.y*blockDim.y + threadIdx.y;
08     if (col < width && row < height) {
09         // Get 1D offset for the grayscale image
10         int grayOffset = row*width + col;
11         // One can think of the RGB image having CHANNEL
12         // times more columns than the gray scale image
13         int rgbOffset = grayOffset*CHANNELS;
14         unsigned char r = Pin[rgbOffset]; // Red value
15         unsigned char g = Pin[rgbOffset + 1]; // Green value
16         unsigned char b = Pin[rgbOffset + 2]; // Blue value
17         // Perform the rescaling and store it
18         // We multiply by floating point constants
19         Pout[grayOffset] = 0.21f*r + 0.71f*g + 0.07f*b;
20     }
21 }

```

FIGURE 3.4

Source code of `colorToGrayscaleConversion` with 2D thread mapping to data.



16×16 blocks

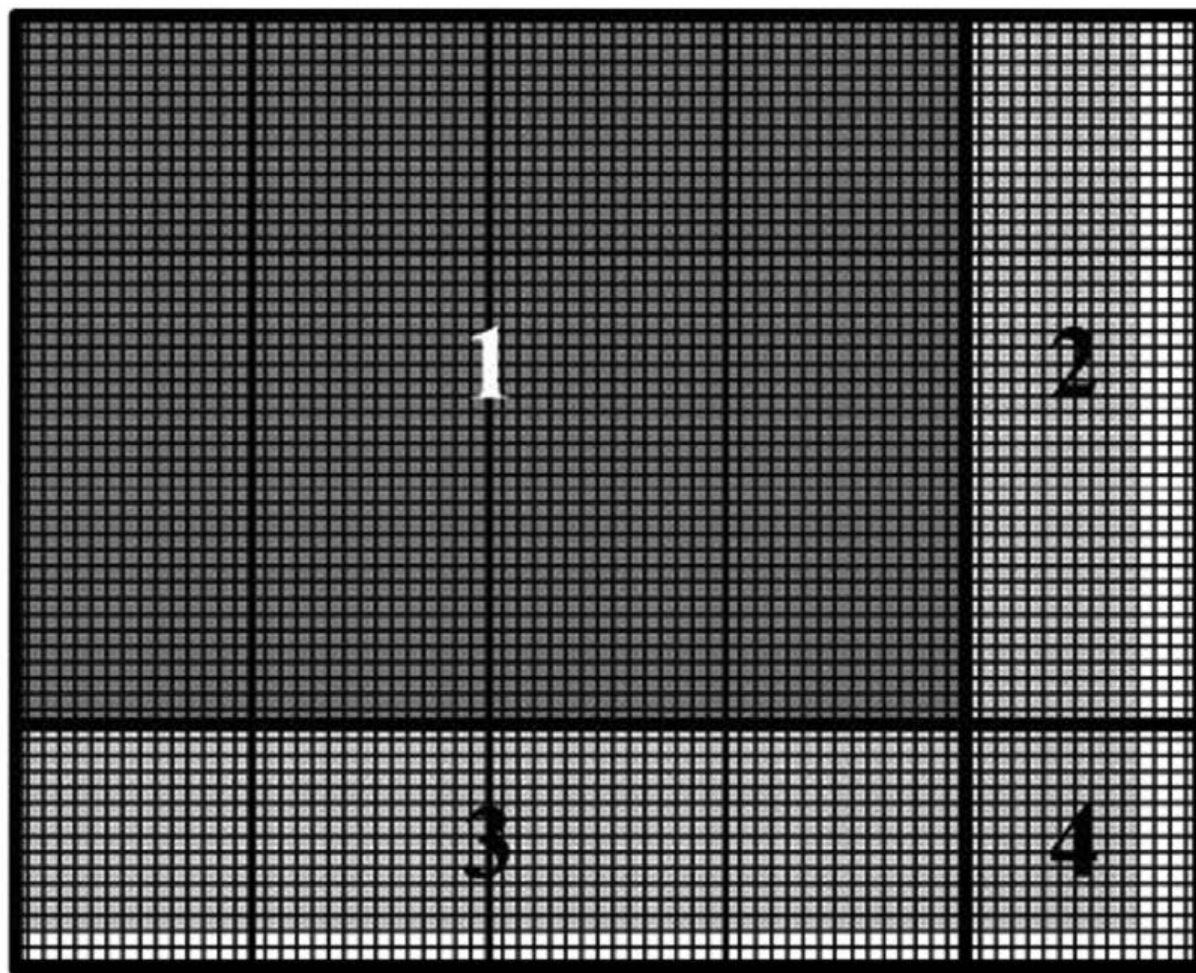


FIGURE 3.5

Covering a 76×62 picture with 16×16 blocks.



FIGURE 3.6

An original image (*left*) and a blurred version (*right*).

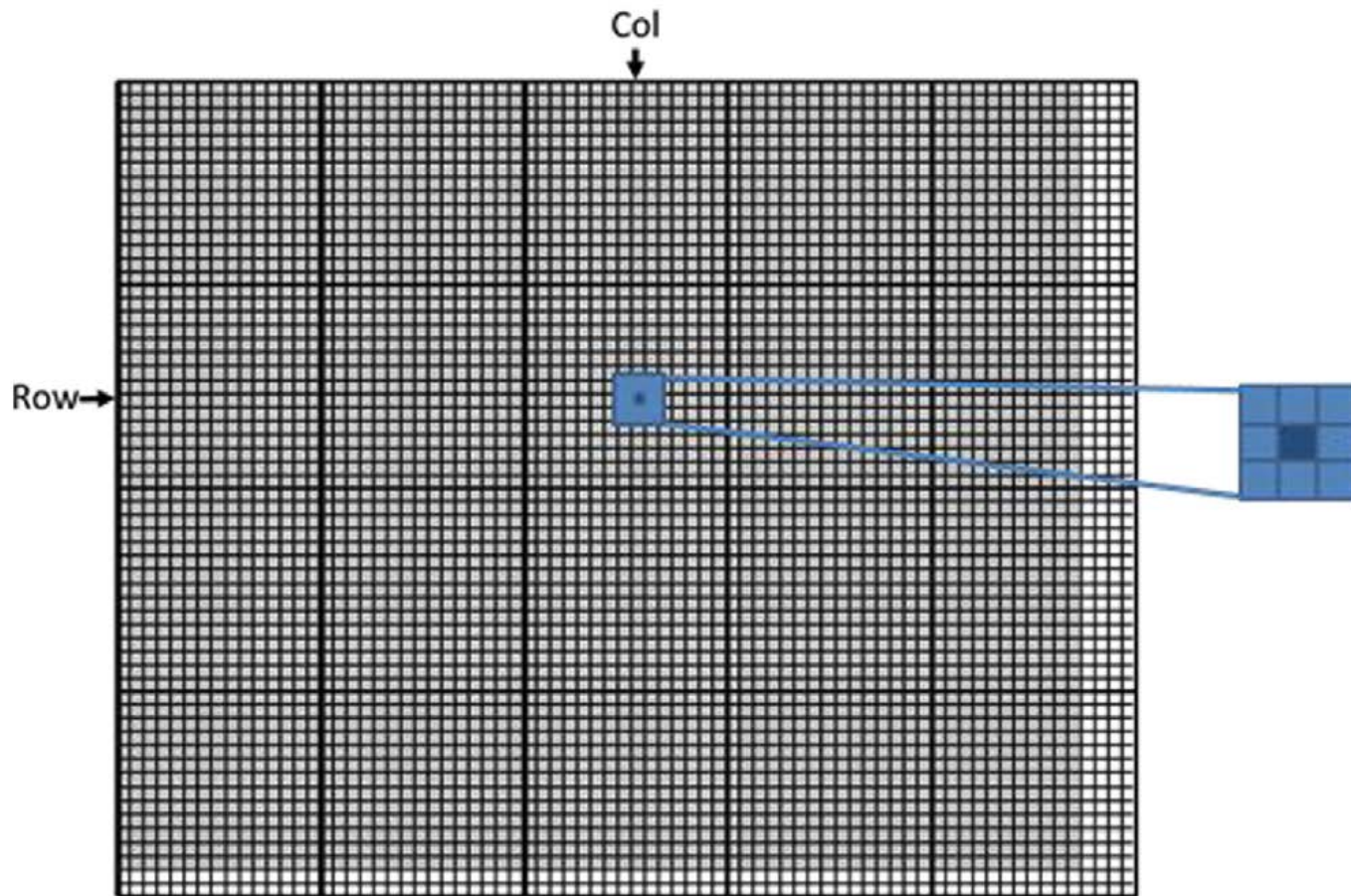


FIGURE 3.7

Each output pixel is the average of a patch of surrounding pixels and itself in the input image.


```

01  __global__
02  void blurKernel(unsigned char *in, unsigned char *out, int w, int h){
03      int col = blockIdx.x*blockDim.x + threadIdx.x;
04      int row = blockIdx.y*blockDim.y + threadIdx.y;
05      if(col < w && row < h) {
06          int pixVal = 0;
07          int pixels = 0;
08          // Get average of the surrounding BLUR_SIZE x BLUR_SIZE box
09          for(int blurRow=-BLUR_SIZE; blurRow<BLUR_SIZE+1; ++blurRow){
10              for(int blurCol=-BLUR_SIZE; blurCol<BLUR_SIZE+1; ++blurCol){
11                  int curRow = row + blurRow;
12                  int curCol = col + blurCol;
13                  // Verify we have a valid image pixel
14                  if(curRow>=0 && curRow<h && curCol>=0 && curCol<w) {
15                      pixVal += in[curRow*w + curCol];
16                      ++pixels; // Keep track of number of pixels in the avg
17                  }
18              }
19          }
20      }
21      // Write our new pixel value out
22      out[row*w + col] = (unsigned char)(pixVal/pixels);
23  }
24  }

```

FIGURE 3.8

An image blur kernel.

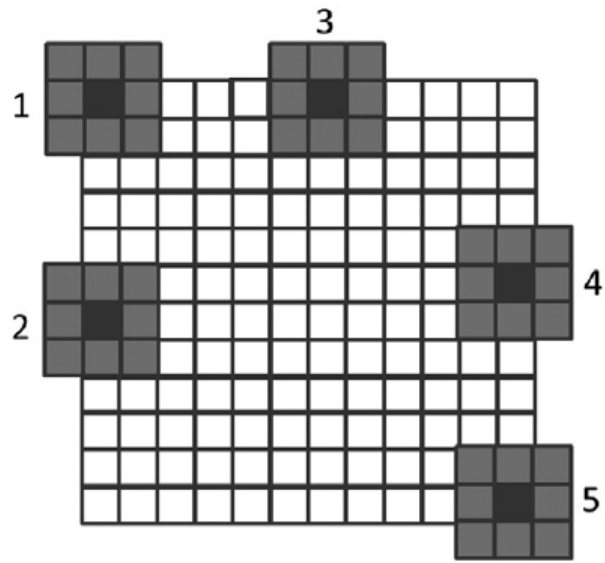


FIGURE 3.9

Handling boundary conditions for pixels near the edges of the image.

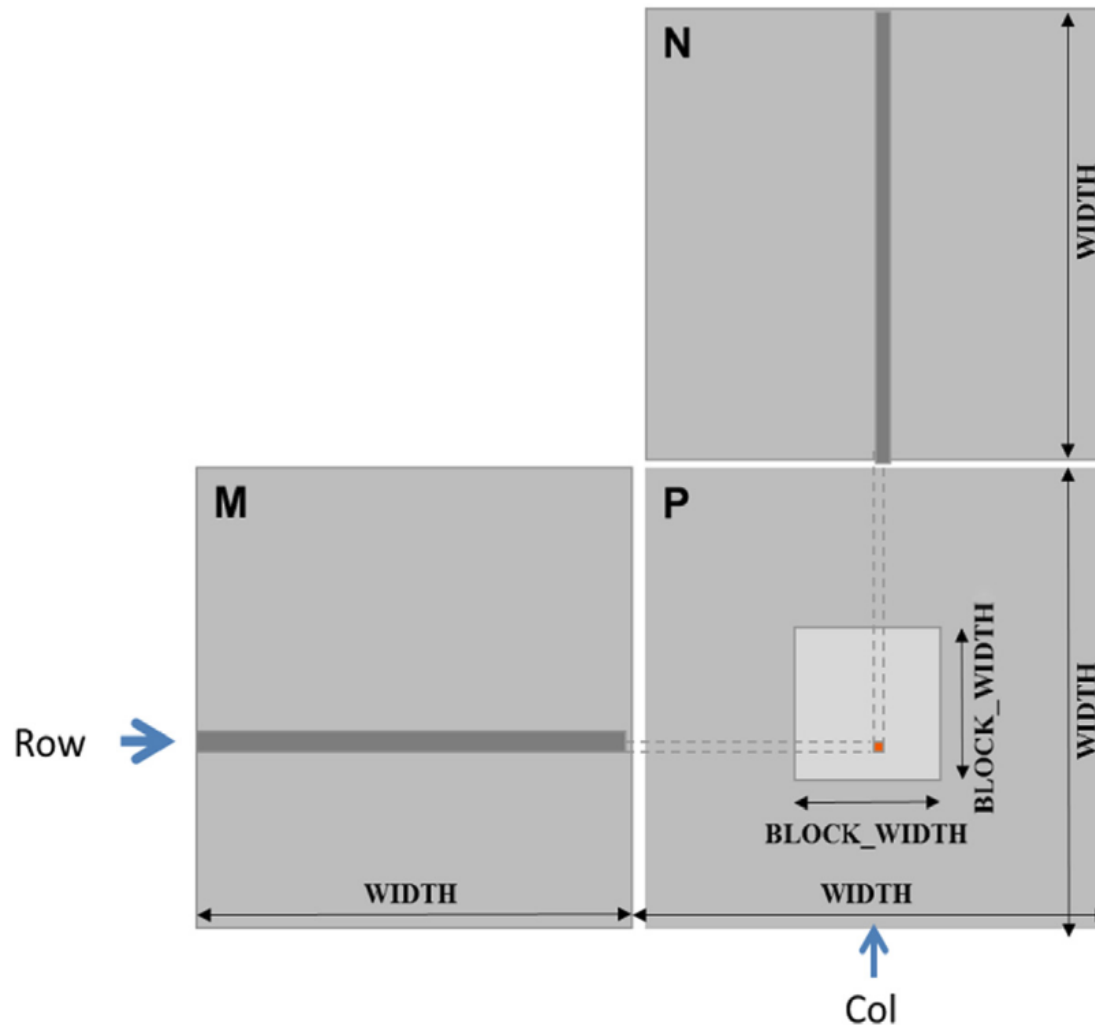


FIGURE 3.10

Matrix multiplication using multiple blocks by tiling P.

```
01  __global__ void MatrixMulKernel(float* M, float* N,  
02                                float* P, int Width) {  
03      int row = blockIdx.y*blockDim.y+threadIdx.y;  
04      int col = blockIdx.x*blockDim.x+threadIdx.x;  
05      if ((row < Width) && (col < Width)) {  
06          float Pvalue = 0;  
07          for (int k = 0; k < Width; ++k) {  
08              Pvalue += M[row*Width+k]*N[k*Width+col];  
09          }  
10          P[row*Width+col] = Pvalue;  
11      }  
12  }
```

FIGURE 3.11

A matrix multiplication kernel using one thread to compute one P element.

BLOCK_WIDTH = 2

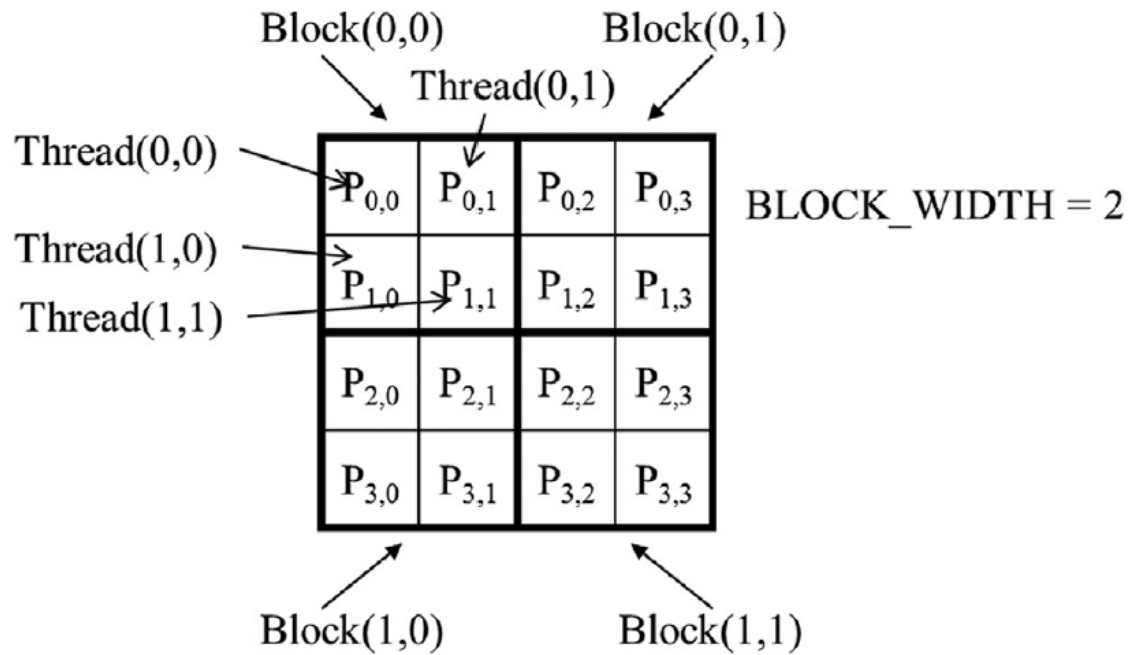


FIGURE 3.12

A small execution example of `matrixMulKernel`.

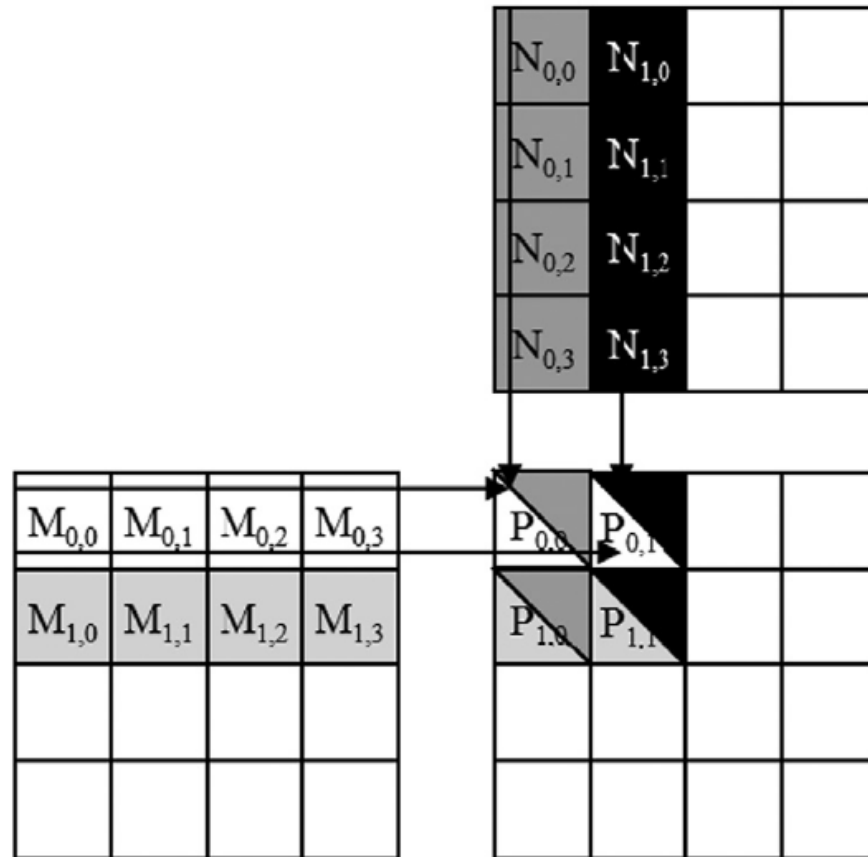


FIGURE 3.13

Matrix multiplication actions of one thread block.

```
dim3 dimGrid(32, 1, 1);  
dim3 dimBlock(128, 1, 1);  
vecAddKernel<<<dimGrid, dimBlock>>>(...);
```

In-text figure 1

```
dim3 dog(32, 1, 1);  
dim3 cat(128, 1, 1);  
vecAddKernel<<<dog, cat>>>(...);
```

In-text figure 2

```
dim3 dimGrid(ceil(n/256.0), 1, 1);  
dim3 dimBlock(256, 1, 1);  
vecAddKernel<<<dimGrid, dimBlock>>>(...);
```

In-text figure 3

```
vecAddKernel<<<ceil(n/256.0), 256>>> (...);
```

In-text figure 4


```
dim3 dimGrid(2, 2, 1);  
dim3 dimBlock(4, 2, 2);  
KernelFunction<<<dimGrid, dimBlock>>>(...);
```

In-text figure 5

```
dim3 dimGrid(ceil(m/16.0), ceil(n/16.0), 1);  
dim3 dimBlock(16, 16, 1);  
colorToGrayscaleConversion<<<dimGrid,dimBlock>>>  
                                (Pin_d, Pout_d, m, n);
```

In-text figure 6

```
col = blockIdx.x*blockDim.x + threadIdx.x
```

In-text figure 7

```
int plane = blockIdx.z*blockDim.z + threadIdx.z
```

In-text figure 8

```
row = blockIdx.y*blockDim.y + threadIdx.y
```

In-text figure 9


```
col = blockIdx.x*blockDim.x + threadIdx.x
```

In-text figure 10

```

01  __global__ void foo_kernel(float* a, float* b, unsigned int M,
unsigned int N) {
02      unsigned int row = blockIdx.y*blockDim.y + threadIdx.y;
03      unsigned int col = blockIdx.x*blockDim.x + threadIdx.x;
04      if(row < M && col < N) {
05          b[row*N + col] = a[row*N + col]/2.1f + 4.8f;
06      }
07  }
08  void foo(float* a_d, float* b_d) {
09      unsigned int M = 150;
10      unsigned int N = 300;
11      dim3 bd(16, 32);
12      dim3 gd((N - 1)/16 + 1, (M - 1)/32 + 1);
13      foo_kernel <<< gd, bd >>>(a_d, b_d, M, N);
14  }

```

In-text figure 11