# CHAPTER 7

Convolution

An introduction to constant memory
and caching
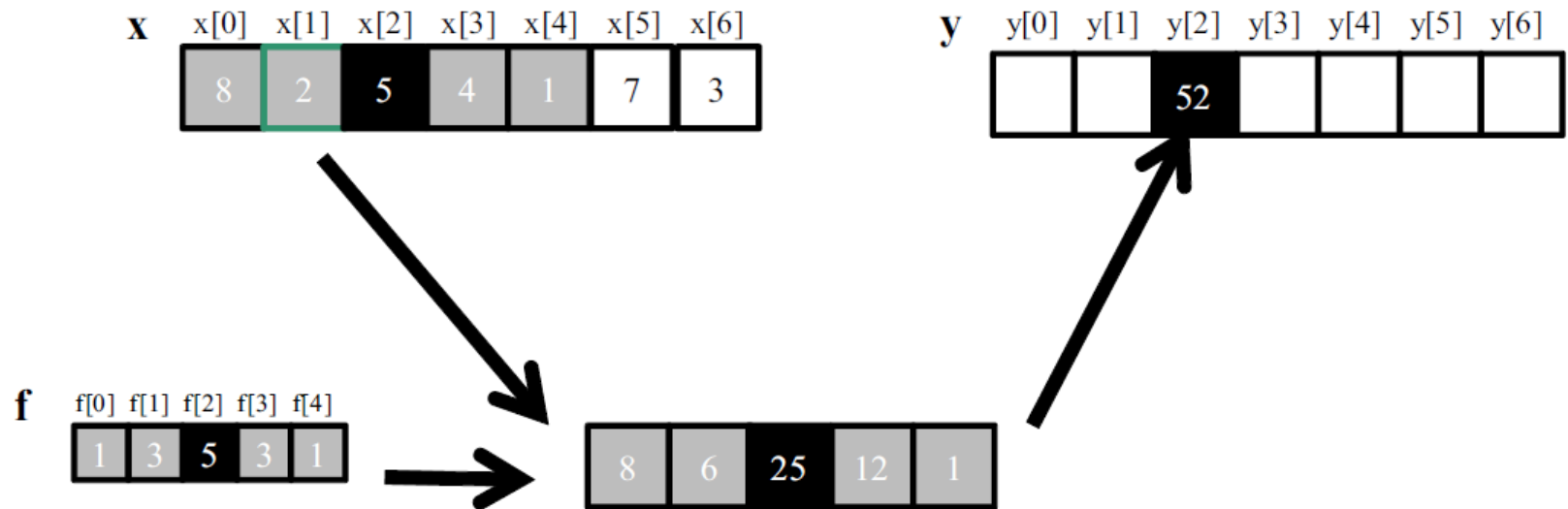
**FIGURE 7.1**

A 1D convolution example, inside elements.
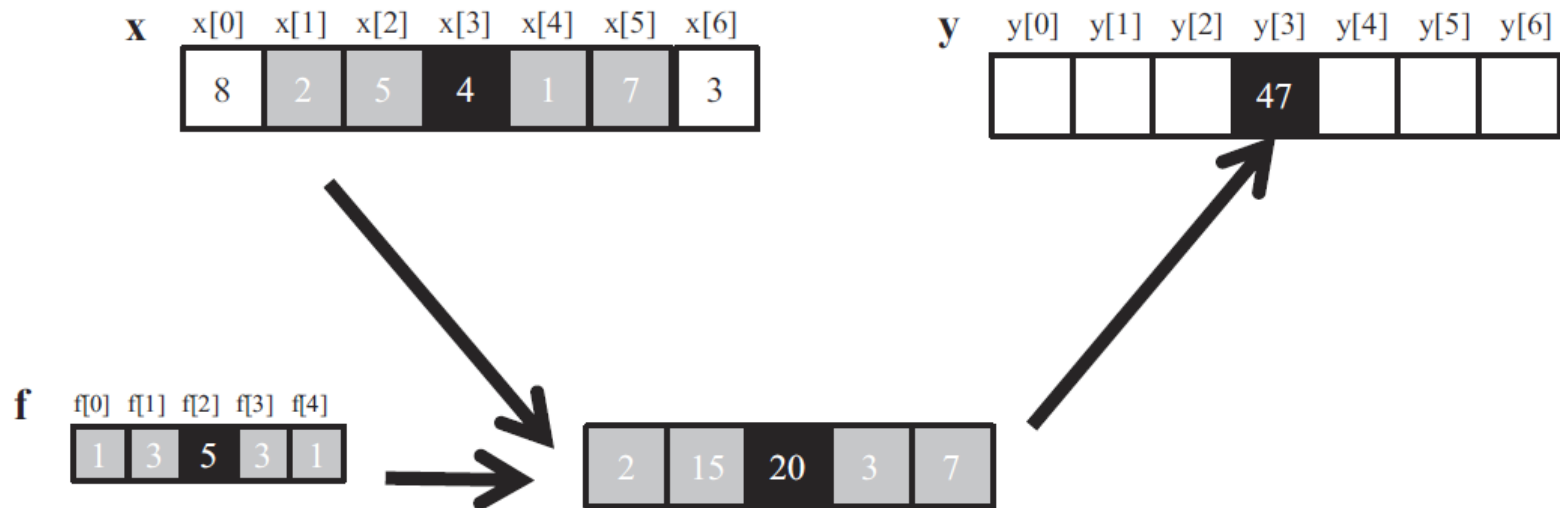
**FIGURE 7.2**

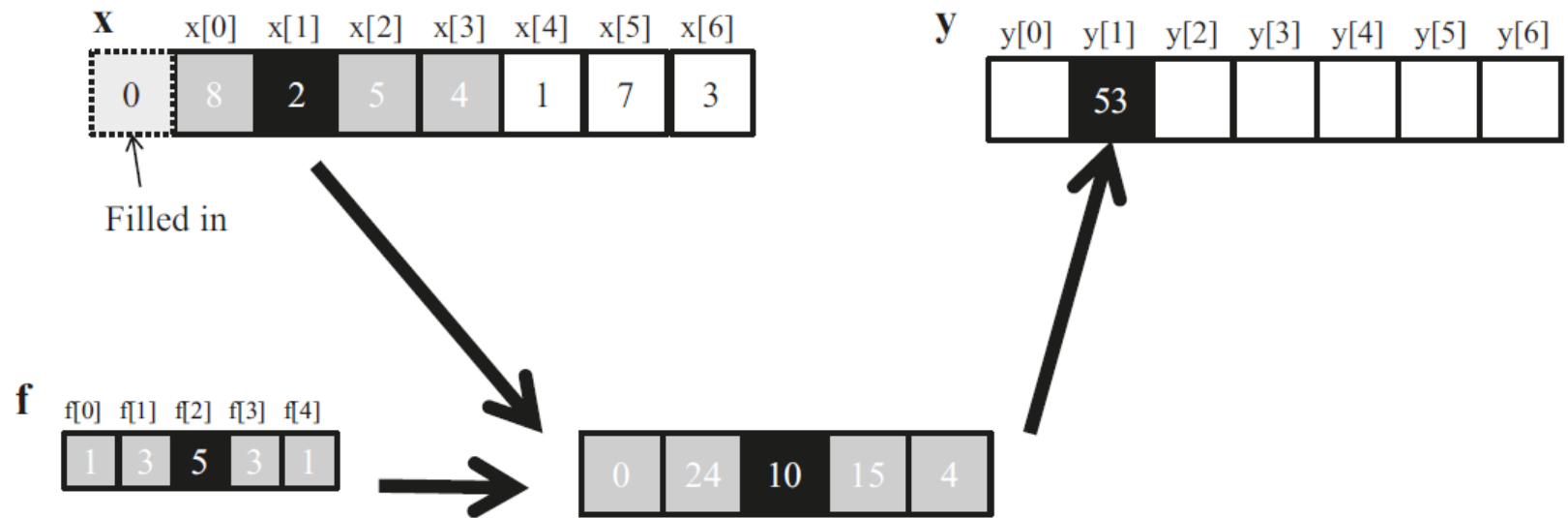1D convolution, calculation of y[3].

**FIGURE 7.3**

A 1D convolution boundary condition.

**FIGURE 7.4**
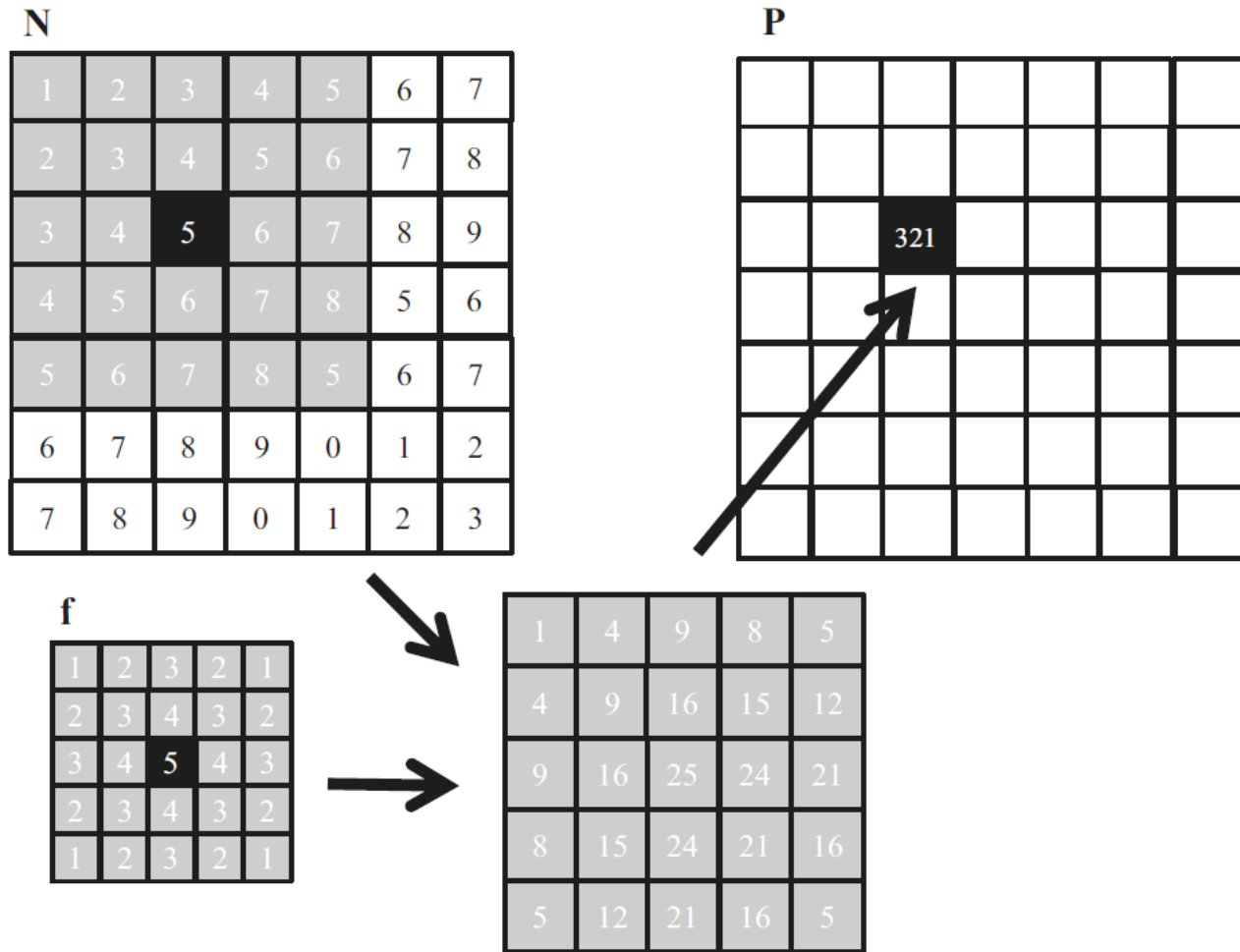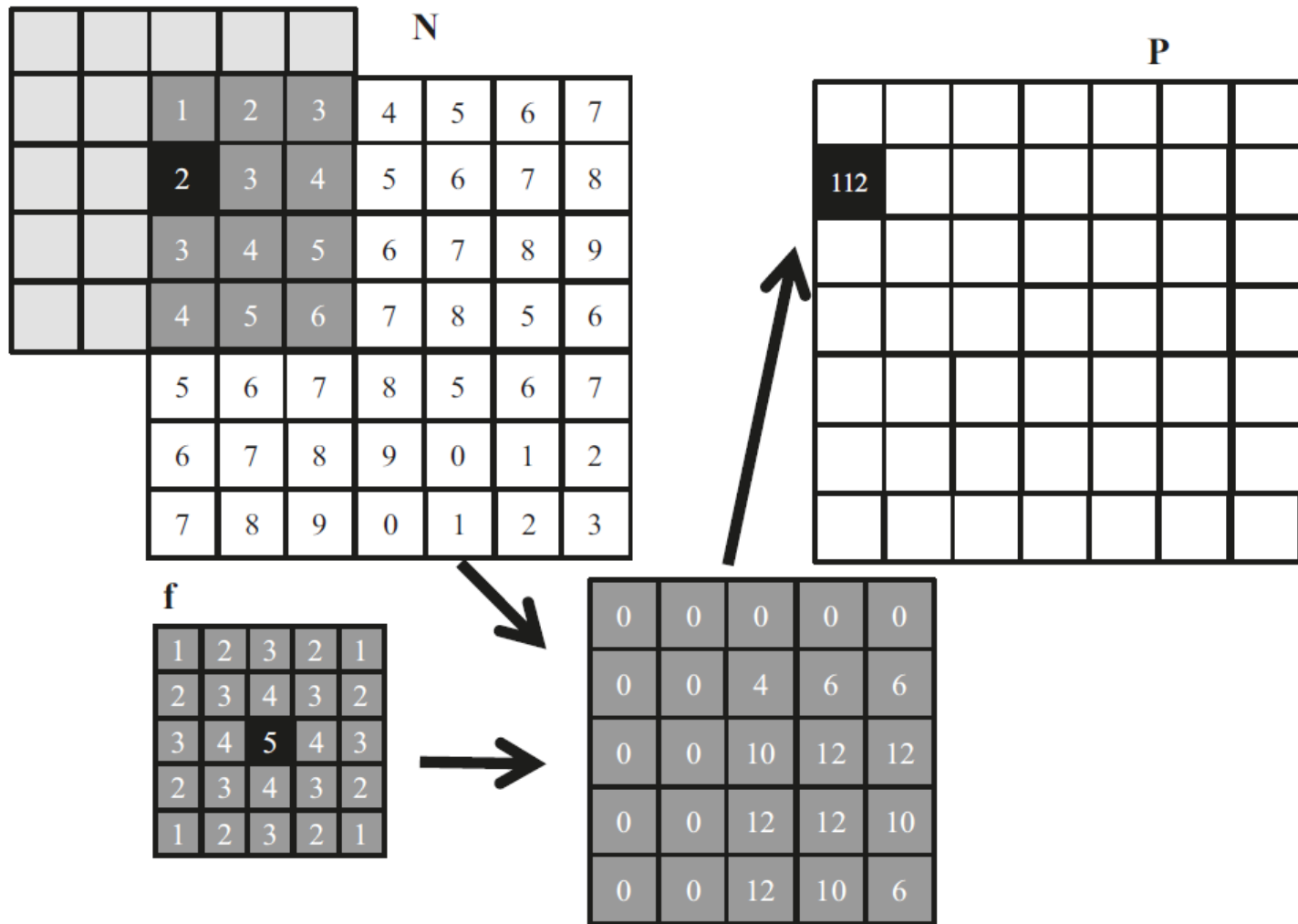
A 2D convolution example.

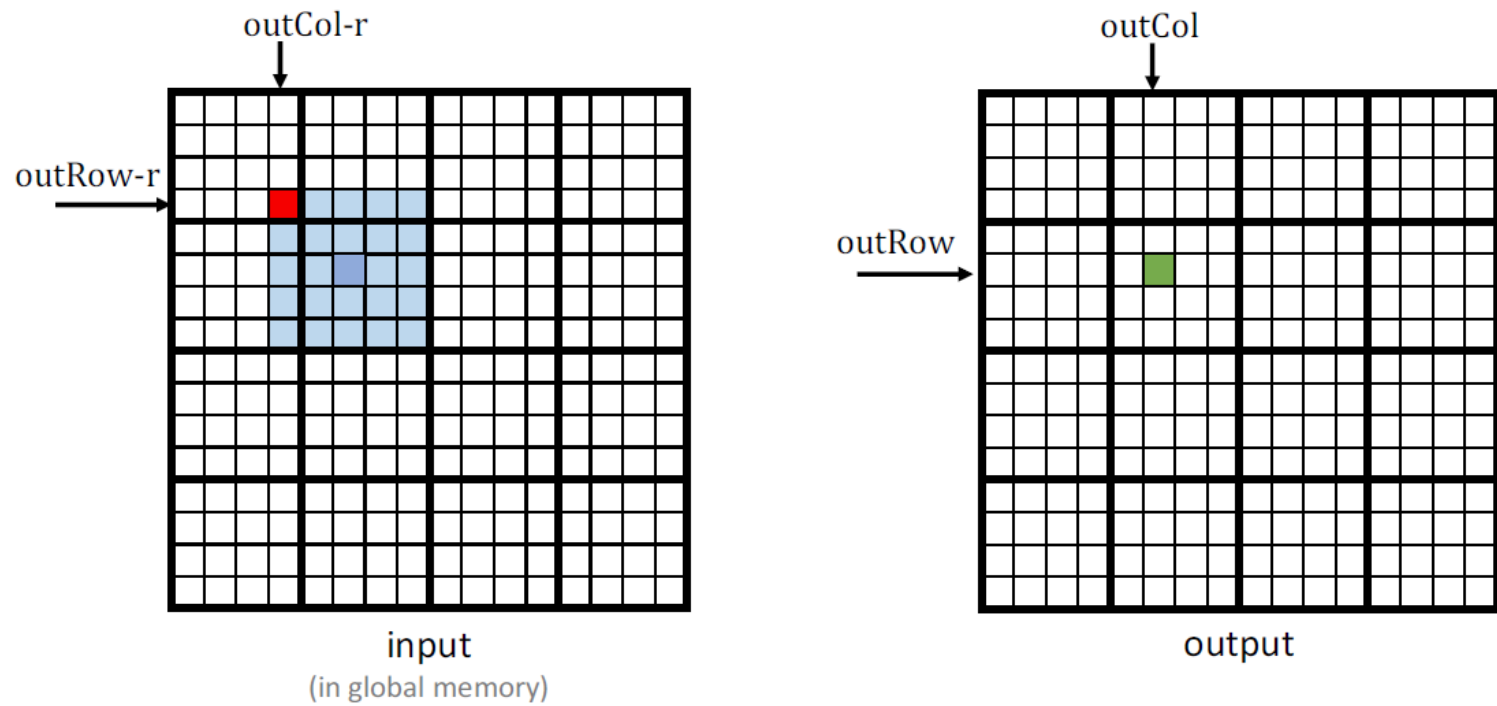**FIGURE 7.5**

A 2D convolution boundary condition.

**FIGURE 7.6**

Parallelization and thread organization for 2D convolution.

```
01  __global__ void convolution_2D_basic_kernel(float *N, float *F, float *P,
      int r, int width, int height) {
02    int outCol = blockIdx.x*blockDim.x + threadIdx.x;
03    int outRow = blockIdx.y*blockDim.y + threadIdx.y;
04    float Pvalue = 0.0f;
05    for (int fRow = 0; fRow < 2*r+1; fRow++) {
06      for (int fCol = 0; fCol < 2*r+1; fCol++) {
07          inRow = outRow - r + fRow;
08          inCol = outCol - r + fCol;
09          if (inRow >= 0 && inRow < height && inCol >= 0 && inCol < width) {
10              Pvalue += F[fRow][fCol]*N[inRow*width + inCol];
11          }
12        }
13      }
14    P[outRow][outCol] = Pvalue;
15 }
```

**FIGURE 7.7**

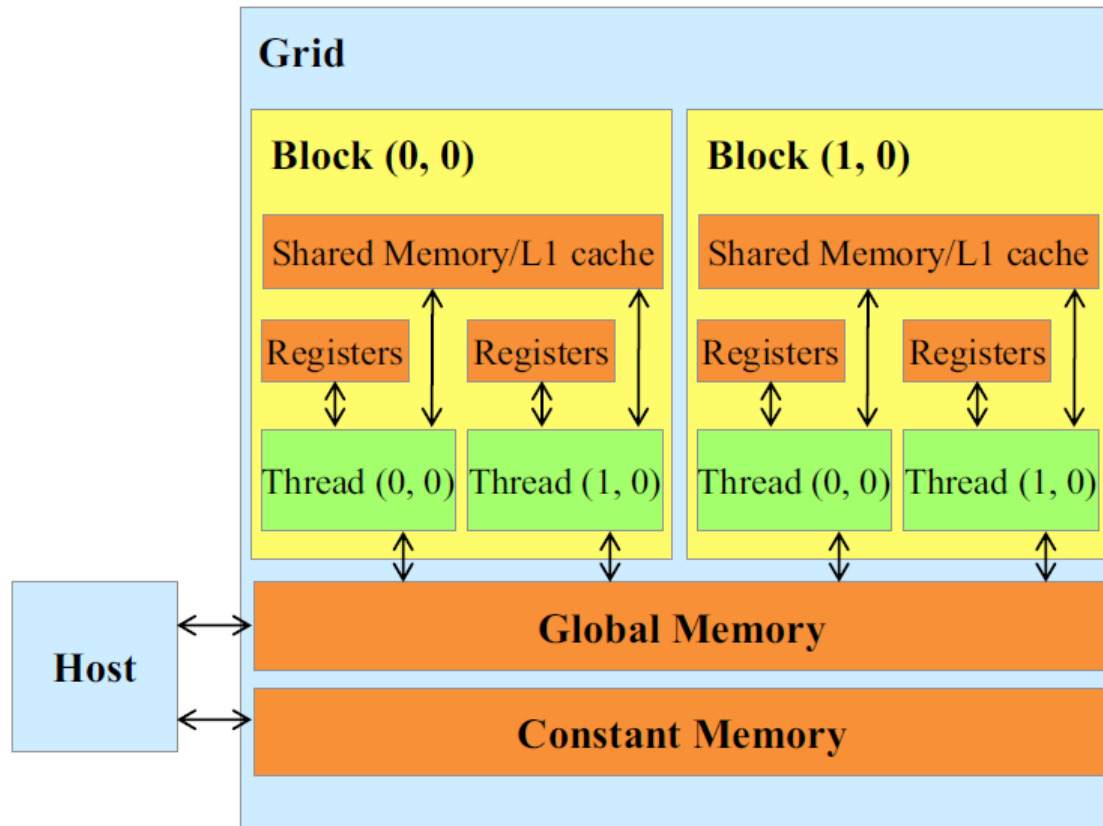A 2D convolution kernel with boundary condition handling.

**FIGURE 7.8**

A review of the CUDA memory model.

```
01  __global__ void convolution_2D_const_mem_kernel(float *N, float *P, int r,
      int width, int height) {
02      int outCol = blockIdx.x*blockDim.x + threadIdx.x;
03      int outRow = blockIdx.y*blockDim.y + threadIdx.y;
04      float Pvalue = 0.0f;
05      for (int fRow = 0; fRow < 2*r+1; fRow++) {
06        for (int fCol = 0; fCol < 2*r+1; fCol++) {
07            inRow = outRow - r + fRow;
08            inCol = outCol - r + fCol;
09            if (inRow >= 0 && inRow < height && inCol >= 0 && inCol < width) {
10                Pvalue += F[fRow][fCol]*N[inRow*width + inCol];
11            }
12          }
13      }
14      P[outRow*width+outCol] = Pvalue;
15  }
```

**FIGURE 7.9**

A 2D convolution kernel using constant memory for F.

**FIGURE 7.10**

A simplified view of the cache hierarchy of modern processors.

input tile dimension

input tile
(in shared memory)

input

filter dimension

filter radius

filter
(in constant
memory)

output tile dimension
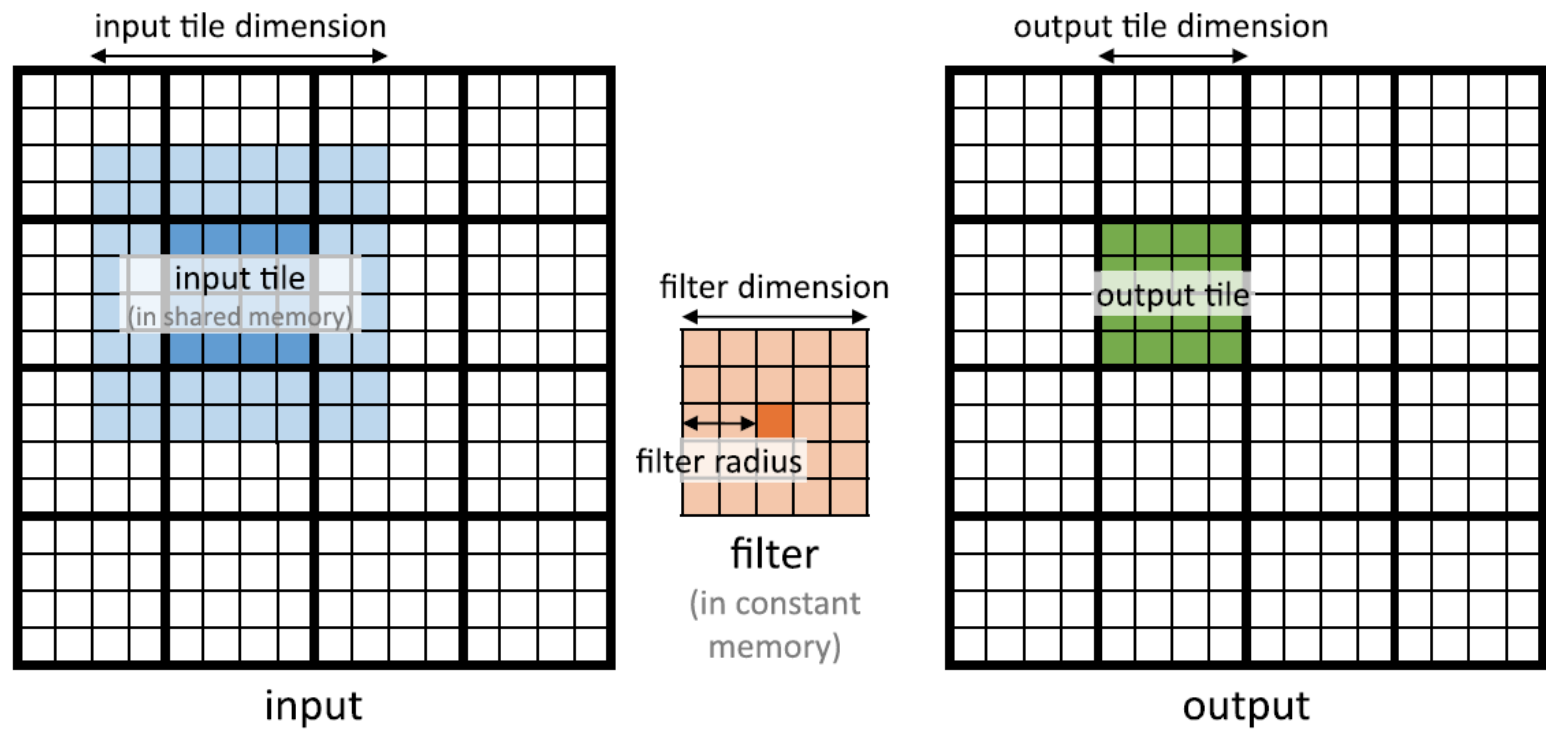
output tile

output

**FIGURE 7.11**

Input tile versus output tile in a 2D convolution.

```
01  #define IN_TILE_DIM 32
02  #define OUT_TILE_DIM ((IN_TILE_DIM) - 2*(FILTER_RADIUS))
03  __constant__ float F_c[2*FILTER_RADIUS+1][2*FILTER_RADIUS+1];
04  __global__ void convolution_tiled_2D_const_mem_kernel(float *N, float *P,
05                                                        int width, int height) {
06    int col = blockIdx.x*OUT_TILE_DIM + threadIdx.x - FILTER_RADIUS;
07    int row = blockIdx.y*OUT_TILE_DIM + threadIdx.y - FILTER_RADIUS;
08    //loading input tile
09     __shared__ N_s[IN_TILE_DIM][IN_TILE_DIM];
10    if(row>=0 && row<height && col>=0 && col<width) {
11      N_s[threadIdx.y][threadIdx.x] = N[row*width + col];
12    } else {
13      N_s[threadIdx.y][threadIdx.x] = 0.0;
14    }
15    __syncthreads();
16    // Calculating output elements
17    int tileCol = threadIdx.x - FILTER_RADIUS;
18    int tileRow = threadIdx.y - FILTER_RADIUS;
19    // turning off the threads at the edges of the block
20    if (col >= 0 && col < width && row >=0 && row < height) {
21      if (tileCol>=0 && tileCol<OUT_TILE_DIM && tileRow>=0
22                    && tileRow<OUT_TILE_DIM){
23        float Pvalue = 0.0f;
24        for (int fRow = 0; fRow < 2*FILTER_RADIUS+1; fRow++) {
25          for (int fCol = 0; fCol < 2*FILTER_RADIUS+1; fCol++) {
26            Pvalue += F[fRow][fCol]*N_s[tileRow+fRow][tileCol+fCol];
27          }
28        }
29        P[row*width+col] = Pvalue;
30      }
31    }
32  }
```

**FIGURE 7.12**

A tiled 2D convolution kernel using constant memory for F.

Input Tile and Thread Block

Output Tile

(0,0)

(1,1)

(6,6)

(0,0)

(5,5)

**FIGURE 7.13**

A small example that illustrates the thread organization for using the input tile elements in the shared memory to calculate the output tile elements.

| IN_TILE_DIM | | 8 | 16 | 32 | Bound |
|---|---|---|---|---|---|
| 5x5 filter (FILTER_RADIUS = 2) | OUT_TILE_DIM | 4 | 12 | 28 | - |
| | Ratio | 3.13 | 7.03 | 9.57 | 12.5 |
| 9x9 filter (FILTER_RADIUS = 4) | OUT_TILE_DIM | - | 8 | 24 | - |
| | Ratio | - | 10.13 | 22.78 | 40.5 |

**FIGURE 7.14**

Arithmetic-to-global memory access ratio as a function of tile size and filter size for a 2D tiled convolution.

```
01   #define TILE_DIM 32
02   __constant__ float F_c[2*FILTER_RADIUS+1][2*FILTER_RADIUS+1];
03   __global__ void convolution_cached_tiled_2D_const_mem_kernel(float *N,
                                              float *P, int width, int height) {
04     int col = blockIdx.x*TILE_DIM + threadIdx.x;
05     int row = blockIdx.y*TILE_DIM + threadIdx.y;
       //loading input tile
06     __shared__ N_s[TILE_DIM][TILE_DIM];
07     if(row<height && col<width) {
08       N_s[threadIdx.y][threadIdx.x] = N[row*width + col];
09     } else {
10       N_s[threadIdx.y][threadIdx.x] = 0.0;
11     }
12     __syncthreads();
       // Calculating output elements
       // turning off the threads at the edges of the block
13     if (col < width && row < height) {
14       float Pvalue = 0.0f;
15       for (int fRow = 0; fRow < 2*FILTER RADIUS+1; fRow++) {
16         for (int fCol = 0; fCol < 2*FILTER_RADIUS+1; fCol++) {
17           if (threadIdx.x-FILTER_RADIUS+fCol >= 0 &&
18               threadIdx.x-FILTER_RADIUS+fCol < TILE_DIM &&
19               threadIdx.y-FILTER_RADIUS+fRow >= 0 &&
20               threadIdx.y-FILTER_RADIUS+fRow < TILE_DIM){
21             Pvalue += F[fRow][fCol]*N_s[threadIdx.y+fRow][threadIdx.x+fCol];
22           }
23           else {
24             if (row-FILTER_RADIUS+fRow >= 0 &&
25                 row-FILTER_RADIUS+fRow < height &&
26                 col-FILTER_RADIUS+fCol >=0 &&
27                 col-FILTER RADIUS+fCol < width) {
24               Pvalue += F[fRow][fCol]*
25                                        N[(row-FILTER_RADIUS+fRow)*width+col-
FILTER_RADIUS+fCol];
26             }
27           }
28         }
29         P[row*width+col] = Pvalue;
30       }
31     }
32   }
```

**FIGURE 7.15**

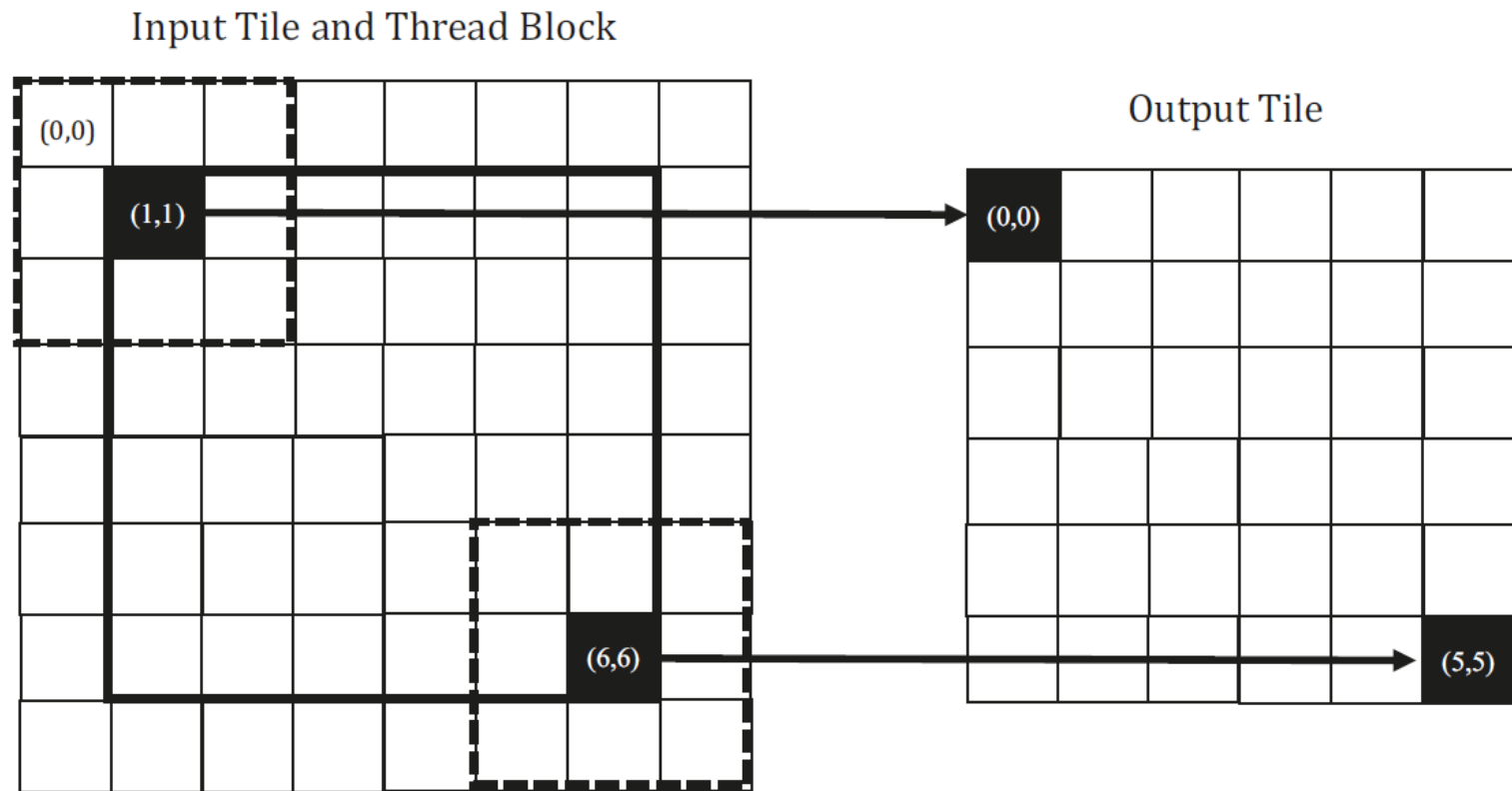A tiled 2D convolution kernel using caching for halos and constant memory for F.

$$
\begin{aligned}
P_{2,2} =\ & N_{0,0}*M_{0,0} + N_{0,1}*M_{0,1} + N_{0,2}*M_{0,2} + N_{0,3}*M_{0,3} + N_{0,4}*M_{0,4} \\
& + N_{1,0}*M_{1,0} + N_{1,1}*M_{1,1} + N_{1,2}*M_{1,2} + N_{1,3}*M_{1,3} + N_{1,4}*M_{1,4} \\
& + N_{2,0}*M_{2,0} + N_{2,1}*M_{1,1} + N_{2,2}*M_{2,2} + N_{2,3}*M_{2,3} + N_{2,4}*M_{2,4} \\
& + N_{3,0}*M_{3,0} + N_{3,1}*M_{3,1} + N_{3,2}*M_{3,2} + N_{3,3}*M_{3,3} + N_{3,4}*M_{3,4} \\
& + N_{4,0}*M_{4,0} + N_{4,1}*M_{4,1} + N_{4,2}*M_{4,2} + N_{4,3}*M_{4,3} + N_{4,4}*M_{4,4} \\
=\ & 1*1 + 2*2 + 3*3 + 4*2 + 5*1 \\
& + 2*2 + 3*3 + 4*4 + 5*3 + 6*2 \\
& + 3*3 + 4*4 + 5*5 + 6*4 + 7*3 \\
& + 4*2 + 5*3 + 6*4 + 7*3 + 8*2 \\
& + 5*1 + 6*2 + 7*3 + 8*2 + 5*1 \\
=\ & 1 + 4 + 9 + 8 + 5 \\
& + 4 + 9 + 16 + 15 + 12 \\
& + 9 + 16 + 25 + 24 + 21 \\
& + 8 + 15 + 24 + 21 + 16 \\
& + 5 + 12 + 21 + 16 + 5 \\
=\ & 321
\end{aligned}
$$

In-text figure 1

```
__global__ void
convolution_2D_basic_kernel(float *N, float *F, float *P, int r,
                            int width, int height) {
    // kernel body
}
```

In-text figure 2

```
#define FILTER_RADIUS 2
__constant__ float F[2*FILTER_RADIUS+1][2*FILTER_RADIUS+1];
```

In-text figure 3

```
cudaMemcpyToSymbol(F,F_h,(2*FILTER_RADIUS+1)*(2*FILTER_RADIUS+1)*sizeof(float));
```

In-text figure 4

```
cudaMemcpyToSymbol(dest, src, size)
```

In-text figure 5