

CHAPTER 5

Memory architecture and data locality

```
07     for (int k = 0; k < Width; ++k) {  
08         Pvalue += M[row*Width+k] * N[k*Width+col];  
09     }
```

FIGURE 5.1

The most executed part of the matrix multiplication kernel in Fig. 3.11.

Device code can:

- R/W per-thread **registers**
- R/W per-thread **local memory**
- R/W per-block **shared memory**
- R/W per-grid **global memory**
- Read only per-grid **constant memory**

Host code can

- Transfer data to/from per grid **global** and **constant** memories

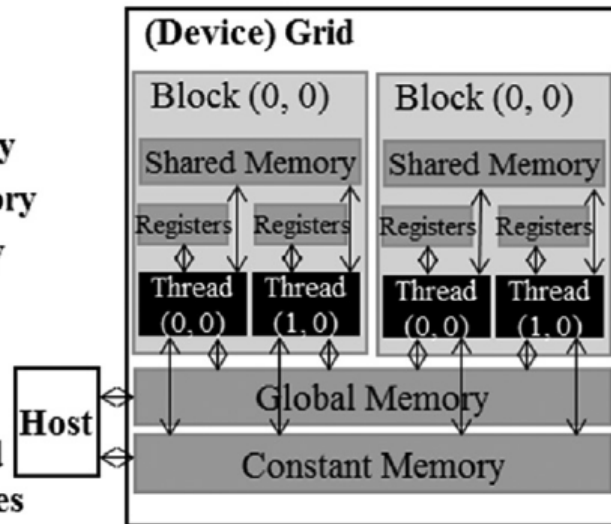


FIGURE 5.2

An (incomplete) overview of the CUDA device memory model. An important type of CUDA memory that is not shown in this figure is the texture memory, since its use is not covered in this textbook.

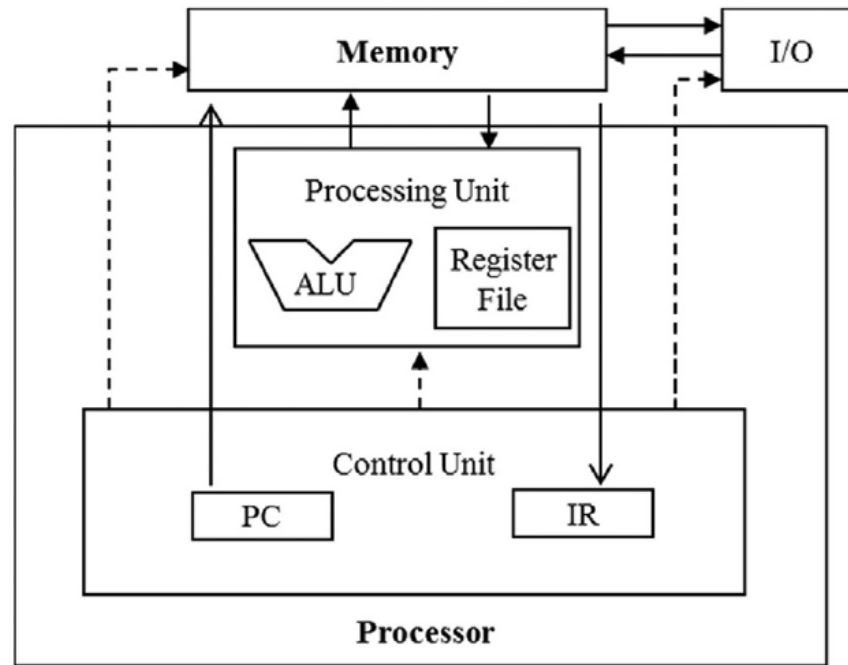


FIGURE 5.3

Memory versus registers in a modern computer based on the von Neumann model.

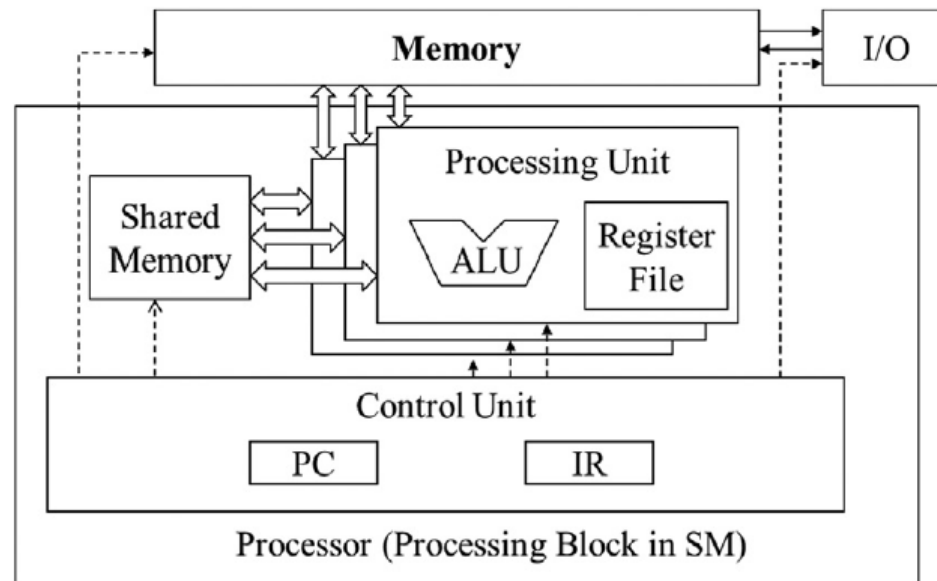


FIGURE 5.4

Shared memory versus registers in a CUDA device SM.

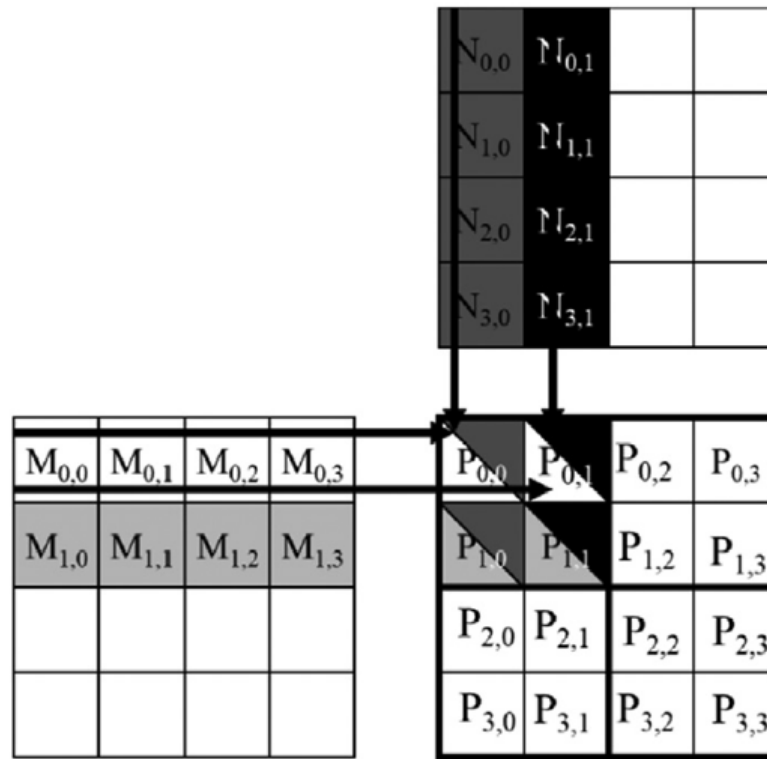



FIGURE 5.5

A small example of matrix multiplication. For brevity we show $M[y*\text{Width}+x]$, $N[y*\text{Width}+x]$, $P[y*\text{Width}+x]$ as $M_{y,x}$, $N_{y,x}$, $P_{y,x}$, respectively.

Access order 

thread _{0,0}	$M_{0,0} * N_{0,0}$	$M_{0,1} * N_{1,0}$	$M_{0,2} * N_{2,0}$	$M_{0,3} * N_{3,0}$
thread _{0,1}	$M_{0,0} * N_{0,1}$	$M_{0,1} * N_{1,1}$	$M_{0,2} * N_{2,1}$	$M_{0,3} * N_{3,1}$
thread _{1,0}	$M_{1,0} * N_{0,0}$	$M_{1,1} * N_{1,0}$	$M_{1,2} * N_{2,0}$	$M_{1,3} * N_{3,0}$
thread _{1,1}	$M_{1,0} * N_{0,1}$	$M_{1,1} * N_{1,1}$	$M_{1,2} * N_{2,1}$	$M_{1,3} * N_{3,1}$

FIGURE 5.6

Global memory accesses performed by threads in block_{0,0}.

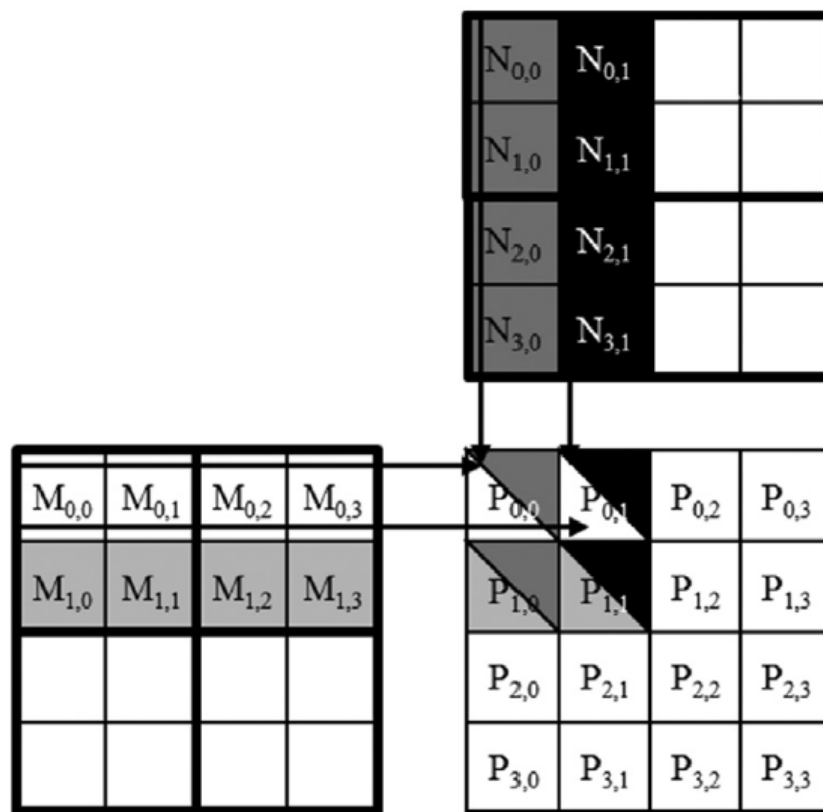


FIGURE 5.7

Tiling M and N to utilize shared memory.

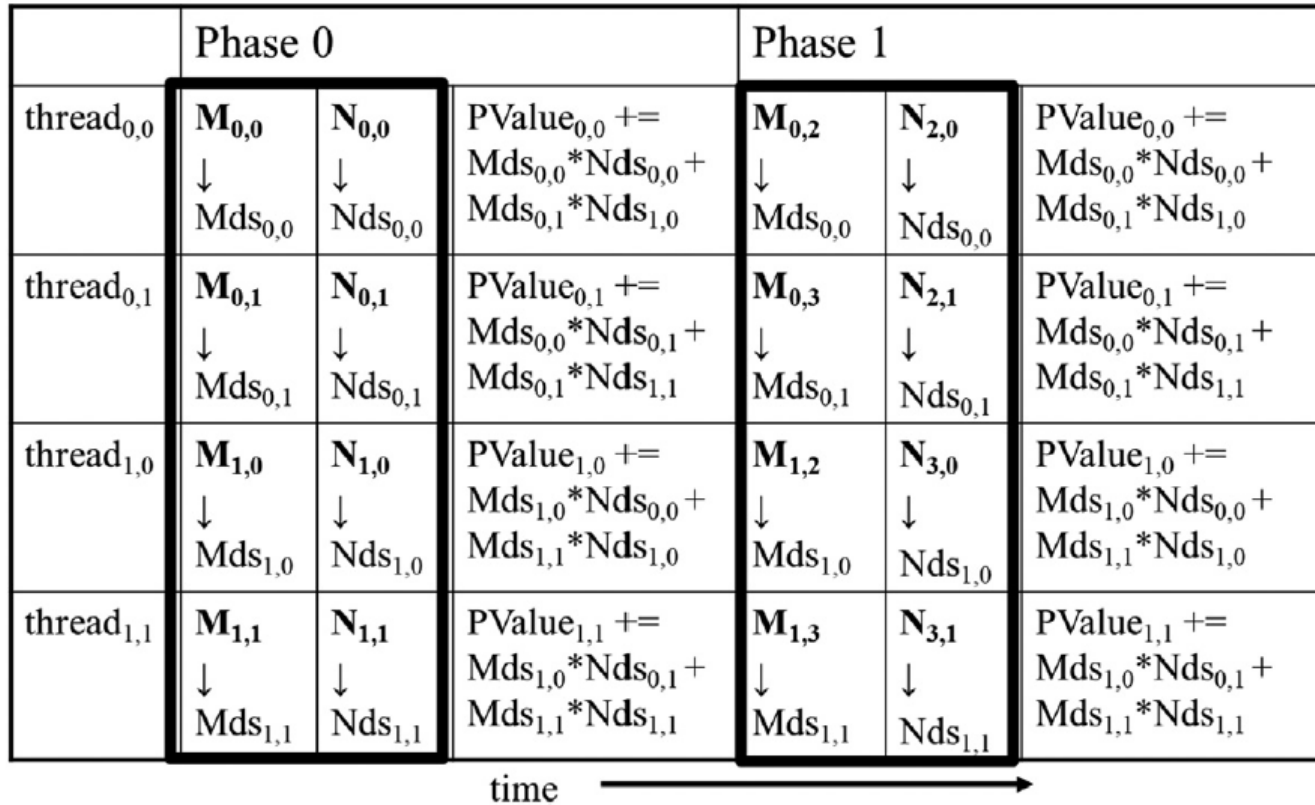


FIGURE 5.8

Execution phases of a tiled matrix multiplication.

```

01  #define TILE_WIDTH 16
02  __global__ void matrixMulKernel(float* M, float* N, float* P, int Width) {
03
04      __shared__ float Mds[TILE_WIDTH][TILE_WIDTH];
05      __shared__ float Nds[TILE_WIDTH][TILE_WIDTH];
06
07      int bx = blockIdx.x;  int by = blockIdx.y;
08      int tx = threadIdx.x; int ty = threadIdx.y;
09
10      // Identify the row and column of the P element to work on
11      int Row = by * TILE_WIDTH + ty;
12      int Col = bx * TILE_WIDTH + tx;
13
14      // Loop over the M and N tiles required to compute P element
15      float Pvalue = 0;
16      for (int ph = 0; ph < Width/TILE_WIDTH; ++ph) {
17
18          // Collaborative loading of M and N tiles into shared memory
19          Mds[ty][tx] = M[Row*Width + ph*TILE_WIDTH + tx];
20          Nds[ty][tx] = N[(ph*TILE_WIDTH + ty)*Width + Col];
21          __syncthreads();
22
23          for (int k = 0; k < TILE_WIDTH; ++k) {
24              Pvalue += Mds[ty][k] * Nds[k][tx];
25          }
26          __syncthreads();
27      }
28      P[Row*Width + Col] = Pvalue;
29
30
31  }

```

FIGURE 5.9

A tiled matrix multiplication kernel using shared memory.

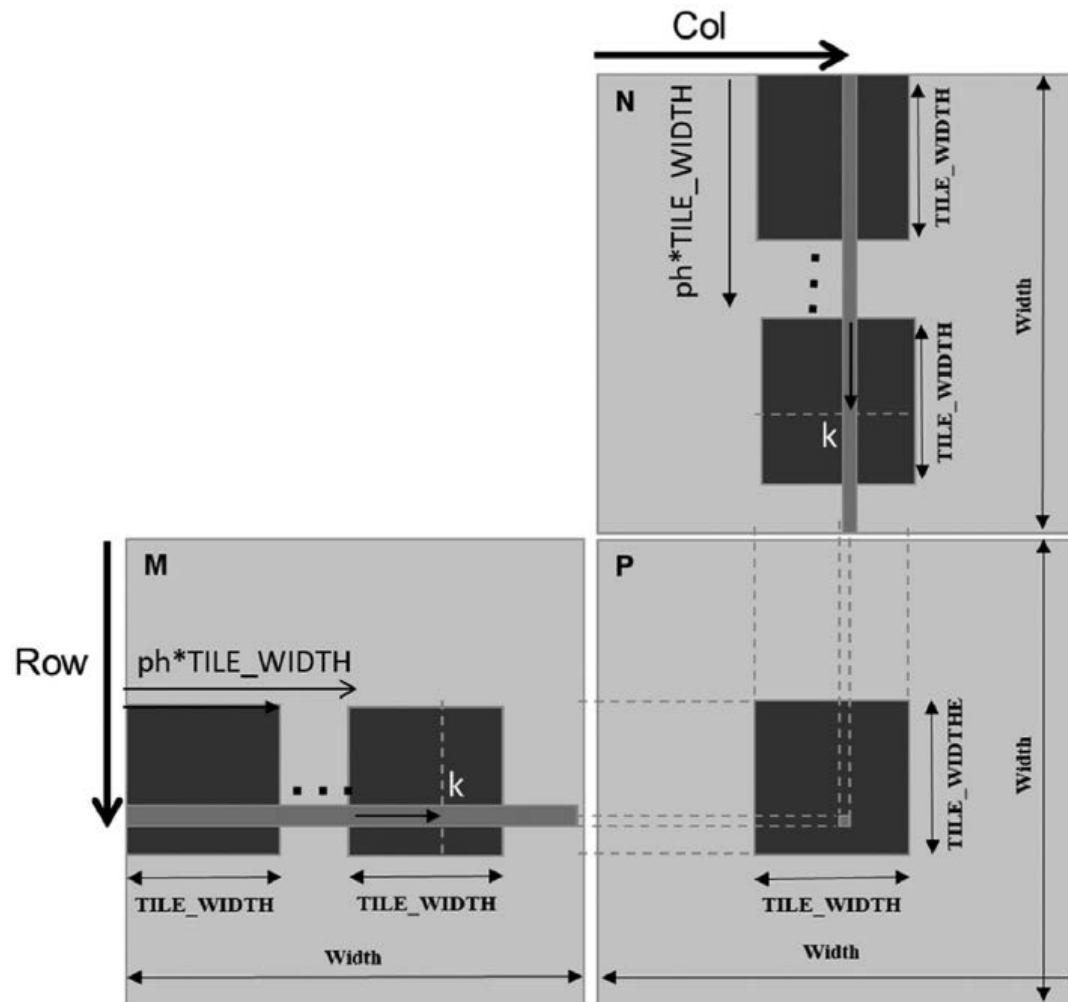


FIGURE 5.10

Calculation of the matrix indices in tiled multiplication.

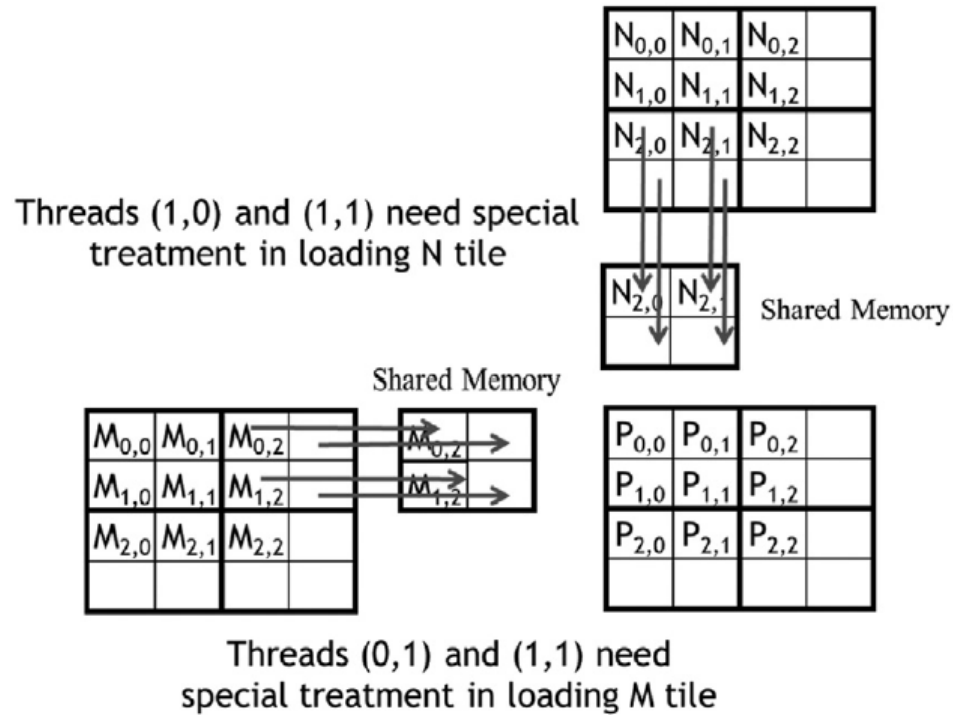


FIGURE 5.11

Loading input matrix elements that are close to the edge: phase 1 of block_{0,0}.

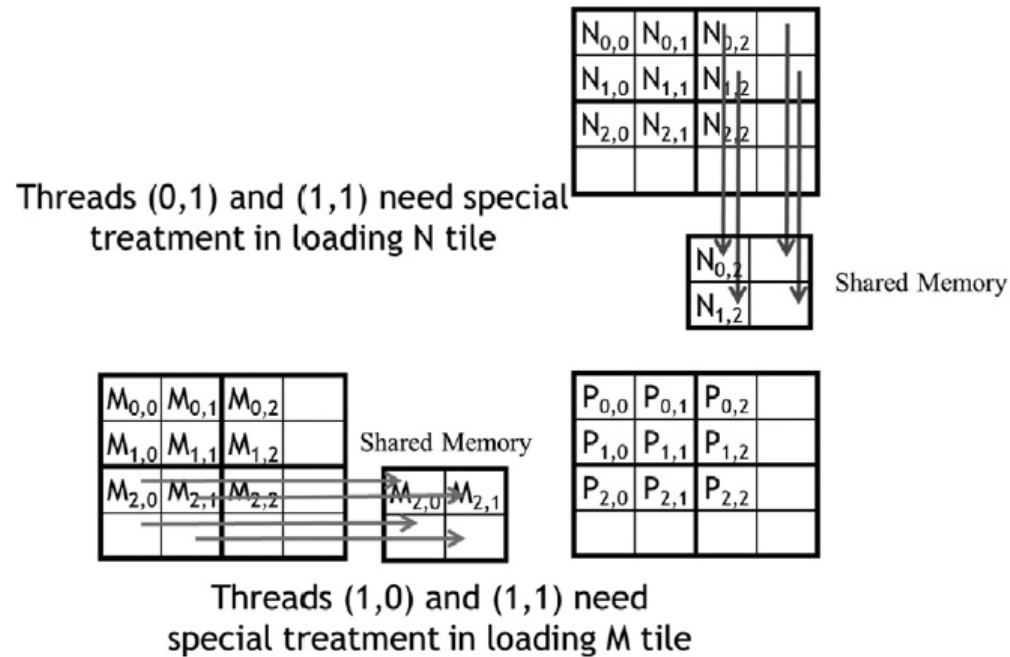


FIGURE 5.12

Loading input elements during phase 0 of block_{1,1}.

```

14 // Loop over the M and N tiles required to compute P element
15 float Pvalue = 0;
16 for (int ph = 0; ph < ceil(Width/(float)TILE_WIDTH); ++ph) {
17
18     // Collaborative loading of M and N tiles into shared memory
19     if ((Row < Width) && (ph*TILE_WIDTH+tx) < Width)
20         Mds[ty][tx] = M[Row*Width + ph*TILE_WIDTH + tx];
21     else Mds[ty][tx] = 0.0f;
22     if ((ph*TILE_WIDTH+ty) < Width && Col < Width)
23         Nds[ty][tx] = N[(ph*TILE_WIDTH + ty)*Width + Col];
24     else Nds[ty][tx] = 0.0f;
25     __syncthreads();
26
27     for (int k = 0; k < TILE_WIDTH; ++k) {
28         Pvalue += Mds[ty][k] * Nds[k][tx];
29     }
30     __syncthreads();
31 }
32 if (Row < Width) && (Col < Width)
33     P[Row*Width + Col] = Pvalue;

```

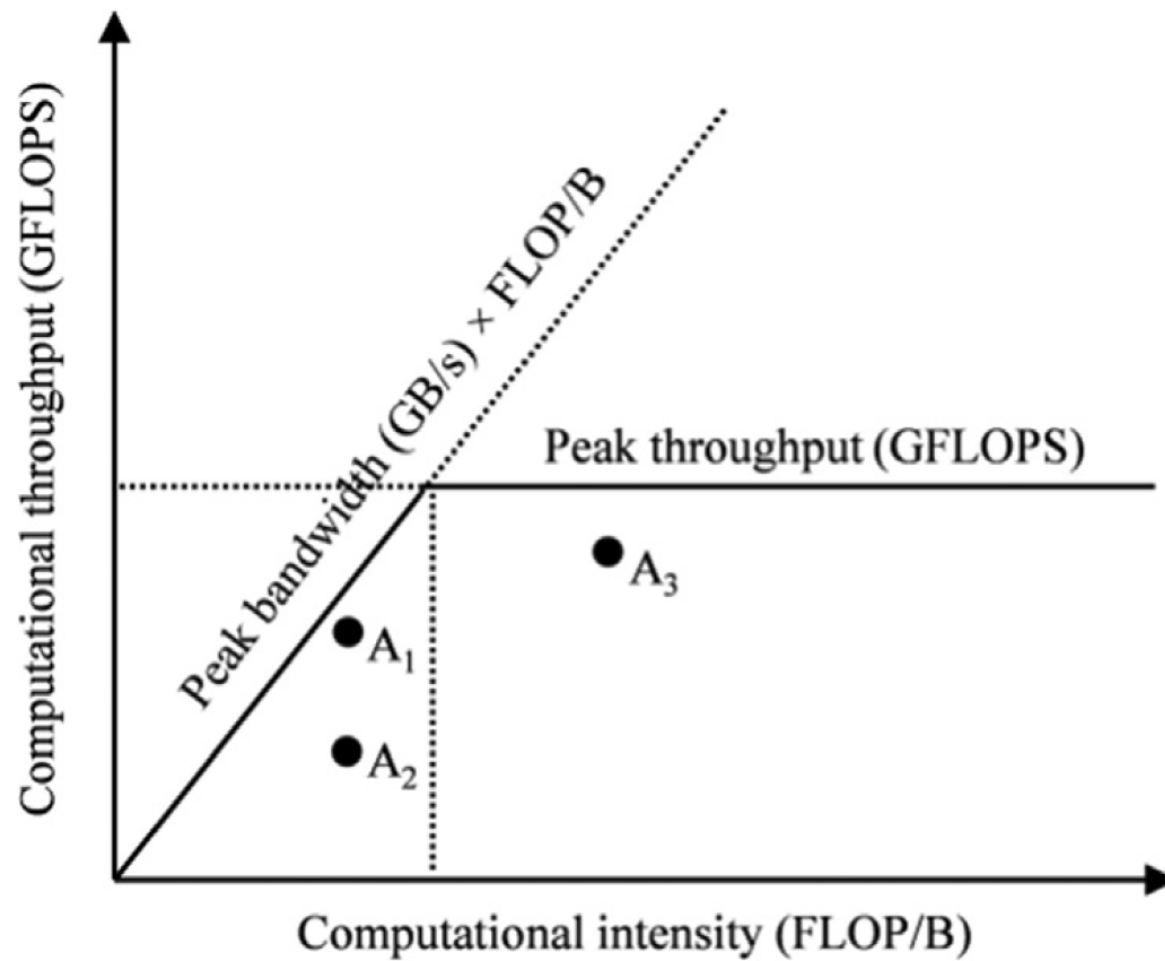
FIGURE 5.13

Tiled matrix multiplication kernel with boundary condition checks.

```
01  #define TILE_WIDTH 16
02  __global__ void matrixMulKernel(float* M, float* N, float* P, int Width,
                                unsigned Mdz_sz, unsigned Nds_sz) {
03
04      extern __shared__ char float Mds_Nds[];
05
06      float *Mds = (float *) Mds_Nds;
07      float *Nds = (float *) Mds_Nds + Mds_sz;
```

FIGURE 5.14

Tiled matrix multiplication kernel with dynamically sized shared memory usage.



In-text figure 1


```
fadd r1, r2, r3
```

In-text figure 2

```
load r2, r4, offset  
fadd r1, r2, r3
```

```
__shared__ float Mds[TILE_WIDTH][TILE_WIDTH];  
__shared__ float Nds[TILE_WIDTH][TILE_WIDTH];
```

In-text figure 4

```
#define TILE_WIDTH 16
```

In-text figure 5

```
extern __shared__ Mds_Nds[];
```

In-text figure 6

```
size_t size =  
calculate_appropriate_SM_usage(devProp.sharedMemPerBlock,  
...);  
  
matrixMulKernel<<<dimGrid,dimBlock,size>>>(Md, Nd, Pd,  
Width, size/2, size/2);
```

In-text figure 7

```

01 dim3 blockDim(BLOCK_WIDTH,BLOCK_WIDTH);
02 dim3 gridDim(A_width/blockDim.x,A_height/blockDim.y);
03 BlockTranspose<<<gridDim, blockDim>>>(A, A_width, A_height);

04 __global__ void
05 BlockTranspose(float* A_elements, int A_width, int A_height)
06 {
07     __shared__ float blockA[BLOCK_WIDTH][BLOCK_WIDTH];

08     int baseIdx = blockIdx.x * BLOCK_SIZE + threadIdx.x;
09     baseIdx += (blockIdx.y * BLOCK_SIZE + threadIdx.y) * A_width;

10     blockA[threadIdx.y][threadIdx.x] = A_elements[baseIdx];

11     A_elements[baseIdx] = blockA[threadIdx.x][threadIdx.y];
12 }

```

In-text figure 8

```

01  __global__ void foo_kernel(float* a, float* b) {
02      unsigned int i = blockIdx.x*blockDim.x + threadIdx.x;
03      float x[4];
04      __shared__ float y_s;
05      __shared__ float b_s[128];
06      for(unsigned int j = 0; j < 4; ++j) {
07          x[j] = a[j*blockDim.x*gridDim.x + i];
08      }
09      if(threadIdx.x == 0) {
10          y_s = 7.4f;
11      }
12      b_s[threadIdx.x] = b[i];
13      __syncthreads();
14      b[i] = 2.5f*x[0] + 3.7f*x[1] + 6.3f*x[2] + 8.5f*x[3]
15            + y_s*b_s[threadIdx.x] + b_s[(threadIdx.x + 3)%128];
16  }
17  void foo(int* a_d, int* b_d) {
18      unsigned int N = 1024;
19      foo_kernel <<< (N + 128 - 1)/128, 128 >>>(a_d, b_d);
20  }

```

In-text figure 9

Table 5.1 CUDA variable declaration type qualifiers and the properties of each type.

Variable declaration	Memory	Scope	Lifetime
Automatic variables other than arrays	Register	Thread	Grid
Automatic array variables	Local	Thread	Grid
__device__ __shared__ int SharedVar;	Shared	Block	Grid
__device__ int GlobalVar;	Global	Grid	Application
__device__ __constant__ int ConstVar;	Constant	Grid	Application