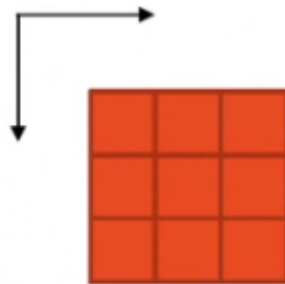


## Matrix Configurations

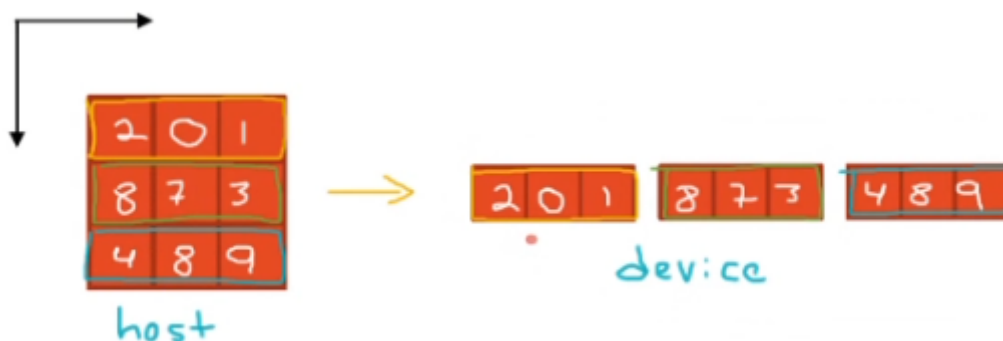
- The ideal is to configure the block/grid according to the data you need to process: when working with vectors, configure grid/blocks as vectors, when working with images, configure grid/blocks as matrix.
  - When working with one image, the block config must be a matrix too.



**Figure 1:** Image

- The kernel must always receive a **vector of information** as a parameter even though our data is a matrix, like an image. Inside the kernel we then need to unfold the matrix so that we can access an image like a vector. Thus, we need to transform the matrix to a vector only in the process of **data transference** between host and device, so that the device receives it as a vector then. But the kernel can be configured as a matrix:

```
1 dim3 grid(1);
2 dim3 block(3,3); // (3,3,1)
```



**Figure 2:** Image

- The idea is to send  $n \times m$  threads for a  $n \times m$  image.
- The coordinates to locate a thread in a single-block grid inside its 2D block are `(threadIdx.x, threadIdx.y)`. So, if we have a config in the form of a matrix, we need

to unfold this 2D matrix block config in order to calculate the global ID, and this id will be used to access the vector we have as a param. The globalId now is calculated as:

```
gId = threadIdx.x + threadIdx.y * blockDim.x
```

- The `threadIdx.y * blockDim.x` tells you how many rows to skip downwards through the 'Y' component.



Figure 3: Image

### Example 01

Sum of matrices

```

1  #include "cuda_runtime.h"
2  #include "device_launch_parameters.h"
3
4  #include <stdio.h>
5  #include <stdlib.h>
6
7  __host__ void checkCUDAError(const char* msg) {
8      cudaError_t error;
9      cudaDeviceSynchronize();
10     error = cudaGetLastError();
11     if (error != cudaSuccess) {
12         printf("ERROR %d: %s (%s)\n", error,
13             cudaGetErrorString(error), msg);
14     }
15 }
16
17 __global__ void matrixSum(int* dev_a, int* dev_b, int* dev_c)
18 {
19     int gId = threadIdx.x + threadIdx.y * blockDim.x;
20     dev_c[gId] = dev_a[gId] + dev_b[gId];
21 }
22
23 int main() {
24     const int N = 3; // if 32 ok, if 33 ERROR 9: invalid
25                       // configuration argument (matrixSum kernel error) and c
26                       // mat is zeroed

```

```
24     int* host_a = (int*)malloc(sizeof(int) * N * N);
25     int* host_b = (int*)malloc(sizeof(int) * N * N);
26     int* host_c = (int*)malloc(sizeof(int) * N * N);
27
28     int* dev_a, * dev_b, * dev_c;
29     cudaMalloc((void**)&dev_a, sizeof(int) * N * N);
30     cudaMalloc((void**)&dev_b, sizeof(int) * N * N);
31     cudaMalloc((void**)&dev_c, sizeof(int) * N * N);
32
33     // init data
34     for (int i = 0; i < N * N; i++) {
35         host_a[i] = (int)(rand() % 10);
36         host_b[i] = (int)(rand() % 10);
37     }
38
39     cudaMemcpy(dev_a, host_a, sizeof(int) * N * N,
40               cudaMemcpyHostToDevice);
41     cudaMemcpy(dev_b, host_b, sizeof(int) * N * N,
42               cudaMemcpyHostToDevice);
43
44     dim3 block(N, N);
45     dim3 grid(1);
46
47     matrixSum << < grid, block >> > (dev_a, dev_b, dev_c);
48     checkCUDAError("matrixSum kernel error");
49
50     cudaMemcpy(host_c, dev_c, sizeof(int) * N * N,
51               cudaMemcpyDeviceToHost);
52
53     printf("\nMatrix A: \n");
54     for (int i = 0; i < N; i++) {
55         for (int j = 0; j < N; j++) {
56             printf("%d ", host_a[j + i * N]);
57         }
58         printf("\n");
59     }
60
61     printf("\nMatrix B: \n");
62     for (int i = 0; i < N; i++) {
63         for (int j = 0; j < N; j++) {
64             printf("%d ", host_b[j + i * N]);
65         }
66         printf("\n");
67     }
68
69     printf("\nMatrix C: \n");
70     for (int i = 0; i < N; i++) {
71         for (int j = 0; j < N; j++) {
72             printf("%d ", host_c[j + i * N]);
73         }
74         printf("\n");
75     }
```

```

72     }
73
74     free(host_a);
75     free(host_b);
76     free(host_c);
77     cudaFree(dev_a);
78     cudaFree(dev_b);
79     cudaFree(dev_c);
80 }

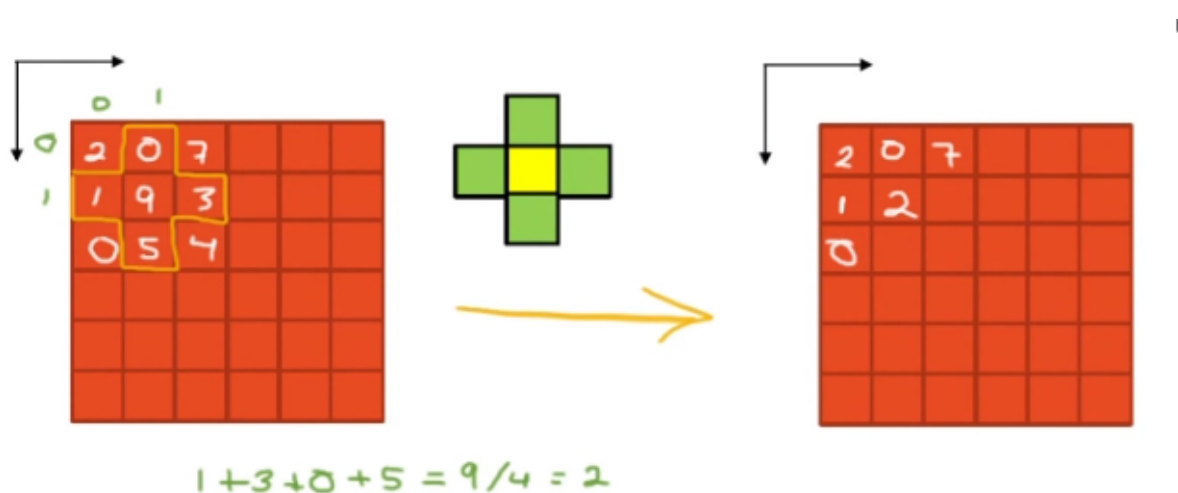
```

- The property `maximumThreadsPerBlock` will tell us how big the matrix can be in order to have a thread per cell. For a 1024 limit, the matrix would be 32 x 32.

## Image Processing: blur mask

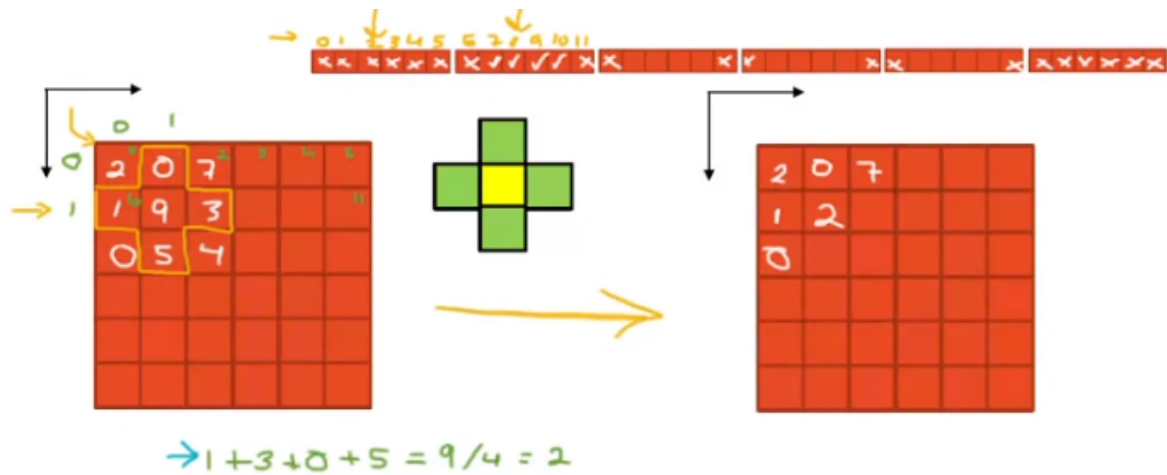
The objective is, given a matrix of information, apply a blur filter using that matrix.

- Without considering the borders:



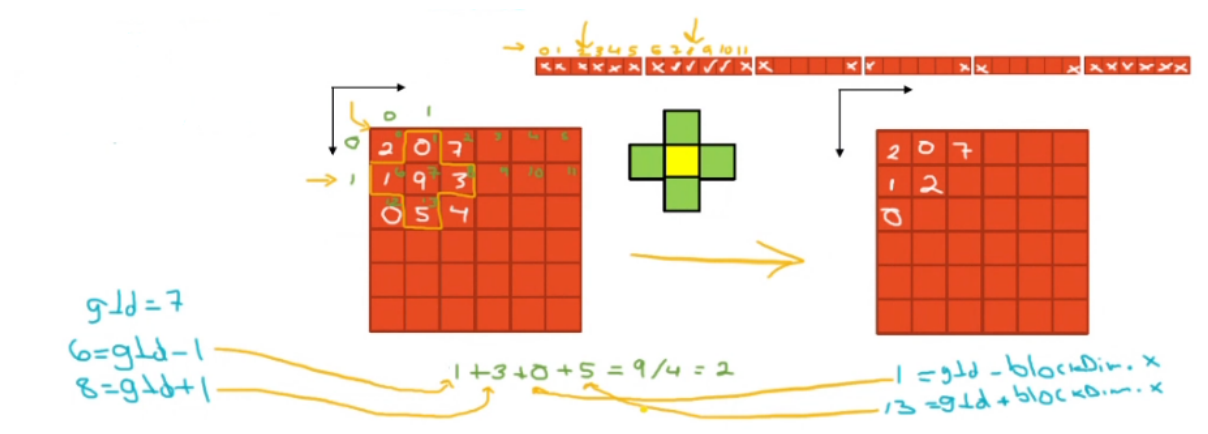
**Figure 4:** Image

- The border threads will not do anything, and this is managed by knowing its global Id. In order to do that, we need to unfold the matrix block config as a vector



**Figure 5:** Image

- Now, we will use the calculated gId to get the element of the vector of information that we need to sum as a neighbour, so that each cell can now contain the average of its four neighbours.



**Figure 6:** Image