# Error Management in CUDA

CUDA provides a way to handle errors that involve **exceeded GPU capacities** or **GPU malfunctioning**: the hardest errors to find out. These are not logic nor syntax errors.

- `cudaError_t` is a CUDA type given to handle errors, that is really an integer, and this number gives us a hint about the possible error. Every CUDA function returns an error that can be stored in this type. The only CUDA function that doesnt return a `cudaError_t` variable is a kernel itself: it must return void.

```
1  cudaError_t error;
2  error = cudaMalloc((void**)&ptr, size);
```

## Types of Errors (Most Common)

- `cudaSuccess` = 0: The API call returned with no errors. In query calls, this also means the query is complete. Successful execution.

- `cudaErrorInvalidValue` = 1: This indicates that one or more parameters passed to the API function call is not within an acceptable range of values. An enum param in a CUDA function that you didnt match, or a different data type sent.

- `cudaErrorMemoryAllocation` = 2: the API call failed because it was unable to allocate enough memory to perform the requested operation. When you do not have/allocate enough space in kernel memory for a requested instruction: you would normally write on memory outside of your array, but if there is no more mem left, this happens.

- `cudaErrorInvalidConfiguration` = 9: This indicates that a kernel launch is requesting resources that can never be satisfied with the current device. Requesting too many shared memory per block than supported, as well as requesting too many threads or blocks. This happens when you have an invalid kernel config (grid/blocks): when you exceed the max num of blocks per grid or threads of the GPU card.

- `cudaErrorInvalidMemcpyDirection` = 21: the direction of the memcpy passed to API call is not one of the types specified by cudaMemcpyKind. You put another word, basically.

## Process

```
1  cudaError_t error;
2  error = cudaMalloc((void**)&ptr, size);
3  cudaGetErrorString(error);
4
5  // after a kernel launch
```

```
6  error = cudaGetLastError();
7  cudaGetErrorString(error);
```

- `error` would be an intger, but in order to avoid checking in the docs, CUDA provides the function `cudaGetErrorString(error)` that, given an integer, it returns a string of the error details/explanation.

- To get a kernel error (otherwise a kernel is just void return), we can catch the last integer of error with `cudaGetLastError();`.

Open any project an type this funtion that will be used after any CUDA function call:

```
1  // host because every function call must be from host
2  __host__ void checkCUDAError(const char* msg){
3      cudaError_t error;
4      cudaDeviceSynchronize(); // avoid catching another error
           that is not the next we want
5      error = cudaGetLastError(); // status of the last CUDA API
           call (maybe 0 or success, not error)
6      if (error != cudaSuccess){
7          printf("ERROR %d: %s (%s)\n", error,
               cudaGetErrorString(error), msg);
8      }
9  }
```

Because the host/device execution is asynchronous (both at the same time), we need to take care of the sequence and sometimes you need to **pause** and wait for the kernel to finish in order to come back to the host: we need to synchronize.

## Example

```
1  #include "cuda_runtime.h"
2  #include "device_launch_parameters.h"
3  #include <stdio.h>
4  #include <stdlib.h> /* srand, rand */
5  #include <time.h> /* time */
6
7  #include<iostream>
8  using namespace std;
9
10 __host__ void checkCUDAError(const char* msg) {
11     cudaError_t error;
12     cudaDeviceSynchronize();
13     error = cudaGetLastError();
14     if (error != cudaSuccess) {
15         printf("ERROR %d: %s (%s)\n", error,
               cudaGetErrorString(error), msg);
16     }
17 }
```

```
18
19  __global__ void idKernel(int* vecA, int* vecB, int* vecC) {
20      int gId = threadIdx.x + blockDim.x * blockIdx.x;
21
22      vecA[gId] = threadIdx.x;
23      vecB[gId] = blockIdx.x;
24      vecC[gId] = gId;
25  }
26
27  void printArray(int* arr, int size, char* msg) {
28      cout << msg << ": ";
29      for (int i = 0; i < size; i++) {
30          printf("%d ", arr[i]);
31      }
32      printf("\n");
33  }
34
35  int main()
36  {
37      const int vectorSize = 64;
38      int* host_a = (int*)malloc(sizeof(int) * vectorSize);
39      int* host_b = (int*)malloc(sizeof(int) * vectorSize);
40      int* host_c = (int*)malloc(sizeof(int) * vectorSize);
41
42      int* dev_a, * dev_b, * dev_c;
43
44      cudaMalloc((void**)&dev_a, sizeof(int) * vectorSize);
45      checkCUDAError("Error at cudaMalloc for dev_a");
46      cudaMalloc((void**)&dev_b, sizeof(int) * vectorSize);
47      checkCUDAError("Error at cudaMalloc for dev_b");
48      cudaMalloc((void**)&dev_c, sizeof(int) * vectorSize);
49      checkCUDAError("Error at cudaMalloc for dev_c");
50
51      srand(time(NULL));
52
53      for (int i = 0; i < vectorSize; i++) {
54          host_a[i] = 0;
55          host_b[i] = 0;
56          host_c[i] = 0;
57      }
58
59      cudaMemcpy(dev_a, host_a, sizeof(int) * vectorSize,
              cudaMemcpyHostToDevice);
60      checkCUDAError("Error at cudaMemcpy for host_a to dev_a");
61      cudaMemcpy(dev_b, host_b, sizeof(int) * vectorSize,
              cudaMemcpyHostToDevice);
62      checkCUDAError("Error at cudaMemcpy for host_b to dev_b");
63      cudaMemcpy(host_c, dev_c, sizeof(int) * vectorSize,
              cudaMemcpyHostToDevice); // error 1
64      checkCUDAError("Error at cudaMemcpy for host_c to dev_c");
65
```

```
66      dim3 grid(1, 1, 1);
67      dim3 block(2000, 1, 1); // max num is 1024, so here we
            will force an error
68      idKernel << < grid, block >> > (dev_a, dev_b, dev_c);
69      checkCUDAError("Error at idKernel execution no. 1");
70      cudaDeviceSynchronize(); // wait until kernel finishes and
            then come back to following code // not needed to
            check
71      cudaMemcpy(host_a, dev_a, sizeof(int) * vectorSize,
            cudaMemcpyDeviceToHost);
72      //check also here
73      cudaMemcpy(host_b, dev_b, sizeof(int) * vectorSize,
            cudaMemcpyDeviceToHost);
74      cudaMemcpy(host_c, dev_c, sizeof(int) * vectorSize,
            cudaMemcpyDeviceToHost);
75
76      printf("Execution 1: 1 block 64 threads \n");
77      printArray(host_a, vectorSize, "threadIdx.x");
78      printArray(host_b, vectorSize, "blockIdx.x");
79      printArray(host_c, vectorSize, "globalId");
80
81      grid.x = 64; // (64, 1, 1)
82      block.x = 1; // (1, 1, 1)
83      idKernel << < grid, block >> > (dev_a, dev_b, dev_c);
84      cudaDeviceSynchronize();
85      cudaMemcpy(host_a, dev_a, sizeof(int) * vectorSize,
            cudaMemcpyDeviceToHost);
86      cudaMemcpy(host_b, dev_b, sizeof(int) * vectorSize,
            cudaMemcpyDeviceToHost);
87      cudaMemcpy(host_c, dev_c, sizeof(int) * vectorSize,
            cudaMemcpyDeviceToHost);
88
89      printf("\nExecution 2: 64 blocks 1 thread \n");
90      printArray(host_a, vectorSize, "threadIdx.x");
91      printArray(host_b, vectorSize, "blockIdx.x");
92      printArray(host_c, vectorSize, "globalId");
93
94      grid.x = 4;
95      block.x = 16;
96      idKernel << < grid, block >> > (dev_a, dev_b, dev_c);
97      cudaDeviceSynchronize();
98      cudaMemcpy(host_a, dev_a, sizeof(int) * vectorSize,
            cudaMemcpyDeviceToHost);
99      cudaMemcpy(host_b, dev_b, sizeof(int) * vectorSize,
            cudaMemcpyDeviceToHost);
100     cudaMemcpy(host_c, dev_c, sizeof(int) * vectorSize,
            cudaMemcpyDeviceToHost);
101
102     printf("\nExecution 3: 4 block 16 threads \n");
103     printArray(host_a, vectorSize, "threadIdx.x");
104     printArray(host_b, vectorSize, "blockIdx.x");
```

```
105       printArray(host_c, vectorSize, "globalId");
106
107       free(host_a);
108       free(host_b);
109       free(host_c);
110       cudaFree(dev_a);
111       cudaFree(dev_b);
112       cudaFree(dev_c);
113       return 0;
114 }
```

**Output**

```
1  ERROR 9: invalid configuration argument (Error at idKernel
       execution no. 1)
2  Execution 1: 1 block 64 threads
3  threadIdx.x: 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
       0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
       0 0 0 0 0 0 0 0
4  blockIdx.x: 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
       0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
       0 0 0 0 0 0 0 0
5  globalId: 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
       0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
       0 0 0 0 0 0 0 0
6
7  Execution 2: 64 blocks 1 thread
8  threadIdx.x: 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
       0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
       0 0 0 0 0 0 0 0
9  blockIdx.x: 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19
       20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39
       40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58
       59 60 61 62 63
10 globalId: 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20
       21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39
       40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59
       60 61 62 63
11
12 Execution 3: 4 block 16 threads
13 threadIdx.x: 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 0 1 2 3 4 5
       6 7 8 9 10 11 12 13 14 15 0 1 2 3 4 5 6 7 8 9 10 11 12 13
       14 15 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15
14 blockIdx.x: 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1
       1 1 1 1 1 1 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 3 3 3 3 3 3 3 3
       3 3 3 3 3 3 3 3
15 globalId: 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20
       21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39
       40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59
       60 61 62 63
```

*ERROR 9: Kernel no. 1 is not really executed*