# Warp Divergence

**Warp Divergence** happens when you launch the threads in a warp to do **different things** among them. For example, if you, inside a kernel, differentiate the activities of the threads using their global ID and send, for example, the even threads to do this and the odds to do something else. Last time, in parallel reduction, you provoked this divergence.

The problem with this divergence is that it theoretically reduces **velocity by half** than the one your kernel would have if threads in a warp didn't diverge, **per if condition**: if you have one if statement, you reduce the speed by half; two if statements, you reduce that half speed by half again, and so on.

What can we do to avoid this? Instead of using the **gId**, use the **warpId** to divide the tasks among the threads. This is not warp divergence, because **all threads of each warp are doing the same thing**.

- Example

This has divergence:

```
1  if (gId % 2 == 0){
2      // activity 1
3  } else {
4      // activity 2
5  }
```

This has no divergence:

```
1  int warpId = gId / 32;
2  if (warpId % 2 == 0){
3      // activity 1
4  } else {
5      // activity 2
6  }
```

Especifically, in parallel reduction we cannot use this to avoid divergence. When you cannot avoid divergence, another alternative is to fix this divergence by using the space of memory called **Shared Memory**. Up until now, we have just used **global memory** when we have variables inside the kernel. The memory we used with `cudaMalloc()` was the memory called **Global Memory**.

**Shared Memory** is independent memory per block, where all threads in a block can access the memory of *only* the block they belong to. It is a faster memory in terms of access (read/write) when compared to global memory, because it is closer to the processor: therefore it is recommended when you have data that is constantly used or referenced

inside a kernel. Parallel Redutcion is apt to this, because the elements of a vector are constantly referenced. The advise would be to store the vector in **Shared Memory**.

Nevertheless, we have less Shared Memory than the memory we have as Global Memory. You can check the capacity of these memories using the `cudaGetDeviceProperties()`, and it is given by block. Take care that your program does not surpass the capacity of Shared Memory when you write a kernel that requires it.

```
1  #include "cuda_runtime.h"
2  #include "device_launch_parameters.h"
3  #include <stdio.h>
4
5  int main() {
6      cudaDeviceProp prop;
7      cudaGetDeviceProperties(&prop, 0);
8      size_t sharedMemory = prop.sharedMemPerBlock;
9      printf("sharedMemPerBlock: %zd bytes\n", sharedMemory); //
           %zd is used to print size_t values
10 }
```

Which outputs:

```
1  sharedMemPerBlock: 49152 bytes
```

## Parallel Reduction With Divergence But Using Shared Memory Fix
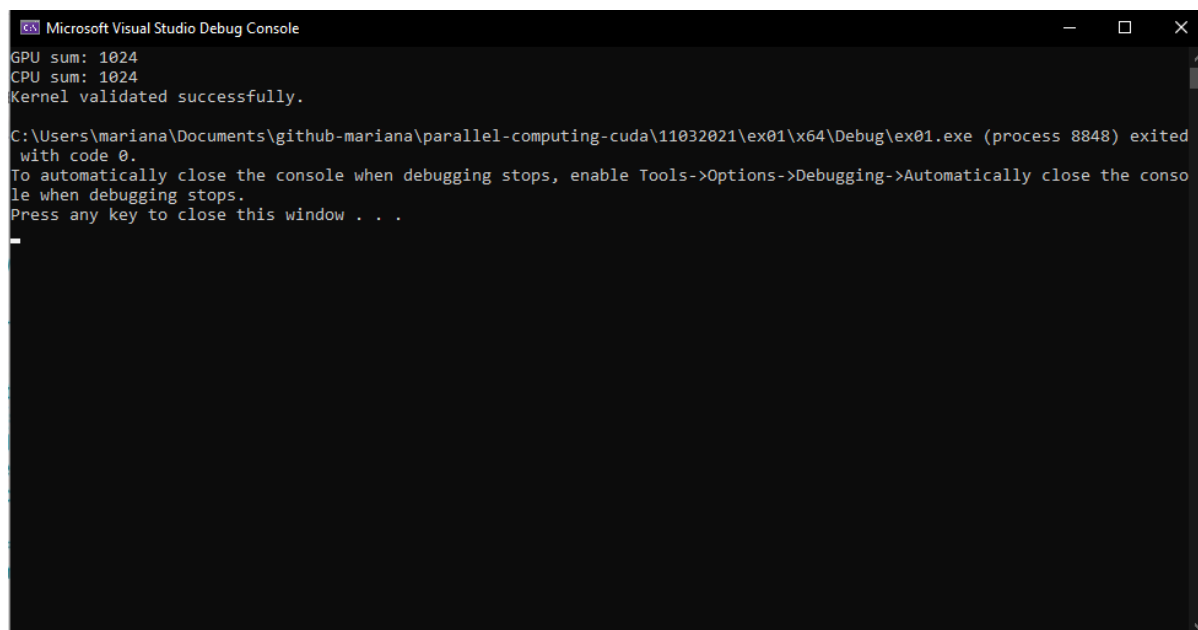
### Solution

- 1D Grid with 1D block of 1024 threads in the x axis.

```
1  #include "cuda_runtime.h"
2  #include "device_launch_parameters.h"
3
4  #include <stdio.h>
5  #include <stdlib.h>
6  #define vecSize 1024
7
8  __host__ void checkCUDAError(const char* msg) {
9      cudaError_t error;
10     cudaDeviceSynchronize();
11     error = cudaGetLastError();
12     if (error != cudaSuccess) {
13         printf("ERROR %d: %s (%s)\n", error,
               cudaGetErrorString(error), msg);
14     }
15 }
16
17 __host__ void validate(int* result_CPU, int* result_GPU, int
       size) {
```

```
18      if (*result_CPU != *result_GPU) {
19          printf("The results are not equal.\n");
20          return;
21      }
22      printf("Kernel validated successfully.\n");
23      return;
24  }
25
26  __host__ void CPU_fn(int* v, int* sum, const int size) {
27      for (int i = 0; i < size; i++) {
28          *sum += v[i];
29      }
30  }
31
32  __global__ void kernel_divergent(int* v, int* sum) {
33      int gId = threadIdx.x;
34      int step = vecSize / 2;
35
36      while (step) {
37          if (gId < step) {
38              v[gId] = v[gId] + v[gId + step];
39          }
40          step = step / 2;
41      }
42      if (gId == 0) {
43          *sum = v[gId];
44      }
45  }
46
47  __global__ void kernel_divergent_fixed(int* v, int* sum) {
48      int gId = threadIdx.x;
49      // vector in shared memory
50      __shared__ int vectorShared[vecSize];
51      vectorShared[gId] = v[gId];
52      // old GPU's need thread synchronize when using shared
            memory: shared memory is visible for all blocks
53      __syncthreads();
54      int step = vecSize / 2; // avoid using blockDim to avoid
            unexpected behaviours
55
56      while (step) {
57          //printf("%d\n", step);
58          if (gId < step) {
59              vectorShared[gId] = vectorShared[gId] +
                    vectorShared[gId + step];
60          }
61          step = step / 2;
62      }
63      if (gId == 0) {
64          *sum = vectorShared[gId];
65      }
```

```
66  }
67
68  int main() {
69
70      const int size = 1024;
71      int* v = (int*)malloc(sizeof(int) * size);
72      int sumCPU = 0;
73      int sumGPU = 0;
74
75      int* dev_v, * sum;
76      cudaMalloc((void**)&dev_v, sizeof(int) * size);
77      cudaMalloc((void**)&sum, sizeof(int));
78
79      for (int i = 0; i < size; i++) {
80          v[i] = 1;
81      }
82
83      cudaMemcpy(dev_v, v, sizeof(int) * size,
              cudaMemcpyHostToDevice);
84      cudaMemcpy(sum, &sumGPU, sizeof(int),
              cudaMemcpyHostToDevice);
85
86      dim3 grid(1);
87      dim3 block(size);
88
89      kernel_divergent_fixed << < grid, block >> > (dev_v, sum);
90      cudaMemcpy(&sumGPU, sum, sizeof(int),
              cudaMemcpyDeviceToHost);
91      printf("GPU sum: %d\n", sumGPU);
92
93      CPU_fn(v, &sumCPU, size);
94      printf("CPU sum: %d\n", sumCPU);
95
96      validate(&sumCPU, &sumGPU, size);
97
98      return 0;
99  }
```

**Output**



**Figure 1:** img