



AGH

**AGH UNIVERSITY OF SCIENCE
AND TECHNOLOGY**

Introduction to CUDA and Open CL

Lab 3

Michał Kunkel
Wiktor Żychowicz

1. Page fault.

This is an exception when running program trying to access a memory page isn't mapped by the memory management unit. When it happens, kernel tries to reach the required page at the location in physical memory, or crushes the program due to an illegal memory access.

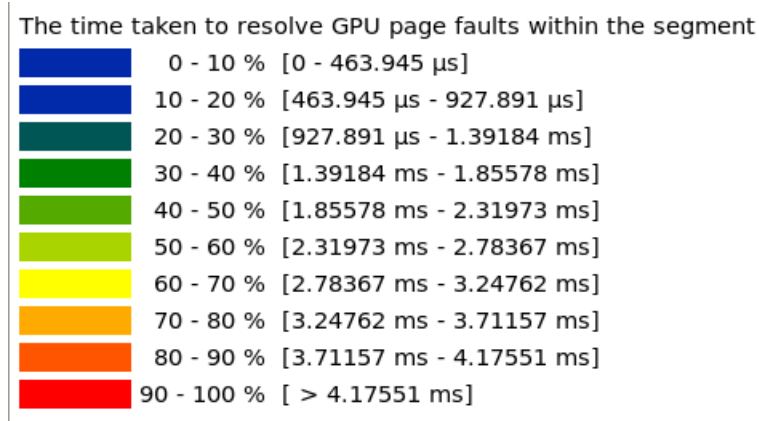
Naming may wrongly suggest that page fault is something bad, but in fact page fault are not always errors. This is the common way for hardware to enlarge memory available to programs in any operating system that uses virtual memory.

On labs we were supposed to research how CPU and GPU behaves when they meet page faults(there are consequence of using cudaMallocManager and "Unified memory") by modifying this code:

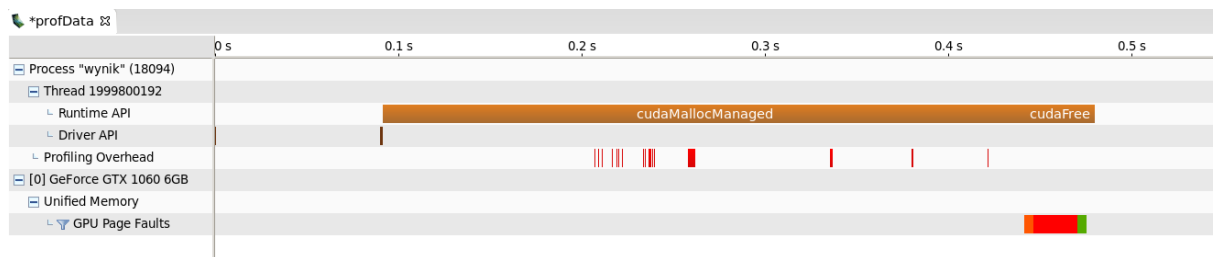
```
1 global
2 void deviceKernel(int *a, int N)
3 {
4     int idx = threadIdx.x + blockIdx.x * blockDim.x;
5     int stride = blockDim.x * gridDim.x;
6     for (int i = idx; i < N; i += stride)
7     {
8         a[i] = i;
9     }
10 }
11
12 void hostFunction(int *a, int N)
13 {
14     for (int i = 0; i < N; ++i)
15     {
16         a[i] = i;
17     }
18 }
19
20 int main()
21 {
22     int N = 7<<10;
23     size_t size = N * sizeof(int);
24     int *a;
25     cudaMallocManaged(&a, size);
26
27     /*
28      * Conduct experiments to learn more about the behavior of
29      * 'cudaMallocManaged'.
30      *
31      * What happens when unified memory is accessed only by the GPU?
32      * What happens when unified memory is accessed only by the CPU?
33      * What happens when unified memory is accessed first by the GPU then the CPU?
34      * What happens when unified memory is accessed first by the CPU then the GPU?
35      *
36      * Hypothesize about UM behavior, page faulting specifically, before each
37      * experiment, and then verify by running 'nvprof'.
38      */
39
40     cudaFree(a);
41 }
42
43 memory_page_faults.cu" 44L, 953C
```

Picture 1: Source code of a program on which we tested.

2. Examples.

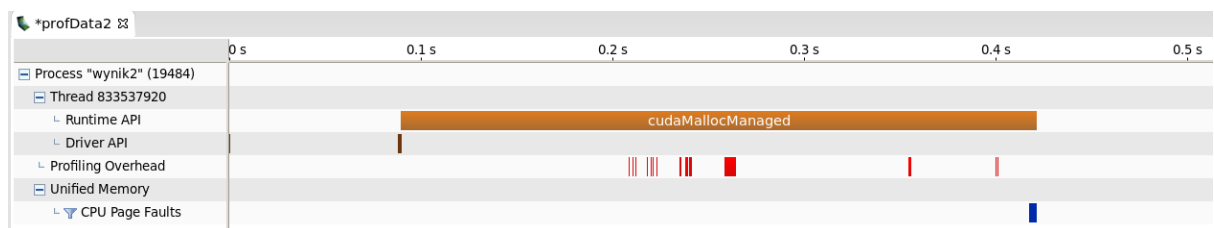


Picture 2.0: Legend of colours in NVIDIA profiler.



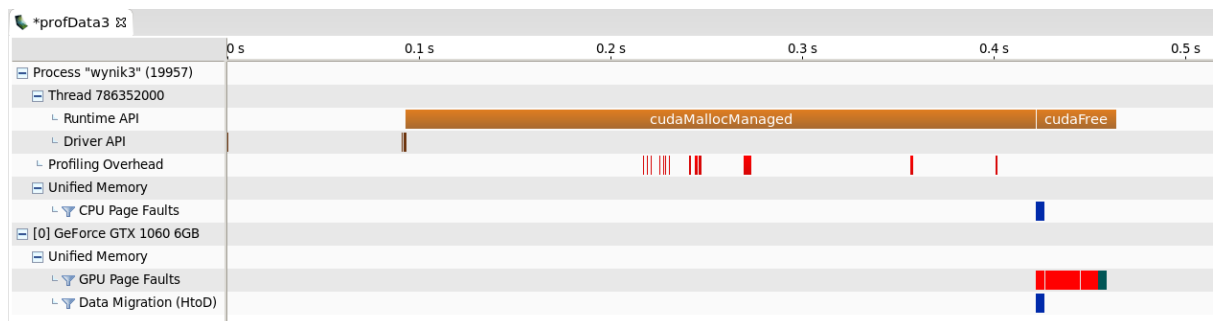
Picture 2.1: Unified memory accessed only by the GPU.

In first case kernel takes a lot of time to overcome (over 4 ms) page fault problems. Despite of them program compiles and works as it should. They do harm only to performance.



Picture 2.2: Unified memory accessed only by the CPU.

In second case kernel managed to reached to problematic area of memory and managed to deal with page fault and for CPU it took less then 1ms. When we compare this example to others we can clearly see that CPU deals with page faults faster than GPU.



Picture 2.3: Unified memory accessed first by GPU then the CPU

In third case we can observe that CPU can help GPU in reaching inaccessible location in physical memory for GPU. After 0.01-0.05 seconds. As we can see to resolve this page fault some data movement is required.

No picture 2.4: Unified memory accessed first by CPU then the GPU

In that case when we execute file we get a segmentation fault which causes it to crash. Sometimes page faults can be fatal to our programs.

3. Some optimization.

Finally we play with some ways we can optimize code. We took simple program which sums up vectors. In all cases we use 'stride'. It is a strategy of doing loops in parallel manner where we instead of launching one thread to one iteration we use in-thread loop to take next elements when it end executing previous(good size of step is $\text{gridDim} * \text{blockDim}$) without waiting for other threads. We tested asynchronized prefetching of data and result. We also check how place where we initialize vectors affect efficiency.