
REPORT: PROJECT WEBCAM STREAMING & FILTERS

January 5, 2022

Matteo Battilana, Salvatore Gabriele La Greca, Giovanni Pollo

Contents

1	Introduction	1
2	Video Streaming	1
2.1	CPU Implementation	2
2.2	GPU Implementation	3
2.2.1	Naif version	4
2.2.2	Asynchronous APIs and removal of useless operations	7
2.2.3	Pinned or page-locked memory	9
2.2.4	Kernel global memory access: the cherry on the cake	10
2.2.5	Using of shared memory for frequent accesses: is it avoidable?	15
3	Image Text Overlay	17
3.1	GPU Implementation	18
3.1.1	Naif version	18
3.1.2	Exploiting the global memory access optimization	20
3.1.3	Constant Memory instead of Global Memory: is it worth?	21
4	HeatMap	22
4.1	Heat map pixel mapping - CPU Naif implementation	23
4.2	CUDA implementation	24
4.2.1	CUDA Naif implementation	25
4.2.2	CUDA switching frames	27
4.2.3	CUDA Global memory access granularity	27
4.3	Evaluation of the number of threads	29
4.3.1	Heatmap Conclusions	31
5	Noise visualizer	31
5.0.1	Naif CPU implementation	32
5.0.2	CUDA implementation	33
5.0.3	CUDA reduction of global memory accesses	34
5.0.4	Overlapping implementation	37
5.1	Fastest implementation	39
6	Noise filter	40
6.1	Convolution filters	41
6.2	Mean filter	41
6.3	Gaussian filter	42
6.4	Tiled CUDA 2D implementation	43
6.5	2D CUDA Tiled Median filter	47

1 Introduction

A classical video streaming algorithm over internet or over any other communication channel is based on the concept of sending not frame by frame as they are but instead they are based on sending the first frame and then the difference between the new one and the previous, where difference here is to be intended pixel-by-pixel difference.

If this can be arguable for video with a lower resolution, the heaviness of sending each frame as it is for high video resolution is notable. Let's take as an example a FULL HD video, means that each frame is composed of 1920x1080 pixels where 1920 is the width of the frame while 1080 is the height. Supposing the frame is in RGB24 format, means that each pixel is represented by 3 byte (one for each channel R, G, and B).

By doing a rapid computation, each frame measures $3 \cdot 1920 \cdot 1080B = 6220800B = 5.93MB$. Supposing now the video used by example is a 30 fps video, means that each second we have 30 frame, each of one measures 5.93 MB: each second we are sending about 178 MB. To send 178 MB/s we need a transfer link bandwidth of 1492 Mbps, that is unfeasible.

So the solution is to send the difference, and this means sending only the pixels that change or, better, pixels where their difference is above a certain threshold.

The purpose of the project is to demonstrate the performances that are obtainable by computing the difference on a CPU and to compare them with the ones obtainable by using a General Purpose GPU or *GPGPU*. A Nvidia GPU will be used for the benchmark and therefore the code will be based on CUDA. Among all these considerations, different filters will be added in order to demonstrate the potentiality of GPGPU computations on video elaboration and streaming.

2 Video Streaming

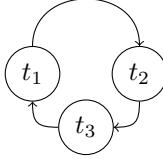
Before looking at the algorithm that computes the difference between frames, is important to give a look at how frames are captured and the overall architecture of the software. It's a client server architecture where the server transmits the difference to the client via a socket.

All the different variants of implementation that will be analyzed produces three outputs:

- `h_pos`: the number of pixels that are different and greater than a certain threshold.
- `h_diff`: in code it is actually the array of the current frame that after the application of the algorithm contains the difference of each byte.
- `h_xs`: it's a mapping vector for the `h_diff`. This means that the `h_diff[0]` is the difference of byte element at position `h_xs[0]`.

In order to capture frame by frame from the webcam and to visualize them, OpenCV is used. It's not so efficient in terms of performances, especially on the platform used (Nvidia Jetson Nano with 4 ARM cores @ 1.5 GHz) but for the purpose of this project it will be fine. The most important thing to underline is that a frame will be represented by the OpenCV's object `Mat` that contains, among other informations, the dimension of the image (that is fixed to a FULL HD resolution) and an array of `uint8_t` items, each representing a channel of a pixel for each pixels of the image. The array can be allocated automatically at the creation of the `Mat` object or an external array can be used and later on this feature will be exploited.

The aim is to have an efficient software, so a multi thread approach is adopted. There are, in fact, 3 different threads each of them with a different purpose: capture, elaborate and send. They work in a circular way. Means that the capture thread is a producer for the elaborate one, the elaborate one is a producer for the send one and the last one is a producer for the capture one. In this way we have all the threads working at the same time on a different task. The t2 thread is where the magic happens so where the elaboration of the difference is executed.



In the following pages, for metric considerations these terms will be used:

- **fps**: number of frames per second.
- **read**: time of execution of the capture thread.
- **for**: time of execution of the elaboration thread. This is what is under discussion in this project.

2.1 CPU Implementation

The CPU implementation is the easier one and the most basic implementation of the algorithm. It consists on a loop between each byte that compose the two frames (the current one and the previous one) and compute the difference, storing it in a new vector.

The C++ implementation is the following:

```

1 int total = 3 * ctx.sampleMat->rows * ctx.sampleMat->cols; // no. bytes
2 Mat pvs = pready->pframe->clone();
3 pready->h_pos = 0;
4
5 for (int i = 0; i < total; i++) {
6     int df = pready->pframe->data[i] - previous.data[i];
7     if (df < -LR_THRESHOLDS || df > LR_THRESHOLDS) {
8         pready->pframe->data[pready->h_pos] = df;
9         pready->h_xs[pready->h_pos] = i;
10        pready->h_pos++;
11    } else {
12        pvs.data[i] -= df;
13    }
14}
15
16 previous = pvs;
17

```

The code here is pretty simple. For each byte, the difference **df** is computed. If this difference is greater than a fixed threshold **LR_THRESHOLDS** means that it is a valid difference so it can be sent. For all the benchmarks proposed in the following, the value of **LR_THRESHOLDS** is fixed to 20. The result of the reconstructed frame at the client side can be seen at Figure 1.



Figure 1: *The result of the reconstructed frame*

The important point here is that if the difference is too low, it can't be simply discarded, so a kind of negative feedback is needed. This is the purpose of the line of code at line 7, where the value of the byte of the current frame (that at the end of the loop will be the previous for the next iteration) is itself minus the value of the difference. This means that at the next iteration, if that value changes again and its difference increases it will take under consideration as a big difference. Without this, there will be a sum of errors in the reconstructed image, leading in a wrong visualization. The result if the error is not considered can be seen at Figure 2.

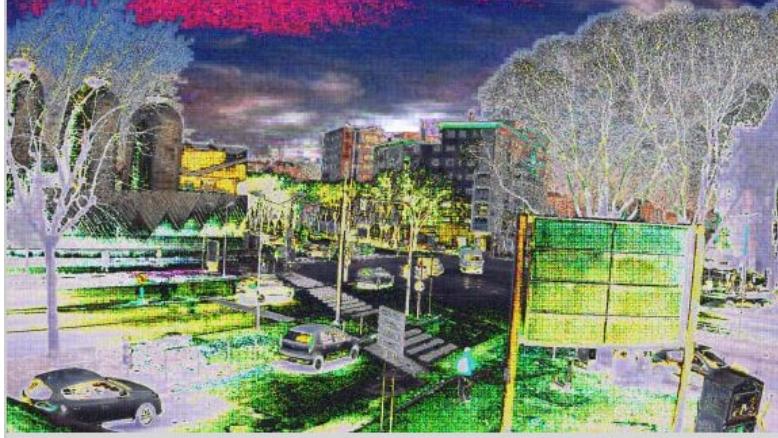


Figure 2: *What happens after a certain time if the error is not take under consideration*

The performances here are pretty bad. By means of in-code time measurements, the video streaming is stable at 7 fps with an average `for` time of 140.0 ms and an average `read` time of 0.0 ms. This means that here the big bottleneck is due to the elaboration part of the software.

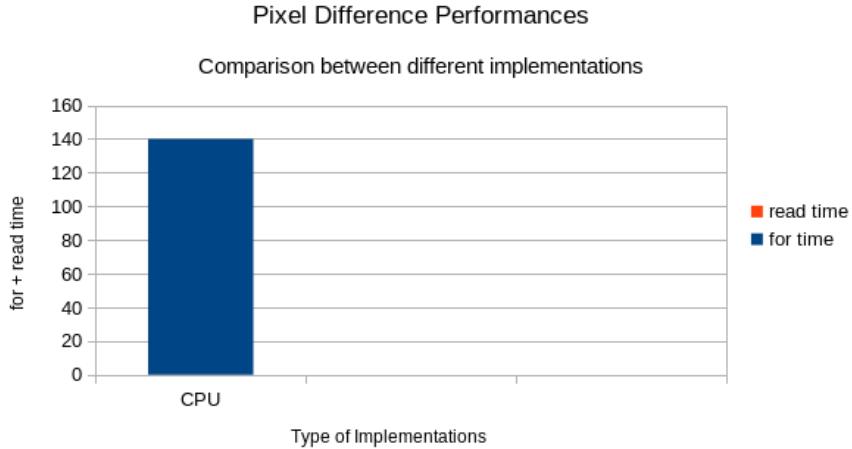


Figure 3: *Pixel Difference Comparison*

2.2 GPU Implementation

The idea is now to port this in the GPU. The GPU can be seen as an accelerator for the CPU, an accelerator that is capable of a high degree of parallelism by executing a lot of simple threads all in parallel in a SIMD way.

In order to port this in CUDA and execute the computation on the GPU, a series of considerations needs to be done. First of all, it's not said that by executing this in GPU there is a direct improvement in performance. This is due to the fact that there is a big bottleneck between the CPU and the GPU and this is the bus that connects them. A GPU in order to do some work

on some data needs to have them in its memory (so the GPU can't access directly the CPU's RAM).

The first thing to do, before defining the kernel (the piece of code that is offloaded to the GPU), is to understand how split the 6220800 bytes among the concurrent threads of the GPU. In the Mat object, the frame is represented row-major, means that rows are at consecutive address, as shown in Table 1.

$\begin{bmatrix} 0 & 1 & 2 \\ 3 & 4 & 5 \\ 6 & 7 & 8 \end{bmatrix}$	$\begin{bmatrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ \hline 0x0 & 0x1 & 0x2 & 0x3 & 0x4 & 0x5 & 0x6 & 0x7 & 0x8 \end{bmatrix}$
---	---

Table 1: Row major representation in memory. On top the value, on bottom the address

Threads can be organized differently according how they access the memory, and in this case there can be two possible cases of memory access per thread:

Thread 0	0	3	6
Thread 1	1	4	7
Thread 2	2	5	8

Table 2: Non consecutive access

Thread 0	0	1	2
Thread 1	3	4	5
Thread 2	6	7	8

Table 3: Consecutive access

So the entire array can be divided into $<N>$ chunks where $<N>$ is the number of threads to launch on the GPU. By organizing the kernels in a way they can access the memory in a consecutive access, the kernel will achieve the so called *memory coalesced access*.

How many threads? For a first implementation, and in order to have a modular implementation of the code (so it can be executed on different GPUs with different capabilities), the following code is adopted:

```

1 struct cudaDeviceProp prop;
2 cudaGetDeviceProperties(&prop, 0); // retrieves device infos
3
4 int total = 3 * ctx.sampleMat->rows * ctx.sampleMat->cols; // no. bytes
5 int nMaxThreads = prop.maxThreadsPerBlock;
6 int maxAtTime = total / nMaxThreads;
7

```

According to this, each thread will work on $N = maxAtTime$ consecutive bytes, dividing the memory into $nMaxThreads$ chunks, and overlapping is totally avoided. The GPU used for the experiments allows a maximum of 1024 threads to be executed concurrently.

For what concerns the memory allocation, obviously the CPU address space and the GPU one are two separated things, so a specific allocation on the device side must be done. This can be accomplished thanks to some CUDA's API as follows:

```

1 uint8_t *d_current, *d_previous, *d_diff;
2 int *d_xs;
3 unsigned int *d_pos;
4
5 cudaMalloc((void**)&d_diff, total * sizeof *d_diff);
6 cudaMalloc((void**)&d_xs, total * sizeof *d_xs);
7 cudaMalloc((void**)&d_current, total * sizeof *d_current);
8 cudaMalloc((void**)&d_previous, total * sizeof *d_previous);
9 cudaMalloc((void**)&d_pos, sizeof *d_pos);
10

```

2.2.1 Naif version

The naif version of the implementation is to port as it is the CPU code into a kernel. From the host side (the CPU) that asks to the device (the GPU) to execute the kernel, there are three kind of operations to be done:

1. copy the frame into the GPU's memory
2. launch the kernel
3. copy back the results from the GPU's memory to the HOST memory

In the naif implementation, both the previous and current frames are copied into the GPU's memory. At the end of the execution, the results are copied back. This implies 2 transfers HostToDevice (HtoD) for current and previous and 3 transfers DeviceToHost (DtoH) for h_xs, h_diff and h_pos. A memset on the device side is required too to set the initial value of d_pos to 0.

The host will therefore execute the following operations:

```

1 cudaMemset(d_pos, 0, sizeof *d_pos);
2
3 // launch
4 cudaMemcpy(d_previous, previous.data, total, cudaMemcpyHostToDevice);
5 cudaMemcpy(d_current, pready->pframe->data, total, cudaMemcpyHostToDevice);
6 kernel<<<1,nMaxThreads>>>(d_current, d_previous, d_diff, maxAtTime, d_pos, d_xs);
7
8 // copy back
9 cudaMemcpy(pready->pframe->data, d_diff, total, cudaMemcpyDeviceToHost);
10 cudaMemcpy(pready->h_xs, d_xs, total*sizeof *d_xs, cudaMemcpyDeviceToHost);
11 cudaMemcpy(&pready->h_pos, d_pos, sizeof *d_pos, cudaMemcpyDeviceToHost);
12
13 // copy previous
14 cudaMemcpy(previous.data, d_current, total, cudaMemcpyDeviceToHost);
15

```

It's possible to notice a pretty similar structure of the CPU implementation, where the *for loop* is replaced by the kernel. The kernel is launched with a grid dimension of (1, 0, 0) and a block dimension of (*nMaxThreads*, 0, 0), thus only the x dimension is used because the kernel is coded in such a way that the chunk assigned to it is analyzed in a linear way.

In order to define the code of the kernel, there are two main problems to be analyzed first:

1. How does the kernel knows at which address it must work?
2. How is possible to fill a variable size array in a concurrent environment in CUDA?

The first problem is easily solved, it's only a matter of indexes. Each kernel needs to know which chunk of the memory is assigned to it so first of all it computes its own unique index that will space from 0 to *nMaxThreads* – 1. Once this is defined, the kernel knows that it must work on the i-th chunk but this needs to be translated into the real index to be used to address the memory. Because each size is *maxAtTime* large, means that the the real index that identifies the start of the chunk is *threadIndex* · *maxAtTime*.

The second problem is a bit more complicated to solve. In the CPU implementation a variable h_pos is used in order to indicize the resulting vectors an incremental way and at the end this variable is used to know the length of the result. The problem when this concept is trasposed to the GPU world is that there is a global variable d_pos that is used by *nMaxThreads* concurrently so there is the need to read and increment this variable in a safe way, so only a kernel at a time can ready and modify it. Luckily, the CUDA framework offers a series of functions that ensures atomic operations between kernels, and in this case the needed one is `atomicInc(unsigned int *address, unsigned int val)` where *address* is the pointer to the value that needs to be atomic incremented, *val* is the wraparound value, so the max value that the variable can reach before going again to 0, and it returns as result the value atomically incremented.

After solved those problems, a look at the real implementation of the kernel can be given:

```

1 __global__ void kernel(
2     uint8_t *current, uint8_t *previous, uint8_t *diff,
3     int maxSect, unsigned int *pos, int *xs) {
4     int x = threadIdx.x + blockDim.x * blockIdx.x;
5     unsigned int npos;
6     int df;
7

```

```

8
9     int start = x * maxSect;
10    int max = start + maxSect;
11
12 #pragma unroll
13    for (int i = start; i < max; i++) {
14
15        df = current[i] - previous[i];
16        if (df < -LR_THRESHOLDS || df > LR_THRESHOLDS) {
17            npos = atomicInc(pos, 6220801);
18            diff[npos] = df;
19            xs[npos] = i;
20        } else {
21            current[i] -= df;
22        }
23    }
24}
25
26}
27

```

As can be seen, apart the computation of the start and end index of the computation, the for loop is pretty much the same as the CPU implementation. The only difference is the one concerning the use of the `atomicInc` function as explained before.

The performances here are slightly better than the CPU implementation. The video streaming is elaborated at around 13 fps with an average `for` time of 70.0 ms and an average `read` time of 0.0 ms, and the comparison with the previous case can be seen in Figure 4.

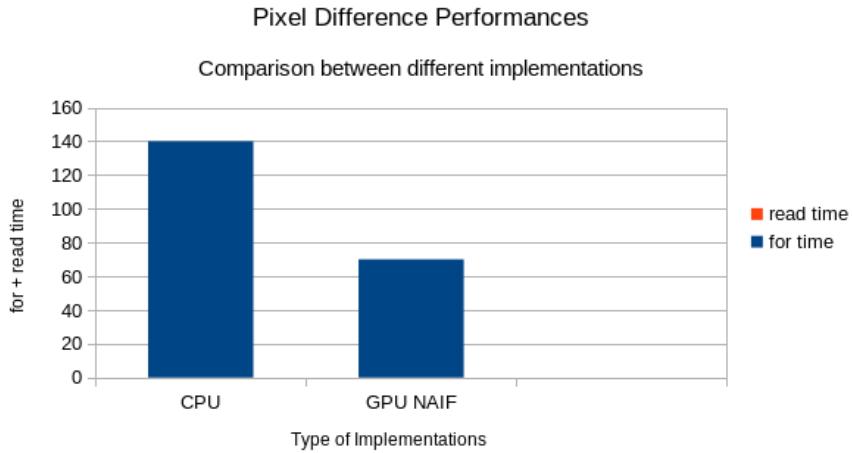


Figure 4: *Pixel Difference Comparison*

This means that the big bottleneck is still due to the elaboration part of the software. Using the profiler offered by Nvidia `nvprof` executed for 30 seconds, it's possible to get a lot more details about the time required by each GPU operation, and in particular what concerns this project is the transfer time and the kernel execution time. The results are shown in table 4.

Operation	Avg (ms)	Min (ms)	Max (ms)
Kernel	40.295	35.912	47.033
DtoH	4.2975	1.5100	52.878
HtoD	3.3280	2.4544	25.808
Memset	0.0008	0.0007	0.0017

Table 4: Profiling results

2.2.2 Asynchronous APIs and removal of useless operations

In the previous implementation, a big portion of the execution time is covered by the kernel execution but the two data transfers are not negligible. The total average time is about 46 ms means that the *for* time is for the 60% due to the GPU computation and the remaining part is related to host tasks.

CUDA offers a variant of the APIs that are asynchronous with respect to the device execution. This means that when one of these APIs is called, the operation starts on the device side but the control returns immediately to the host despite the end of the operation on the GPU.

This means that it is possible to give to the device a series of operations to execute and meanwhile the host can work on its own operations. In this specific case the host doesn't have any task to complete so a function `cudaDeviceSynchronize()` is used to wait for the device to complete its operations. However this solution is slightly more efficient because it is possible to assign to the GPU all the jobs in a consecutive way than waiting instead of assign a job at a time.

```

1 cudaMemsetAsync(d_pos, 0, sizeof *d_pos);
2
3 cudaMemcpyAsync(d_previous, previous.data, total, cudaMemcpyHostToDevice);
4 cudaMemcpyAsync(d_current, pready->pframe->data, total, cudaMemcpyHostToDev);
5 kernel <<<1,nMaxThreads>>>(d_current, d_previous, d_diff, maxAtTime, d_pos, d_xs);
6
7 cudaMemcpyAsync(pready->pframe->data, d_diff, total, cudaMemcpyDeviceToHost);
8 cudaMemcpyAsync(pready->h_xs, d_xs, total*sizeof *d_xs, cudaMemcpyDevToHost);
9 cudaMemcpyAsync(&pready->h_pos, d_pos, sizeof *d_pos, cudaMemcpyDeviceToHost);
10
11 cudaMemcpyAsync(previous.data, d_current, total, cudaMemcpyDeviceToHost);
12 cudaDeviceSynchronize();
13

```

Another optimization that can be done is to remove a particular useless copy both HtoD and DtoH, and that is the one regarding the previous frame. In fact, after the first launch of the kernel, the previous frame is already stored in the GPU's memory so it's avoidable the back copy to the host and gain to the device. What can be done is to implement a simple pointer swap, in fact after the first execution the current frame will become the previous frame at the next round. So there is a pool of two pointers that alternatively works one as the current and the other one as the previous.

The host code becomes like this:

```

1 // current-previous swap
2 uint8_t *d_prev = d_current;
3 d_current = d_previous;
4 d_previous = d_current;
5
6 cudaMemsetAsync(d_pos, 0, sizeof *d_pos);
7
8 // Copy in the current pointer and run
9 cudaMemcpyAsync(d_current, pready->pframe->data, total, cudaMemcpyHostToDev);
10 kernel <<<1,nMaxThreads>>>(d_current, d_previous, d_diff, maxAtTime, d_pos, d_xs);
11
12 cudaMemcpyAsync(pready->pframe->data, d_diff, total, cudaMemcpyDeviceToHost);
13 cudaMemcpyAsync(pready->h_xs, d_xs, total*sizeof *d_xs, cudaMemcpyDevToHost);
14 cudaMemcpyAsync(&pready->h_pos, d_pos, sizeof *d_pos, cudaMemcpyDeviceToHost);
15
16 cudaDeviceSynchronize();
17

```

Another optimization in terms of copies regards the download of the two arrays forming the result. In fact, instead of copying all the *total* bytes is smarter to cpy only *d_pos* bytes, the ones that are really useful.

So the first download copy regards *d_pos*, then a `cudaDeviceSynchronize()` is executed in order to wait for the complete of all the previous operation and getting a valid *d_pos*, then the remaining data are copied.

The host code becomes:

```

1 // current-previous swap
2 uint8_t *d_prev = d_current;
3 d_current = d_previous;
4 d_previous = d_current;
5
6 cudaMemsetAsync(d_pos, 0, sizeof *d_pos);
7
8 // Copy in the current pointer and run
9 cudaMemcpyAsync(d_current, pready->pframe->data, total, cudaMemcpyHostToDevice);
10 kernel<<<1,nMaxThreads>>>(d_current, d_previous, d_diff, maxAtTime, d_pos, d_xs);
11
12 cudaMemcpyAsync(&pready->h_pos, d_pos, sizeof *d_pos, cudaMemcpyDeviceToHost);
13 cudaDeviceSynchronize();
14
15 cudaMemcpyAsync(pready->pframe->data, pready->h_pos, total, cudaMemcpyDToHt);
16 cudaMemcpyAsync(pready->h_xs, d_xs, pready->h_pos* sizeof *d_xs, cudaMemcpyDToH);
17
18
19 cudaDeviceSynchronize();
20

```

Last but not least, the *cudaMemsetAsync* at line 6 can be removed completely. In fact, it's possible to move the initial reset of the *d_pos* counter by doing it at the beginning of the kernel and it's easily done even in a concurrent environment by letting only one kernel to perform the reset. This can be achieved by putting a condition on the kernel index. In this case it will be done by the first kernel. In order to have all the kernel executing more or less in parallel the same instructions, after the reset of the counter, a *__syncthreads()* is executed that assures that all threads are synched until that point.

Thus, the kernel code becomes:

```

1 __global__ void kernel(
2     uint8_t *current, uint8_t *previous, uint8_t *diff,
3     int maxSect, unsigned int *pos, int *xs) {
4     int x = threadIdx.x + blockDim.x * blockIdx.x;
5     unsigned int npos;
6     int df;
7
8     // counter reset done by thread with index = 0
9     if (!x) *pos = 0;
10    __syncthreads();
11
12    int start = x * maxSect;
13    int max = start + maxSect;
14
15    #pragma unroll
16    for (int i = start; i < max; i++) {
17
18        df = current[i] - previous[i];
19        if (df < -LR_THRESHOLDS || df > LR_THRESHOLDS) {
20            npos = atomicInc(pos, 6220801);
21            diff[npos] = df;
22            xs[npos] = i;
23        } else {
24            current[i] -= df;
25        }
26    }
27}
28
29
30}
31

```

All these optimizations let achieve a big result in performance terms. The video streaming is elaborated at around 20 fps with an average **for** time of 55.0 ms and an average **read** time of 0.0 ms, and the comparison with the previous case can be seen in Figure 5.

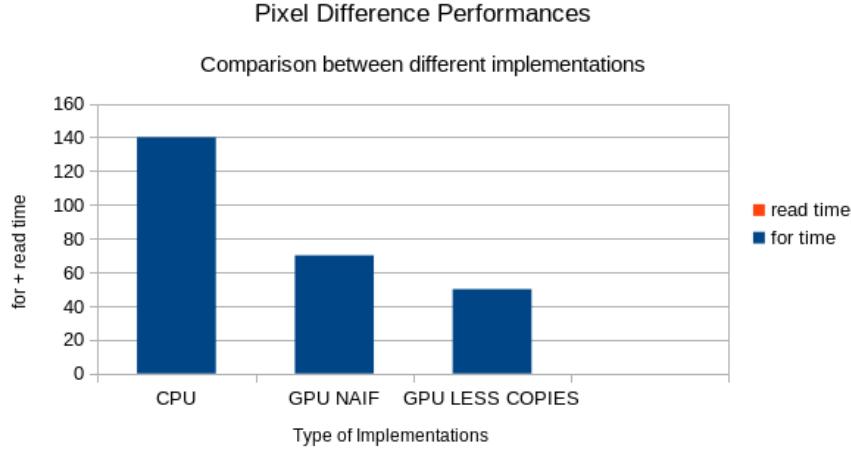


Figure 5: *Pixel Difference Comparison*

The results of the *nuprof* profiler are shown in Table 5 and is noticeable the big improvement in terms of performances for what concerns the DtoH transfer. The HtoD transfer remains more or less the same but it's executed only one time instead of two as the naif implementation.

Operation	Avg (ms)	Min (ms)	Max (ms)
Kernel	41.426	37.278	46.423
D _{to} H	0.0227	0.0009	9.2925
H _{to} D	3.7813	2.1733	14.188

Table 5: Profiling results

2.2.3 Pinned or page-locked memory

Until now a good bottleneck has been removed and it is the transfer from the Device to the Host. What about the Host to Device transfer? In order to solve this issue an understanding on how memory is managed by the CPU and how the GPU transfers data from the CPU's managed memory and the GPU.

What a modern operating system does is to manage the physical memory through the concept of pages and thus the virtual memory. The address space of the virtual memory is different from the physical one, thus a translation is needed. In order to let the hardware translate from a space to another, a special hardware called translation lookaside buffer (TLB) is used. It's a memory cache that is used to reduce the time taken to access a user memory location and it's part of the chip's memory-management unit (MMU) needed in fact by all the modern operating systems. The TLB stores the recent translations of virtual memory to physical memory and can be called an address-translation cache.

When a CUDA's *memcpy* is invoked, the DMA starts transferring data but as soon as the translation is not anymore stored in the cache, the operating system has to solve this situation by loading a new page. Thus this is not really a asynchronous operation and it's highly cpu-dependent. Thus the velocity of the copy is based on the speed of the CPU and in edge devices where the CPU is not so performant as a desktop CPU the overhead of this mechanism can't be neglected.

The solution is to allocate a piece of memory that is pinned or *page locked*. Executing an allocation like this is like telling the operating system virtual memory manager that the memory pages must stay in physical ram so that they can be directly accessed by the GPU across the PCI-express bus. And it's a lot faster thanks to the DMA: when the memory is page locked, the GPU DMA engine can directly run the transfer without requiring the host CPU, which reduces overall latency and decreases transfer times.

CUDA offers an API to do this, called `cudaMallocHost()` that works similar to `cudaMalloc()` but allocates a page-locked memory on the host side. It's required for arrays like `h_xs` and `pready->pframe->data`. If it's pretty easy to be done for `h_xs` that is an array directly managed by the code, how is possible to do that for the pframe's data that is something managed by OpenCV? As said at the beginning, the `Mat` object offers a good option, the one to select an external pointer as data array.

```

1 uint8_t *h_frame;
2 cudaMallocHost((void **)&h_frame, total * sizeof *h_frame);
3 cudaMallocHost((void **)&pready->h_xs, total * sizeof *pready->h_xs);
4
5 // init of Mat object with page-locked h_frame ptr
6 pready->pframe = new Mat(ctx.sampleMat->rows, ctx.sampleMat->cols,
7     ctx.sampleMat->type(), h_frame);
8

```

Without touching anything else, this optimization is able to achieve a good result. The video streaming is elaborated at around 22 fps with an average `for` time of 43.0 ms and an average `read` time of 0.0 ms, and the comparison with the previous case can be seen in Figure 6.

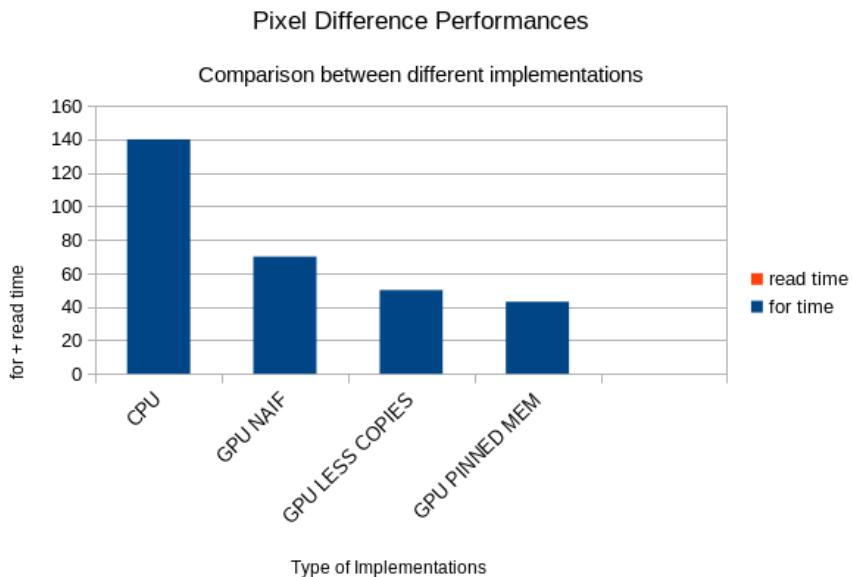


Figure 6: *Pixel Difference Comparison*

The results of the `nvprof` profiler are shown in Table 6 and is noticeable the big improvement in terms of performances for what concerns the HtoD transfer. In fact, it's about 1/4 the original one.

Operation	Avg (ms)	Min (ms)	Max (ms)
Kernel	43.386	38.093	49.607
D _H toH	0.0102	0.0008	1.0466
H _D toH	0.6740	0.6096	2.2818

Table 6: Profiling results

2.2.4 Kernel global memory access: the cherry on the cake

What remains now to further optimize is the kernel execution time. Doing a bunch of analysis, with a maximum number of threads of 1024 each kernel execute a *for loop* that lasts 6075 iterations. Each iteration access in memory three times in the most of the cases (supposing

$N_{pixelChanged} \ll N_{pixelNotChanged}$) and accesses in global memory are a real bit bottleneck for CUDA's kernel because usually the Global memory access time is $400 \sim 600$ cycles. Obviously, each byte needs to be examined so it's not possible to randomly cut accesses in memory. How this can be optimized?

The answer to this question is very simple but at the same time it's brilliant: by increase the size of each read or write in memory. By applying this concept, both the memory accesses and the loop iterations are reduced, so it's an optimization in both ways.

The first think to do is to rewrite the kernel in order to support this. The idea is to load into some GPU's registers some bytes regarding a portion of the current and previous frame under consideration by the for loop of a kernel. Then there is always need of a byte-wise subtraction and comparison so a inner for loop will access byte-wise the word loaded into the register(s) and will do its computations.

The multi-byte memory access is done using the dynamic pointer cast feature, casting the vector pointers to the data type aimed, Therefore the new kernel implementation, using the most common data type that is the `int` one, is:

```

1  typedef int chunk_t;
2
3  __global__ void kernel2(
4      uint8_t *current, uint8_t *previous, uint8_t *diff,
5      int maxSect, unsigned int *pos, int *xs) {
6          int x = threadIdx.x + blockDim.x * blockIdx.x;
7          unsigned int npos;
8          int df;
9          chunk_t cc, pc; // registers storing current copy, previous copy
10
11         if (!x) *pos = 0;
12         __syncthreads();
13
14         int start = x * maxSect;
15         int max = start + maxSect;
16
17         #pragma unroll
18         for (int i = start; i < max; i++) {
19
20             cc = ((chunk_t *)current)[i];
21             pc = ((chunk_t *)previous)[i];
22
23             #pragma unroll
24             for (int j = 0; j < sizeof cc; j++) {
25                 df = ((uint8_t *)&cc)[j] - ((uint8_t *)&pc)[j];
26
27                 if (df < -LR_THRESHOLDS || df > LR_THRESHOLDS) {
28                     npos = atomicInc(pos, 6220801);
29                     diff[npos] = df;
30                     xs[npos] = (i * sizeof cc) + j;
31                 } else {
32                     current[(i * sizeof cc) + j] -= df;
33                 }
34             }
35         }
36     }
37 }
38
39 }
```

The host code changes a little bit because `maxAtTime` now needs to be divided by the size of the `chunk_t` type. Because it's highly probable that the division results in a non integer value, it must be rounded to the nearest greater integer.

This means that by rounding up, each thread will overlap with the next thread of a few bytes. Actually it's not a problem because the probability that while the thread i writes at the very last bytes of its chunk the thread $i+1$ writes at the very beginning of its chunk is very low because

each thread is executed in parallel. The only thing that happens is that the computations on a certain amount of bytes is done two times but each time it's done on the same inputs and the result is always the same. In a low probability case where two writes at the same time happens, as discussed, the overhead in managing the concurrent write is lower than the one of adding instructions to check if the thread is writing outside its chunk (about 1 ms of difference in execution time).

```

1 int max4 = ceil(1.0 * maxAtTime / sizeof(chunk_t));
2
3 // CUDA APIs call
4 ...
5 kernel2<<<1,nMaxThreads>>>(d_current ,d_previous , d_diff , max4, d_pos , d_xs);
6 ...

```

The only problem to manage by choosing the path of not adding any check inside the kernel is that the last kernel can write outside the memory. In order to avoid wrong memory accesses, the allocation of the memory on the GPU is done by considering this margin. In fact, a thread will write at most `sizeof(chunk_t)` bytes outside the memory, so each memory allocation is done by allocating an additional amount of bytes equal to `sizeof(chunk_t)` bytes:

```

1 uint8_t *h_frame;
2 cudaMallocHost((void **)&h_frame, total * sizeof*h_frame + sizeof(chunk_t));
3 cudaMallocHost((void **)&pready->h_xs, tot*sizeof*prdy->hx+sizeof(chunk_t));
4
5 // init of Mat object with page-locked h_frame ptr
6 pready->pframe = new Mat(ctx.sampleMat->rows , ctx.sampleMat->cols ,
7     ctx.sampleMat->type() , h_frame);
8

```

The results of the *nvprof* profiler are shown in Table 7 and is noticeable the big improvement in terms of performances for what concerns the Kernel execution time. In fact, it's almost the 50% of the previous implementation. The average fps are 21, with a *for* time of 26.0 ms and a *read* time of 30 ms.

Operation	Avg (ms)	Min (ms)	Max (ms)
Kernel	23.700	21.590	24.671
DtoH	0.0079	0.0008	1.0009
HtoD	0.6483	0.6074	2.5641

Table 7: Profiling results

A further optimization can be done by looking at the new kernel implementation and noticing that only 2/3 of the memory accesses are optimized. The feedback write in fact is still done at each byte, and it can be easily optimized by modifying the `cc` register variable in the inner loop (byte wise) and save back the result in the global memory at the end of the inner loop only if needed (so if the `cc` variable has been modified).

The implementation is the following:

```

1 typedef int chunk_t;
2
3 __global__ void kernel2(
4     uint8_t *current , uint8_t *previous , uint8_t *diff ,
5     int maxSect , unsigned int *pos , int *xs) {
6         int x = threadIdx.x + blockDim.x * blockIdx.x;
7         unsigned int npos;
8         int df;
9         chunk_t cc , pc; // registers storing current copy , previous copy
10        bool currUpdateRequired = false;
11
12        if (!x) *pos = 0;
13        __syncthreads();
14
15        int start = x * maxSect;
16        int max = start + maxSect;
17

```

```

18 #pragma unroll
19     for (int i = start; i < max; i++) {
20
21         cc = ((chunk_t *)current)[i];
22         pc = ((chunk_t *)previous)[i];
23
24         #pragma unroll
25         for (int j = 0; j < sizeof cc; j++) {
26             df = ((uint8_t *)&cc)[j] - ((uint8_t *)&pc)[j];
27
28             if (df < -LR_THRESHOLDS || df > LR_THRESHOLDS) {
29                 npos = atomicInc(pos, 6220801);
30                 diff[npos] = df;
31                 xs[npos] = (i * sizeof cc) + j;
32             } else {
33                 ((uint8_t *)&cc)[j] -= df;
34                 currUpdateRequired = true;
35             }
36         }
37
38         // storing in global memory only if needed
39         if (currUpdateRequired) {
40             ((chunk_t *)current)[i] = cc;
41             currUpdateRequired = false;
42         }
43     }
44 }
45
46 }
47

```

The results of the *nuprof* profiler are shown in Table 8 and is noticeable another big improvement, reducing the kernel time by another 50%. The average fps are 24, with a *for* time of 11.5 ms and a *read* time of 40 ms. The overall comparison is shown in Figure 7.

Operation	Avg (ms)	Min (ms)	Max (ms)
Kernel	9.7944	8.7887	17.906
D <small>to</small> H	0.0080	0.0008	1.2807
H <small>to</small> D	0.6443	0.5969	2.9633

Table 8: Profiling results

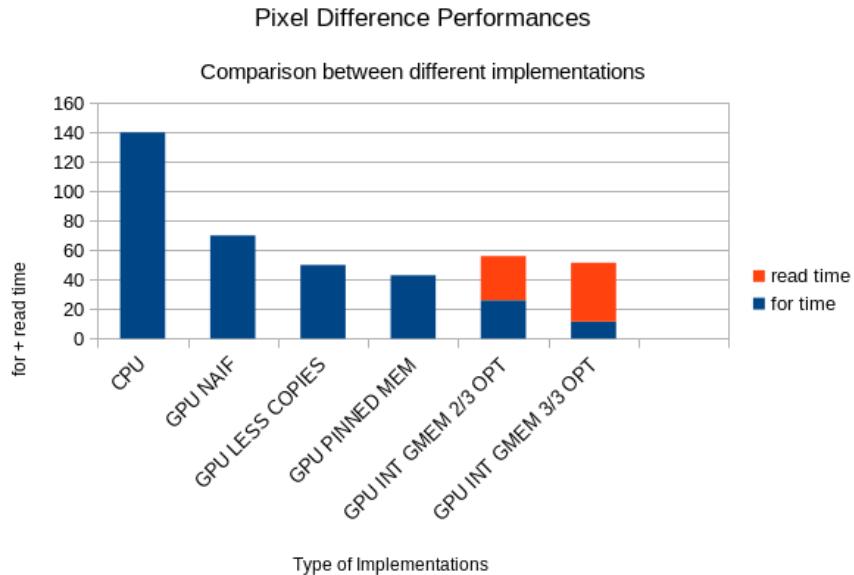


Figure 7: *Pixel Difference Comparison*

The overall situation can be optimized even more. CUDA offers the possibility to perform vectorized memory accesses by using *vectorized data types*. What are them and why they are so important?

First of all, these special data types are of different types depending on their reference type. For example there is `int2`, `float3`, `ulong2` and so on. Looking at the public header files they are defined as follow (taking `int4` as an example):

```

1 struct __device_builtin__ __builtin_align__(16) int4 {
2     int x, y, z, w;
3 };
4
5 typedef __device_builtin__ struct int4 int4;

```

Thanks to the fact they are `__device_builtin__`, they are identified by the CUDA compiler that will optimize the generated instructions by reducing them and using special *vectorized instructions* capable of a larger transaction in the memory. So thanks to these instructions is possible to load big chunks of data at a time with less instructions. Please note that a common struct with 16 int variables as members will not perform well as expected because `int4` is a builtin type so recognized by the compiler and vectorized instructions will be used by activating some internal magic in the compiler itself.

The usage of these vectorized instructions is easily seen at the assembly level. The kernel seen before with a classical `int` data type is translated into these instructions:

```

1 /*00c8*/ MOV R3, R5;
2 /*00d0*/ LD.E R18, [R10];
3 /*00d8*/ LD.E R14, [R2];
4 /*00e0*/ SSY 0x240;

```

The two LD.E instructions are the one corresponding to the two reads from the global memory into the `cc` and `pc`. The LD.E is a scalar instruction that will load a word from memory. Using instead `int4` builtin type, the disassembly code will be:

```

1 /*00c8*/ MOV R3, R17;
2 /*00d0*/ MOV R12, R0;
3 /*00d8*/ LD.E.128 R4, [R2];
4 /*00e0*/ SSY 0x360;
5 /*00e8*/ LD.E.128 R8, [R12];

```

Here, the two LD.E.128 are vectorized instructions. By using the vectorized load and store instructions LD.E.{64,128} (and the counter part ST.E.{64,128}) operations like load and store data can be done in 64- or 128-bit widths. Using vectorized loads reduces the total number of instructions, reduces latency, and improves bandwidth utilization.

At this point, the `nvprof` profiler is run and the results are shown in Table 9 and is noticeable a really big improvement in the computation, reducing the kernel time by another 60%. The average fps are 26, with a *for* time of 4.7 ms and a *read* time of 37 ms.

Operation	Avg (ms)	Min (ms)	Max (ms)
Kernel	3.4197	2.9266	9.0223
DtoH	0.0061	0.0008	1.1601
HtoD	0.6387	0.5957	2.5595

Table 9: Profiling results

Even if the 128 bit width is the maximum width support, it's possible to optimize the resulting code even more at the instruction level by using a larger builtin types like `long4`. As can be intuitive, it will be translated into 2 instructions. This helps introducing a slightly optimization by reducing the memory access frequency, achieving an average reduction of 0.4 ms in the kernel execution time.

The resulting assembly code is the following one, where it is possible to notice two loads at a time with the first one loading at offset 0x0 as the normal one and the second one loads at offset 0x10 that is 128 bit far from the first load, thus with two instructions a total of 256 bit or 32 byte are loaded, that is the size of a `long4` data type.

```

1 /*00c8*/ MOV R3, R13;
2 /*00d0*/ LD.E.128 R8, [R20];
3 /*00d8*/ LD.E.128 R4, [R20+0x10];
4 /*00e0*/ SSY 0x1e0;
5 /*00e8*/ LD.E.128 R12, [R2];
6 /*00f0*/ LD.E.128 R16, [R2+0x10];

```

The overall comparison is shown in Figure 8.

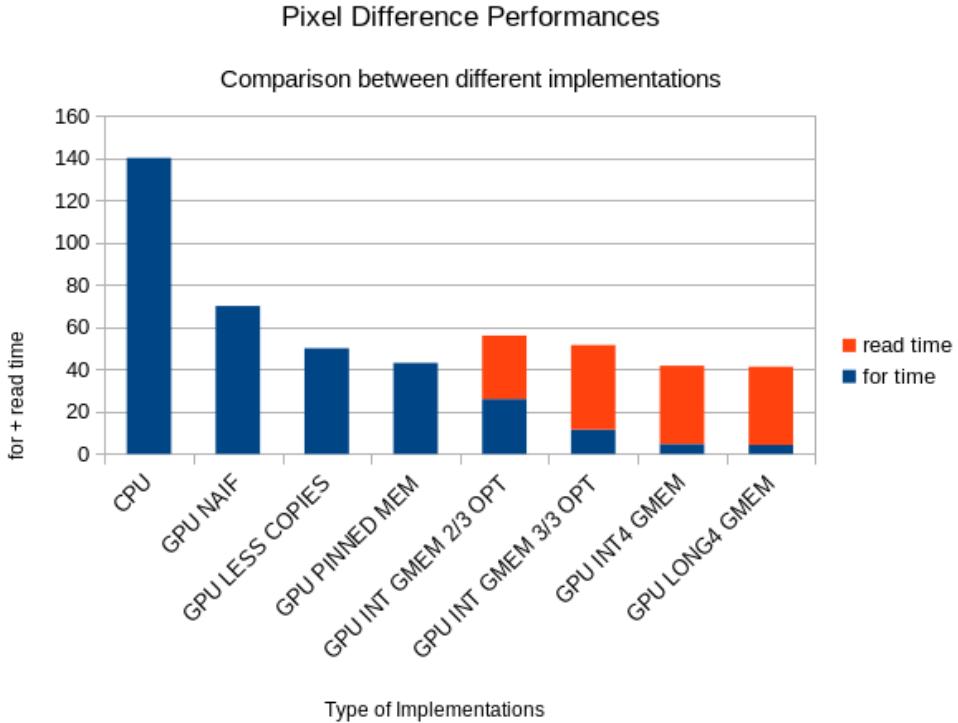


Figure 8: *Pixel Difference Comparison*

In other words, by using *vectorized instructions* the performance improvement comes from efficiency gains, in two ways:

1. At the instruction level, a multi-word vector load or store only requires a single instruction to be issued, so the bytes per instruction ratio is higher and total instruction latency for a particular memory transaction is lower.
2. At the memory controller level, a vector sized transaction request from a warp results in a larger net memory throughput per transaction, so the bytes per transaction ratio is higher. Fewer transaction requests reduces memory controller contention and can produce higher overall memory bandwidth utilization.

2.2.5 Using of shared memory for frequent accesses: is it avoidable?

Unfortunately, the algorithm and the relative CUDA kernel implementation doesn't show up any pattern of frequent access of the same data over the time, this is because the kernel works everytime on two frames and they differs continuously with respect to the previous kernel execution.

Looking at the last and very optimal kernel implementation, the only frequent global memory access where the same data is read and maybe updated is the one regarding the count of pixel

changed or `h_pos`. This means that it's reasonable to think to execute the increment operation `atomicInc()` in the shared memory and then transfer the final result in the global memory in order to have only one access in global memory per thread with respect to the pixel count variable.

A possible implementation can be the following one:

```

1  typedef long4 chunk_t;
2
3  __global__ void kernel12(
4      uint8_t *current, uint8_t *previous, uint8_t *diff,
5      int maxSect, unsigned int *pos, int *xs) {
6          int x = threadIdx.x + blockDim.x * blockIdx.x;
7          unsigned int npos;
8          int df;
9          chunk_t cc, pc; // registers storing current copy, previous copy
10         __shared__ unsigned int s_pos[1];
11
12         if (!x) *s_pos = 0;
13         __syncthreads();
14
15         int start = x * maxSect;
16         int max = start + maxSect;
17
18         #pragma unroll
19         for (int i = start; i < max; i++) {
20
21             cc = ((chunk_t *)current)[i];
22             pc = ((chunk_t *)previous)[i];
23
24             #pragma unroll
25             for (int j = 0; j < sizeof cc; j++) {
26                 df = ((uint8_t *)&cc)[j] - ((uint8_t *)&pc)[j];
27
28                 if (df < -LR_THRESHOLDS || df > LR_THRESHOLDS) {
29                     npos = atomicInc(s_pos, 6220801);
30                     diff[npos] = df;
31                     xs[npos] = (i * sizeof cc) + j;
32                 } else {
33                     current[(i * sizeof cc) + j] -= df;
34                 }
35             }
36         }
37     }
38
39     __syncthreads();
40     if (!x) *pos = *s_pos;
41
42 }
43
44 }
```

Notice the use of the function `__syncthreads()` at the end. It is used to wait for all threads to finish in order to be sure that all the `atomicInc()` instructions have been executed. After this operation, the global memory variable `pos` is updated with the value stored in the shared memory variable `s_pos` by the first thread, the one with $id = 0$.

Unfortunately the result is something unexpected: the performance are exactly the same or a bit worse than the previous implementation. This is mainly for two reasons:

1. As explained in section 2.2.4, the part of code where the `atomicInc()` function is executed is the one with a lower execution frequency because it's supposed that between a frame and another only a bunch of pixels changes so it's not worth trying to improve it that requires adding other instructions (for example the one at the end of the kernel of synchronization) that only adds overhead.
2. Atomic instructions are very particular: they rely on the bottom hardware for executing memory operations in an atomic way.

About the last point, atomic memory operations are very important in parallel programming, allowing concurrent threads to correctly perform read-modify-write operations on shared data.

About the importance of how these kind of operations are implemented in CUDA architectures, as an example the Kepler architecture can be analyzed. It significantly increased the throughput of atomic operations to global memory compared to the earlier Fermi architecture; however, both Fermi and Kepler implemented shared memory atomics using an expensive lock/update/unlock pattern. So they were implemented totally in software.

However, Maxwell architecture (the GPU on the *Jetson Nano* used for this project is based on this architecture) improved this by implementing native hardware support for shared memory atomic operations for 32-bit integers, and native shared memory 32-bit and 64-bit compare-and-swap (CAS), which can be used to implement other atomic functions with reduced overhead (compared to the Fermi and Kepler methods which were implemented in software).

So atomic operations have not been an easy thing. They are very performant for working in the global memory and *only recently* they have been optimized to work in the shared memory with near the same performance compared to the global memory.

3 Image Text Overlay

Another part of the project under discussion is the ability to add a text overlay over an image. In this particular implementation, it supports ASCII characters over a single line of the image. In order to be sent over internet, it's executed before applying the algorithm for sent the frame over a socket, according to Chapter 2.

The final result can be seen in Figure 9:



Figure 9: *Text Overlay*

The basic idea behind the implementation is the ability to generate a frame (made of RGB pixels) of an ASCII character, dynamically and based on an available font. Everything starts from the specification of the supported alphabet. In fact, not all ASCII characters are supported but only a limited set, specified during the compilation time, in order to have an efficient code in terms of memory occupied and execution time without wasting useful resources (both on the host side and on the device side) for characters that are never used.

After defining the alphabet, at the beginning of the program, it's possible to translate each character into a set of pixels composing it. This is done thanks to a series of OpenCV's function like `cv::getTextSize` that is able to retrieve the frame dimensions (width and height) of a character and `cv::putText` that is able to put a text (or a character in this case) over a frame.

First thing that is done is to retrieve the dimension of a character. This is done only for a single character and it is supposed that all the other characters in the alphabet have the same size.

Than, an array storing all the frames one after another is created with a size that is equal to $n_{chars} \cdot Area_{char}$. A loop on the alphabet will fill this array with all the character' frames. The code to implement this is the following:

```

1 #define CHARS_STR "0123456789BFPSWbkps :"
2
3 auto txtsz = cv::getTextSize("A", cv::FONT_HERSHEY_PLAIN, 3, 2, 0);
4 int fullArea = 3 * txtsz.area();
5
6 uint8_t *charsPx = new uint8_t[(sizeof(CHARS_STR) - 1) * fullArea];
7 memset(charsPx, 0x0, (sizeof(CHARS_STR) - 1) * fullArea);
8 for (int i = 0; i < (sizeof(CHARS_STR) - 1); i++) {
9     Mat pxBaseMat(txtsz.height, txtsz.width, base.type(), charsPx+i*fullArea);
10    cv::putText(pxBaseMat, std::string(CHARS_STR).substr(i, 1),
11                cv::Point(0, txtsz.height+1), cv::FONT_HERSHEY_PLAIN, 3,
12                cv::Scalar(0, 255, 0), 2, cv::LINE_AA);
13 }
```

When the array is finally full filled, it's ready to be copied onto the GPU memory. This operation is done only at the beginning, so it's ininfluent from the point of view of the performances:

```

1 uint8_t *d_charsPx;
2 int totcpy = fullArea * sizeof *d_charsPx * (sizeof(CHARS_STR) - 1);
3 cudaMalloc((void **)&d_charsPx, totcpy * sizeof *d_charsPx);
4 cudaMemcpy(d_charsPx, charsPx, totcpy * sizeof *d_charsPx, cudaMemcpyHostToDevice);
```

3.1 GPU Implementation

After the initialization phase, is now time to give a look at how each character frame is placed onto the image. It's a replacement of the pixels of the frame under analysis with the one of the corresponding character.

The first thing to do is to decide how many kernel to launch and how many pixel each kernel needs to manage. By deciding to assign to each kernel a certain number of pixel in a row-major fashion, and by considering the fact that the frame size of each character is fixed but can change if the font is changed, a portion of the code is used to decide in a dynamic way the two values: number of kernels and pixels for each kernel. It's decided by dividing the total number of bytes of a character frame by the number of frame and if it's a multiple of the size of `chunk_t` type that for a first analysis will be a `uint8_t`, then the number of threads found is the right one. The code is the following:

```

1 typedef uint8_t chunk_t;
2
3 int nThreadsToUse = nMaxThreads;
4 int eachThreadDoes = 1;
5 for (int i = nMaxThreads; i > 0; i--) {
6     float frac = fullArea / (i * 1.0);
7     int frac_int = fullArea / i;
8     if ((int)ceil(frac) == frac_int && frac_int % sizeof(chunk_t) == 0) {
9         nThreadsToUse = i;
10        eachThreadDoes = frac_int / sizeof(chunk_t);
11        break;
12    }
13 }
```

3.1.1 Naif version

After the setup time, its the turn of the computation step. This is the core, where each character of the string to print is analyzed before applying the streaming algorithm. This is done on the host side and not on the device side in order to avoid the copy HtoD of the string.

The advantage is given by the usage of async operations, so while the frame is copied on the GPU, the host scans each character and a new kernel is fired in an async way in order to print each character one after the another.

```

1 // Copy in the current pointer and run
2 cudaMemcpyAsync(d_current , pready->pframe->data , total , cudaMemcpyHostToDev );
3
4 for ( int offset = 0 , j = 0; j < text.length (); j++ , offset+=txtsz.width*3) {
5     int idx ;
6     for ( int i = 0; i < ( sizeof(CHARS_STR) - 1); i++) {
7         if (CHARS_STR[i] == text.at(j)) {
8             idx = i;
9             break;
10        }
11    }
12
13    kernel_char<<<1, nThreadsToUse>>>(d_current , d_charsPx + idx * fullArea ,
14        eachThreadDoes , offset , 3 * txtsz.width , 3 * pready->pframe->cols );
15 }
```

The basic implementation of the kernel is the following:

```

1 __global__ void kernel_char(
2 uint8_t *current , uint8_t *matrix , int N, int offset ,
3 int matrixWidth , int currWidth) {
4     int thid = threadIdx.x + blockIdx.x * blockDim.x;
5
6     int start = thid * N;
7     int max = start + N;
8
9     for ( int i = start; i < max; i++) {
10         int x = offset + i % matrixWidth;
11         int y = i / matrixWidth;
12         current[y * currWidth + x] = matrix[i];
13     }
14 }
```

By running the profiler for a fixed time of 30 seconds, the kernel takes an average time of 2.981 us so it occupies the 1.53% of the total time.

A comparison chart is show in Figure 10:

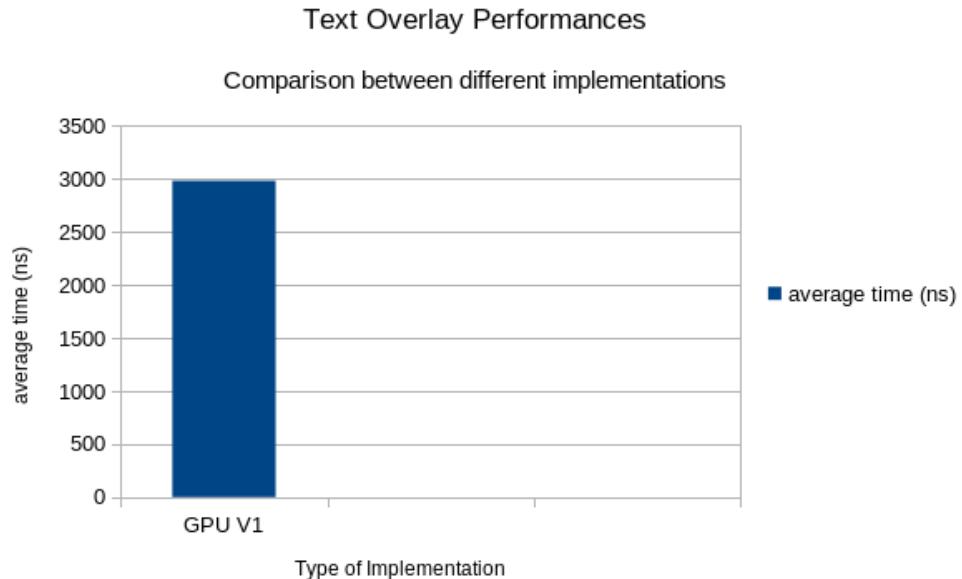


Figure 10: *Text Overlay Comparison*

3.1.2 Exploiting the global memory access optimization

It's possible then to optimize the kernel from the point of view of the global memory access parallelism as explained in Section 2.2.4. By copying and modifying the same kernel structure, the resulting kernel is:

```

1  typedef long4 chunk_t;
2
3  __global__ void kernel2_char(
4      uint8_t *current, uint8_t *matrix, int N, int offset,
5      int matrixWidth, int currWidth) {
6          int thid = threadIdx.x + blockIdx.x * blockDim.x;
7          chunk_t cc;
8
9          int start = thid * N;
10         int max = start + N;
11
12         for (int i = start; i < max; i++) {
13             int reali = i * sizeof(chunk_t);
14             int x = offset + reali % matrixWidth;
15             int y = reali / matrixWidth;
16
17             int idx = (y * currWidth + x) / sizeof(chunk_t);
18             cc = ((chunk_t *)current)[idx];
19
20             #pragma unroll
21             for (int j = 0; j < sizeof cc; j++) {
22                 ((uint8_t *)&cc)[j] = matrix[reali + j];
23             }
24
25             ((chunk_t *)current)[idx] = cc;
26         }
27     }
```

By running the profiler for a fixed time of 30 seconds, the kernel takes an average time of 3.778 us so it occupies the 1.91% of the total time.

A comparison chart is show in Figure 11:

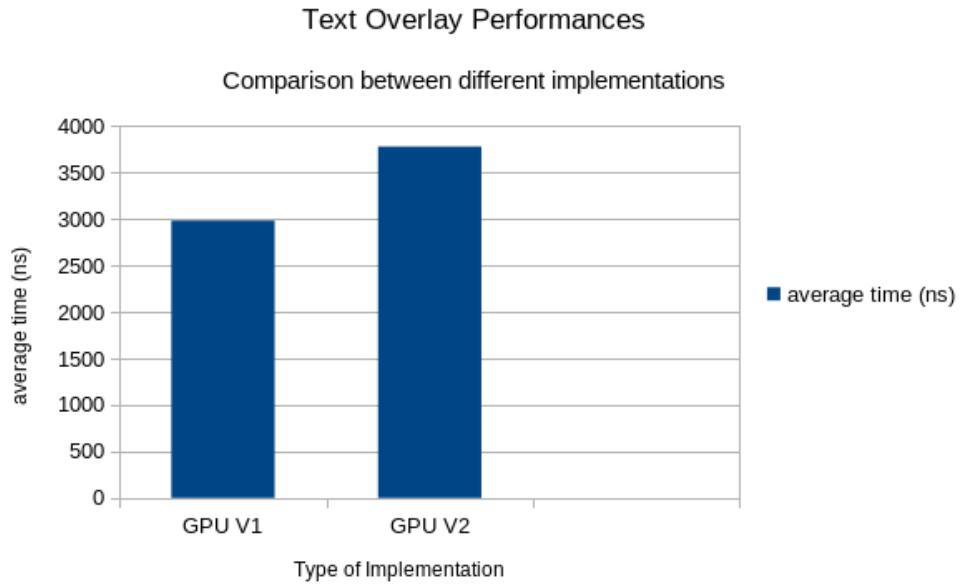


Figure 11: *Text Overlay Comparison*

As noticeable, the performance are worse. This is because the intermediate *for loop* can be compacted in a single instruction because it's only an assignment without other operations in the middle. The new kernel is then shorter and faster, as follow:

```

1 typedef long4 chunk_t;
2
3 __global__ void kernel2_char(
4     uint8_t *current, uint8_t *matrix, int N, int offset,
5     int matrixWidth, int currWidth) {
6     int thid = threadIdx.x + blockIdx.x * blockDim.x;
7     chunk_t cc;
8
9     int start = thid * N;
10    int max = start + N;
11
12    for (int i = start; i < max; i++) {
13        int reali = i * sizeof(chunk_t);
14        int x = offset + reali % matrixWidth;
15        int y = reali / matrixWidth;
16
17        int idx = (y * currWidth + x) / sizeof(chunk_t);
18        ((chunk_t *)current)[idx] = ((chunk_t *)matrix)[i];
19    }
20

```

By running the profiler for a fixed time of 30 seconds, the kernel takes an average time of 1.868 us so it occupies the 1.01% of the total time. This means that the last version of the kernel is faster than the first one.

A comparison chart is show in Figure 12:

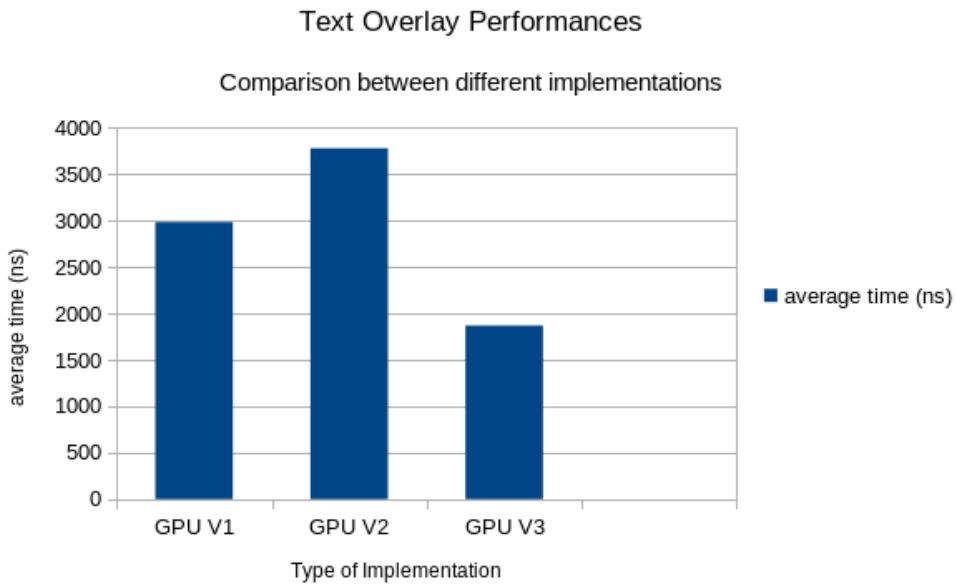


Figure 12: *Text Overlay Comparison*

3.1.3 Constant Memory instead of Global Memory: is it worth?

Because character frames are basically constant during the execution of the program, is it really necessary to use the global memory or another kind of memory can be exploited? If yes, is it really worth?

For this kind of purposes, CUDA devices offer the so called Constant Memory, that is a piece of memory on the GPU side that is constant. Because it's *constant* from the GPU point of view but it can modified by the Host, it seems to be perfect for this part of the project.

Usually it's large *64 KB* so supposing to have an alphabet whose frames doesn't exceed this memory limit, a constant variable (array) can be declared and initialized as follows:

```

1 __constant__ uint8_t c_matrix[( sizeof(CHARS_STR) - 1 ) * 3 * 32 * 28];
2
3 uint8_t *c_matrixPtr;
4 cudaMemcpyToSymbol(c_matrix, charsPx, totcpy);
5 cudaGetSymbolAddress(( void **)&c_matrixPtr, c_matrix);

```

Leaving `kernel2_chars` as it is and using `c_matrixPtr` instead of `d_charsPx` in the kernel call and running the profiler, an unexpected value is returned by the profiler. The execution time is almost the same (faster only for a few nanoseconds) than the previous implementation and this is due to two reasons mainly:

1. the kernel is instruction-bound and not memory-bound
2. according to CUDA's programming guide: "*For all threads of a half-warp, reading from the constant cache is as fast as reading from a register as long as all threads read the same address. The cost scales linearly with the number of different addresses read by all threads.*". This means that the shared memory is really faster than the global memory but because the kernel is accessing every time a different address, the cost of access increases by reaching the cost of the global memory.

A comparison chart is show in Figure 13:

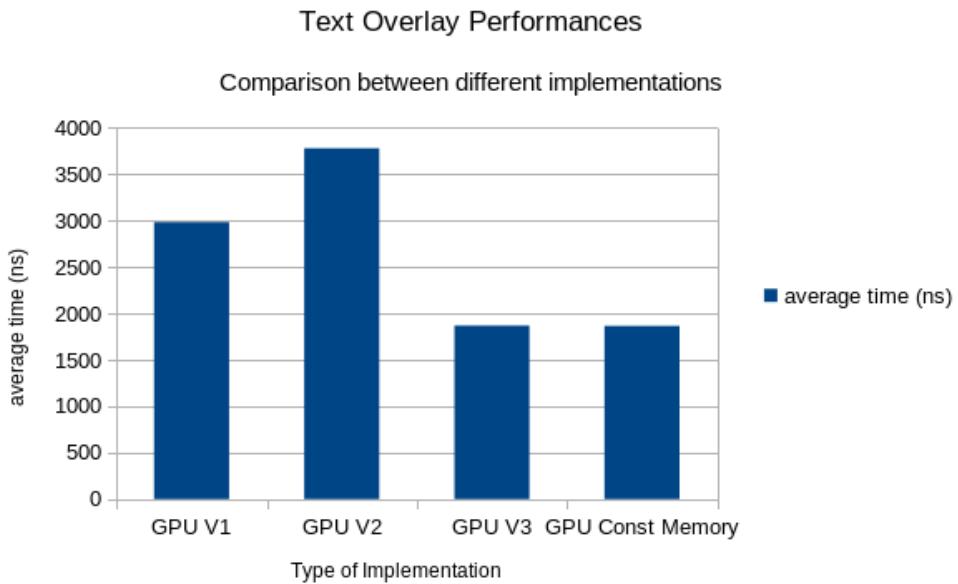


Figure 13: *Text Overlay Comparison*

4 HeatMap

Since the basis of this project is to send the pixel difference, the lower the dissimilarities the higher the bandwidth saving. In order to better visualize which pixels changes the most and the different magnitudes, an heatmap is generated with the current and the previous frames.

A heat map is nothing else than a data visualization technique that represents the magnitude of a measurement as colors in two dimensions, in this case an image. The idea is first compute the difference at the pixel level between two successive frames and then represent it using a color scale from blue to red, where blue means low difference and red high difference. The Image 14 shows a naif implementation of heat map, available at `heat_map_benchmark/v0.cu`, that generates in real time the image on the right.



Figure 14: Example of real-time heat map from webcam

The basic steps for generating an heat map are the following:

- Take two frames from the webcam via OpenCV
- Compute the pixel difference
- Translate the difference into the corresponding color for the heat map
- Copy the result into a support buffer image and visualize it

All programs versions for the heat map are available in the `heat_map_benchmark` folder.

4.1 Heat map pixel mapping - CPU Naif implementation

The first *naif* version has been done on the CPU. This implementation can generate an heatmap in more or less 980ms, that is too much. This is also due to the complexity of the function itself that is used to map a normalized pixel difference to a blue-red scale and because the image is sequentially analyzed. In order to convert the difference of a pixel into the three color component (RGB), the usage of the *sine* function has been done. The difference of the each pixel has been computed as the sum of the absolute value of the difference of the single color component, as:

$$diff = abs(Previous[i, R] - Curr[i, R]) + abs(Previous[i, G] - Curr[i, G]) + abs(Previous[i, B] - Curr[i, B])$$

Where $Previous[i, R]$ is the pixel red color component of the pixel at index i . In the worst case, a pixel can be turned from black to white or vice versa and so, the $0 \leq diff \leq 765$, that is $255 \cdot 3$. Then the $diff$ value is taken and normalize, by using 765 so that, $0 \leq diff_{norm} \leq 1$:

$$diff_{norm} = \frac{diff}{765}$$

This value must be mapped into the tree RGB component of the heat map; the pixel must be more blue if the difference is more toward 0.0, yellow/green if near 0.5 and red if next to 1.0. The smoothed way to perform this task is to use three different sine functions, centered respectively on 0.0, 0.5 and 1.0 as in the following way:

$$\begin{aligned} \text{RED} & \quad \sin(\pi \cdot diff_{norm} - \frac{\pi}{2.0}) \\ \text{GREEN} & \quad \sin(\pi \cdot diff_{norm}) \\ \text{BLUE} & \quad \sin(\pi \cdot diff_{norm} + \frac{\pi}{2.0}) \end{aligned}$$

The plot of the three sine functions is available at Figure 15.

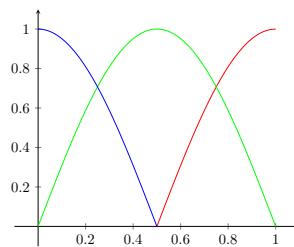


Figure 15: Mapping function from pixel difference to RGB components

The CPU implementation is the following:

```

1 struct HeatElemt {
2     int r;
3     int g;
4     int b;
5 };
6
7 HeatElemt getHeatPixel(int diff){
8     struct HeatElemt h;
9     float diff1 = diff/(255.0*2.0);
10
11    // Map the difference into the three color components
12    h.r = min(max(sin(M_PI*diff1 - M_PI/2.0)*255.0, 0.0),255.0);
13    h.g = min(max(sin(M_PI*diff1)*255.0, 0.0),255.0);
14    h.b = min(max(sin(M_PI*diff1 + M_PI/2.0)*255.0, 0.0),255.0);
15 }
16
17 while (1) {
18     cap >> image2;
19     for (int y = 0; y < H; y++){
20         for (int x = 0; x < W; x++){
21             Vec3b & intensity = image1.at<Vec3b>(y, x);
22             Vec3b a = image1.at<Vec3b>(y, x);
23             Vec3b b = image2.at<Vec3b>(y, x);
24
25             // Compute the absolute difference
26             HeatElemt elem = getHeatPixel(abs(a.val[0] - b.val[0]) +
27                 abs(a.val[1] - b.val[1]) + abs(a.val[2] - b.val[2]));
28
29             intensity.val[0] = elem.b;
30             intensity.val[1] = elem.g;
31             intensity.val[2] = elem.r;
32         }
33     }
34     image1 = image2.clone();
35 }

```

In this CPU implementation, a loop is performed over each pixel of the two frames and the absolute difference computed. Another function, called `getHeatPixel` is used to convert that value to the heat map RGB color space. At this point, the original image is overwritten with the heatmap colors.

4.2 CUDA implementation

The idea is to rewrite what described in the previous section into a CUDA code. The GPU allows to run in parallel multiple instance of the same kernel, so that the execution can be done in parallel in order to speed up the computation.

From the perspective of the interaction, the GPU acts like an accelerator of the CPU, that ask it to execute the kernel by resorting to the following phases:

1. Copy the frames from the Host to the Device memory
2. Execute the kernel and wait its completion
3. Retrieve the result from the Device to the Host memory

This is exactly what the next section will explain. The kernel call allows to configure how many thread per kernel will be executed; obviously, depending on that number, the accessed locations must be defined accordingly. In fact, by defining K the number of thread lunched, each thread will work on a specific portion of the entire image:

$$\text{Thread portion dimension} = \frac{W \cdot H \cdot 3}{K}$$

Where W and H are the width and height of the image. The multiplication by 3 depends on the fact that each pixels is defined by three `uint8_t` datatype that each one of them defines a specific color.

Furthermore, from the memory allocation perspective of the GPU, the memory for the two frames and the heat map must be allocated. This is done only once at the startup of the program, thanks to the CUDA API.

```

1 // Pointer definition
2 uint8_t *d_current, *d_previous;
3 uint8_t *d_heat_pixels;
4
5 // Memory space reservation on GPU
6 cudaMalloc((void **)&d_current, W*H*C * sizeof *d_current);
7 cudaMalloc((void **)&d_previous, W*H*C * sizeof *d_previous);
8 cudaMalloc((void **)&d_heat_pixels, W*H*C * sizeof *d_heat_pixels);

```

Unfortunately, this type of the problem doesn't need any shared or constant memory because each location in the image (all three colors of all pixels) are accessed only once and using a shared memory would have only reduced the performance due to the overhead of the useless copy. Moreover, there are no data that are constant.

CUDA allows to get the information about the maximum number of thread by using the `cudaGetDeviceProperties` command and, by referring to the equation defined before, K can be set to that value. In fact, the kernel is launched with the following configuration:

- *Grid dimensions:* 1, 0, 0
- *Block dimensions:* K , 0, 0

Having the maximum number of threads in the `threads` variable, it's possible to call the kernel in the following way:

```

1 // Kernel per block computation
2 cudaGetDeviceProperties(&prop, 0);
3 int threads = prop.maxThreadsPerBlock;
4 int maxSection = (W*H*C)/threads;

```

An extensive analysis of the correct number of threads has been done at Section 4.3.

4.2.1 CUDA Naif implementation

The first implementation of the algorithm in CUDA is based on a 1:1 transposition of what done in the CPU, in the GPU. Always using OpenCV, two next frames are fetched, send to the kernel and the heat map computed.

```

1 // Copy from Host to Device
2 cudaMemcpy(d_prev, image1, W*H*C * sizeof *image1, cudaMemcpyHostToDevice);
3 cudaMemcpy(d_curr, image2, W*H*C * sizeof *image2, cudaMemcpyHostToDevice);
4
5 kernel<<<1, threads>>>(d_curr, d_prev, (W*H*C)/threads, d_out);
6
7 // Copy heat map from Device to Host
8 cudaMemcpy(heatmap, d_out, W*H*C * sizeof *heatmap, cudaMemcpyDeviceToHost);

```

This naif implementation implies two memory transfers (HostToDevice) for the previous and current frames and one DeviceToHost for the generated heat map.

In order to speed up the data management, instead of copying into a support array of `uint8_t` the entire two frames (previous and current), both frames are directly copied into the device buffers with the `cudaMemcpy` procedure.

From the kernel perspective, there is a non-negligible complication with the respect to the CPU implementation. The CPU code is based on a single thread that iterates over the entire image in a sequential way, accessing one location after the other. For the GPU this is not the case, since now, each thread will work in parallel on a portion of the image that is long `maxSect`.

Due to this, each kernel must know exactly from which pixel start to retrieve the data and where to store the results. This is only a matter of index management; let's suppose to have frame with dimension $1920*1080*3$. So, by using 1024 threads per block, each thread will work on:

$$maxSect = \frac{1920 * 1080 * 3}{1024} = 6075$$

This means that the first thread must work from 0 to 6074, the second one from 6075 to 12149 and so on. This can be simply achieved by giving an univocal index identifier to each kernel, that can be generated as:

```
int x = threadIdx.x + blockDim.x * blockIdx.x;
```

The GPU implementation of the first CUDA kernel is the following:

```

1  __global__ void kernel(uint8_t *current, uint8_t *previous,
2                         int maxSect, uint8_t* d_heat_pixels) {
3
4     // Index relative to the block
5     int x = threadIdx.x + blockDim.x * blockIdx.x;
6
7     // Start of the sector for this thread
8     int start = x * maxSect;
9     int max = start + maxSect;
10    for (int i = start; i < max; i=i+C) {
11
12        // Compute the pixel difference
13        int pixelDiff = fabsf(current[i] - previous[i]) + fabsf(current[i+1]
14                           - previous[i+1]) + fabsf(current[i+2] - previous[i+2]);
15        float diff1 = pixelDiff/(255*3.0);
16
17        // Map different into the three color component
18        int r = fminf(fmaxf(sinf(M_PI*diff1 - M_PI/2.0)*255.0, 0.0),255.0);
19        int g = fminf(fmaxf(sinf(M_PI*diff1)*255.0, 0.0),255.0);
20        int b = fminf(fmaxf(sinf(M_PI*diff1 + M_PI/2.0)*255.0, 0.0),255.0);
21        d_heat_pixels[i] = b;
22        d_heat_pixels[i+1] = g;
23        d_heat_pixels[i+2] = r;
24    }
25 }
```

After having run the `nvprof`, the main contributions to the execution time from the profiler are:

Type	Time (%)	Avg	Name
GPU activities	86.14	49.958	kernel
	9.38	2.5577ms	[CUDA memcpy HtoD]
	4.48	2.4427ms	[CUDA memcpy DtoH]
API calls	93.88	18.820ms	cudaMemcpy
	5.88	181.61ms	cudaMalloc

Table 10: *Profiling result v1.cu*

In fact, the cumulative time needed to copy all two frames into the device, execute the kernel and copy back the image in order to display it, takes approximately 57ms. From the GPU perspective, the average time to execute one single kernel is 49.958ms.

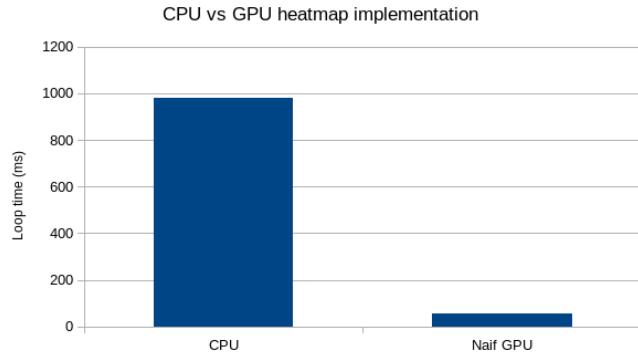


Figure 16: *Loop time CPU vs GPU*

4.2.2 CUDA switching frames

The first idea to reduce the time used for copying the frames from the Host to the Device. Instead of copying each time both the two frame, one of the two can be reused by only switching the two points, so that the one that previously was the current will became the previous; at this point, only the new frame must be copied into the device memory. This means that approximately the *CUDA memcpy HtoD* should be half of the previous time. This led to a *v2.cu* implementation, that exploit this pointer switching to reduce the memory transfer. The below code portion shows only the pointer switching in the loop used to fetch the frames.

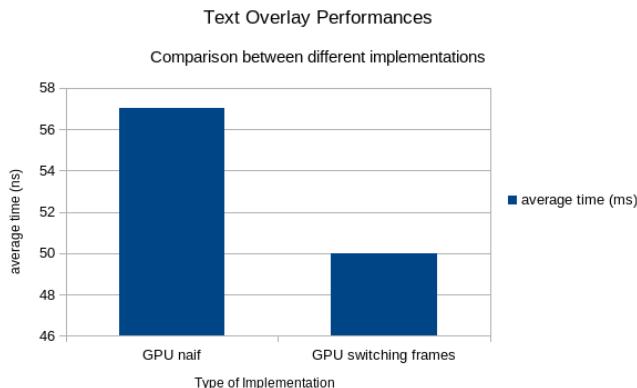
```

1 while(1){
2     // New frame fetch
3     cap >> image2;
4
5     // Pointer switching
6     uint8_t* tmp = d_curr;
7     d_curr = d_prev;
8     d_prev = tmp;
9
10    // Kernel call
11    cudaMemcpy(d_curr, image2, W*H*C * sizeof *image2, cudaMemcpyHostToDevice);
12    kernel<<<1, threads>>>(d_curr, d_prev, (W*H*C)/threads, d_out);
13    cudaMemcpy(heatmap, d_out, W*H*C * sizeof *heatmap, cudaMemcpyDeviceToHost);
14
15    image1 = image2.clone();
16 }
```

Type	Time (%)	Avg	Name
GPU activities	90.54	46.023	kernel
	4.85	2.4400ms	[CUDA memcpy HtoD]
	4.61	2.3453ms	[CUDA memcpy DtoH]
API calls	89.88	26.026ms	cudaMemcpy
	9.90	192.15ms	cudaMalloc

Table 11: *Profiling result v2.cu*

The average kernel execution time is in ms, and it is not a figure of interest, since it is more or less equal to the previous version. What is important is that now, the percentage of time used for the kernel is increase, due to the reduced *CUDA memcpy HtoD* (from 9.38% to 4.85%). This means that the GPU will analyze more frame in the same time frame. As we would expect, now the time needed to copy a frame from the memory to the device, execute the kernel and then retrieve the heat map takes about 50ms (12% faster).



4.2.3 CUDA Global memory access granularity

The problem with the generation of the heat map is that, the computation of the color must be done every time for all the pixels in order to compute the complete image.

Since each thread must perform 6075 iterations and need to write on the Global memory the same amount of times. In order to reduce the number of accesses to the Global memory, the idea was to access at the `int` level instead of the `byte` level. This means that frame information are still copied from host to device as arrays of `bytes` but they are accessed at the `int` level. So, if the current and the previous frames are passed as `uint8_t *current`, `uint8_t *previous`, the access is aligned at the 4 bytes. Another problem arises: the threads now access the memory with a granularity of 4 byte, but since the pixel difference needs only the first three bytes (RGB), in order to optimize and avoid to read twice from the memory, colors are updated only once every 3 bytes. So, when the position of the color in the image is the first, the colors of the heat map are set in the current and next two bytes of the output array, accordingly to the computed difference.

```

1  __global__ void kernel(uint8_t *current, uint8_t *previous,
2                         int maxSect, uint8_t *d_heat_pixels) {
3     int x = threadIdx.x + blockDim.x * blockIdx.x;
4     int start = x * maxSect;
5     int max = start + maxSect;
6     int cc, pc;
7     for (int i = start; i < max; i++) {
8
9        // Access one 4 byte at a time
10       cc = ((int *)current)[i];
11       pc = ((int *)previous)[i];
12       int pixelDiff = 0;
13       for (int j = 0; j < 4; j++) {
14
15          // Conversion from difference to heat map only every 3 bytes
16          if((i*4+j) % 3 == 0){
17              int pixelDiff = fabsf(((uint8_t *)&cc)[j] - ((uint8_t *)&pc)[j]) +
18                  fabsf(((uint8_t *)&cc)[j+1] - ((uint8_t *)&pc)[j+1]) +
19                  fabsf(((uint8_t *)&cc)[j+2] - ((uint8_t *)&pc)[j+2]);
20              float diff1 = pixelDiff/(255*2.0);
21              int r = fminf(fmaxf(sinf(M_PI*diff1 - M_PI/2.0)*255.0, 0.0),255.0);
22              int g = fminf(fmaxf(sinf(M_PI*diff1)*255.0, 0.0),255.0);
23              int b = fminf(fmaxf(sinf(M_PI*diff1 + M_PI/2.0)*255.0, 0.0),255.0);
24              d_heat_pixels[i*4+j] = b;
25              d_heat_pixels[i*4+j+1] = g;
26              d_heat_pixels[i*4+j+2] = r;
27
28             // Reset the pixel difference
29             pixelDiff = 0;
30         }
31     }
32   }
33 }
```

Since now each threads works on 4 bytes at a iteration, the dimension of the data section that each block must work on is reduced by 1/4, as in the following way:

```

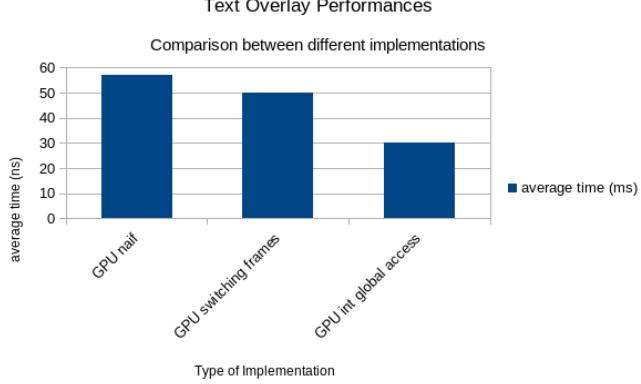
1  cudaGetDeviceProperties(&prop, 0);
2  int threads = prop.maxThreadsPerBlock;
3
4  cudaMemcpy(d_curr, image2, W*H*C * sizeof *image2, cudaMemcpyHostToDevice);
5
6  // Kernel call /4
7  kernel<<<1, threads>>>(d_curr, d_prev, ((W*H*C)/threads)/4, d_out);
8  cudaMemcpy(heatmap, d_out, W*H*C * sizeof *heatmap, cudaMemcpyDeviceToHost);
```

This led to another version, available at `v3.cu` allows to obtain the following results:

Type	Time (%)	Avg	Name
GPU activities	84.08	25.457ms	kernel
	8.14	2.4405ms	[CUDA memcpy HtoD]
	7.78	2.3549ms	[CUDA memcpy DtoH]
API calls	85.73	15.810ms	cudaMemcpy
	13.93	172.12ms	cudaMalloc

Table 12: Profiling result `v3.cu` using `int` access level

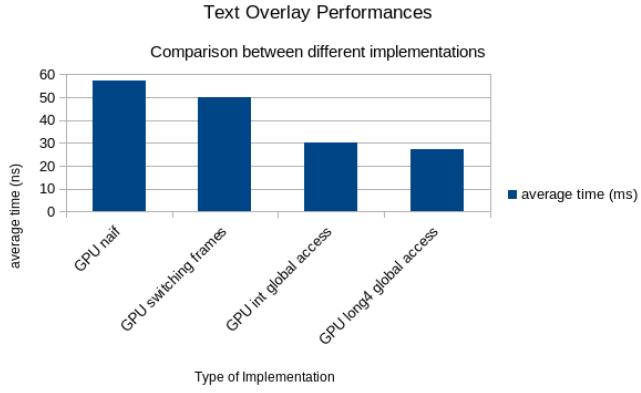
This version allows to copy the next frame from memory to device, execute the kernel and retrieve the result in about 30ms (about 40% of performance increase from `v2.cu`). This is highlighted in the table by the average time needed to execute the kernel itself, we went from 46.023ms to 25.457ms, thanks to the reduce access time to the memory.



This can be pushed even more, so that instead of accessing the images that are in the global memory, at the `int` level, they are accessed via a vectorized access. The complete explanation is available at Section 2.2.4. By the way, it is able to increase even more the performance, allowing to complete the for loop in around 27ms using `long4` data type.

Type	Time (%)	Avg	Name
GPU activities	82.83	23.589ms	kernel
	8.92	2.5160ms	[CUDA memcpy HtoD]
	8.24	2.3478ms	[CUDA memcpy DtoH]
API calls	85.42	14.903ms	cudaMemcpy
	14.18	165.80ms	cudaMalloc

Table 13: Profiling result `v3.cu` using `long4` access level



4.3 Evaluation of the number of threads

For a first implementation, the number of thread has been set to the maximum allowable from the architecture, that is given by the `cudaGetDeviceProperties` CUDA function in order to make the program independent form the device used. For example, for the Jetson Nano, the maximum number of threads are 1024.

In order to understand how the number of threads impacts on the heat map generation, a bash script has been build in order to dynamically change the K parameter via compiler directive.

The number of threads must be a multiple of 4, so that the array of the pixels can be divided into portion in such a way that a pixel is not split between two kernel.

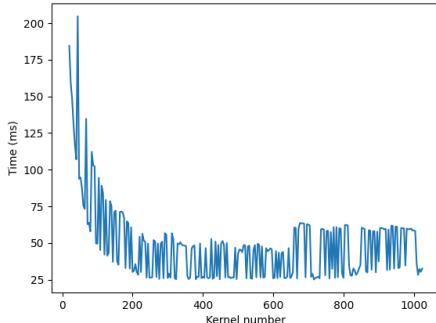
This is the bash code used to extract the *nvprof* information:

```

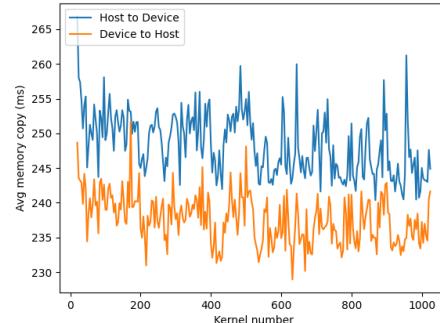
1 #!/bin/bash
2 for i in {1..256}
3 do
4   k=$(( 4*i ))
5
6   # Profiler call
7   avg='sudo /usr/local/cuda/bin/nvprof ./heatMap $k 2>&1'
8
9   # Data extraction
10  kern='echo "$avg" | grep kernel | awk '{print $6}''
11  kernt='echo "$avg" | grep kernel | awk '{print $4}''
12  hd1='echo "$avg" | grep "CUDA memcpy HtoD" | awk '{print $4}''
13  hd2='echo "$avg" | grep "CUDA memcpy HtoD" | awk '{print $2}''
14  dh1='echo "$avg" | grep "CUDA memcpy DtoH" | awk '{print $4}''
15  dh2='echo "$avg" | grep "CUDA memcpy DtoH" | awk '{print $2}''
16  echo "$k $all $kern $kernt $hd1 $hd2 $dh1 $dh2" >> times.txt
17 done

```

In this way, the output of the profiler and the time needed to copy the frame, generate the heat map and retrieve the result is parsed and wrote into a file called `times.txt`. Thanks to another script, the most useful data are plot, as below:



(a) Time needed to perform an heatmap depending on the number of kernel set



(b) Time needed to copy from Host to Device (blue) and from Device to Host (orange)

This is not the behaviour that we would have expected, beside the time needed to copy is more or less constant, by increasing the number of threads for that kernel we would expect that the time needed for the heat map computation would be lower. This is probably due to the fact that the *warp* has a fixed size of 32 threads and, even if by increasing the number of threads its execution time is lower, their management probably introduce too much overhead to obtain benefits.

Beside this, the time needed to perform an heatmap vs the number of threads, shows a peculiar behaviour. After about $N > 280$ the times tend to oscillate between more or less 50ms and 27ms. Even if this seems not a big difference, in the field of real-time image processing, it's a huge improvement. In order to avoid errors, the same script has been run multiple times and the results is always the same. Since the internal infrastructure is a black box, the hardware probably manage in different ways the threads with the respect to their number.

The best observed thread configuration for the Jetson Nano, seems to be 716.

In fact, by running the same exact algorithm describe before but the number of threads is set to 718 (that is the best kernel accordingly to the plot above), the time needed to copy a frame, compute the heat map and then copy back the heat map matrix is about 27ms. This leads to a increasing of performance of 10%, that in this domain is not negligible.

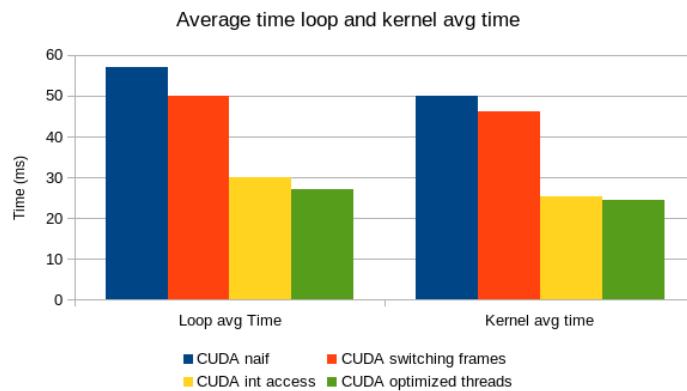
The following table shows some meaningful data extracted from the profiler:

Type	Time (%)	Avg	Name
GPU activities	84.30	24.380ms	kernel
	8.04	2.4599ms	[CUDA memcpy HtoD]
	7.66	2.383ms	[CUDA memcpy DtoH]
API calls	90.06	16.192ms	cudaMemcpy
	9.57	115.27ms	cudaMalloc

Table 14: Profiling result v3.cu

4.3.1 Heatmap Conclusions

By considering only the CUDA algorithm, the optimizations shows a decreasing fashion of the time needed to process the two frames and generate the heat map. The Loop average time, is the time needed to copy the image from the Host to the Device, summed to the heat map time computation and the time needed to copy it back. On the other hand, the plot on the right shows only the average kernel time. All data are in ms.



5 Noise visualizer

As already explained in the first sections, the algorithm is based on sending only the difference of all pixels whose color components are above a certain threshold.

In order to better visualize the noise in each frame taken from the webcam, all colors of all pixels that are above a certain threshold, are colored in red. The threshold is the same used by the kernel that computes all pixel differences. Obviously, the complexity of this computation is much easier than the previous one and the CPU allows to compute that kind of map in around 250ms.

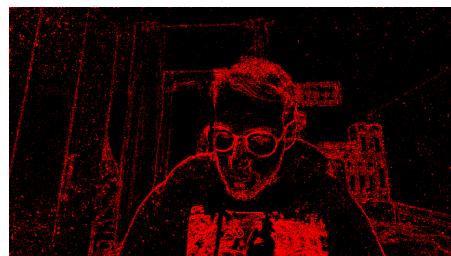


Figure 18: All non normalized pixel difference are turned to red if above a certain threshold, in this case 20. Beside the figure in the center, all red pixels can be identified as noise.

The purpose of this new color mapping is completely different from the heatmap, since the latter one is used to understand the magnitude of the pixel difference while the second one (the black-red) is used to better visualize the noise pixels, by using a threshold.

Why is this necessary? This allows to have a visual evaluation of the noise filter, that is used to reduce smooth the noise in order to reduce the difference between two frames and reduce even

more the bandwidth. The noise filter will be explained in the next section.

First of all, we need to define what is the noise: by means of noise we are referring to all the random variations on the color information in an image, that is visible as grain in film and pixel level. The level of the noise mainly depends:

- Length of the exposure
- Sensitivity of the optical camera
- Brightness of the environment
- Temperature

In most of the cases, the webcam used at home is a cheap one that has an extremely small pixels in the camera sensor. There are mainly three types of noises: fixed pattern noise, random noise and banding noise. The common noise in our webcam is the one shown below. As you can see, the dark and less brightness parts of the image are characterized by some particles that shows different color from the whole context.



Figure 19: *Example of noise*

From our perspective this is a very bad behaviour due to the fact that the noise is random, and can change the value of the pixel at each frame. This leads to an high amount of pixel changed, even if the image is always the same. For this reason, *noise image filtering* is performed over the frame, before being analyzed by the CUDA kernel that computes the pixel differences.

The internal logic of a webcam allows to trim some parameters of the video capture directly by OpenCV. This is not a software level manipulation of the frame, but the internal control unit of the webcam, while sampling the frame, performs different behaviour depending on the selected configuration.

For this reason, the original image has been enhanced with a control panel, so that the user can change the brightness, the contrast, the saturation and the gain of the image. The Figure 20, shows the panel that is shown to the user. Even if it can be used to reduce the noise, by for example reducing the contrast, it allows the user to increase the quality of the send image.

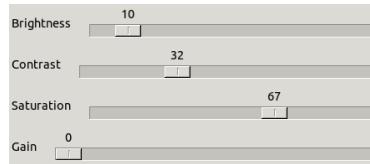


Figure 20: *Image Control panel*

5.0.1 Naif CPU implementation

The *naif* CPU implementation starts from the same program base used for the heat map. The entire image, pixel by pixel, is analyzed and if one of the color component is above the same threshold set for the pixel difference, that pixel is colored in red otherwise in black.

The code below is the C implementation of what explained before, where `LR_THRESHOLDS` is set to 20.

```

1  for (int y = 0; y < H; y++){
2    for (int x = 0; x < W; x++){
3      Vec3b & intensity = image1.at<Vec3b>(y, x);
4      Vec3b a = image1.at<Vec3b>(y, x);
5      Vec3b b = image2.at<Vec3b>(y, x);
6
7      if (abs(a.val[0] - b.val[0]) > LR_THRESHOLDS ||
8          abs(a.val[1] - b.val[1]) > LR_THRESHOLDS ||
9          abs(a.val[2] - b.val[2]) > LR_THRESHOLDS)
10     {
11       intensity.val[0] = 0;
12       intensity.val[1] = 0;
13       intensity.val[2] = 255;
14     } else {
15       intensity.val[0] = 0;
16       intensity.val[1] = 0;
17       intensity.val[2] = 0;
18     }
19   }
20 }
```

The resulting image is the one at Figure 18. From the time perspective, as said before, being the operations for the noise visualizer less complex, the CPU is able to generate the image in around 250ms.

5.0.2 CUDA implementation

The same considerations and optimization steps done for the heat map can be done also for the generation of the frames that allows the user to better visualize the noise in the frames. So, starting from an already optimized memory management, the resulting kernel call is based as always on first copy the frame into the Device Global memory, call the kernel and then retrieve the image, as in the code below:

```

1  cudaGetDeviceProperties(&prop, 0);
2  int threads = prop.maxThreadsPerBlock;
3
4  cudaMemcpy(d_curr, image2, W*H*C * sizeof *image2, cudaMemcpyHostToDevice);
5
6  // Kernel call /4
7  kernel_red<<<1, threads>>>(d_curr, d_prev, ((W*H*C)/threads/
8  (sizeof(chunk_t)), d_out);
9  cudaMemcpy(heatmap, d_out, W*H*C * sizeof *heatmap, cudaMemcpyDeviceToHost);
```

In this case, the optimization described in the pixel difference generation, regarding the *vectorized access* to memory is also applied here, where a new type called `chunk_t` is defined. In this way more byte can be read with a single memory access.

```

1  typedef int4 chunk_t;
2
3  __global__ void kernel_red(uint8_t *current, uint8_t *previous,
4    int maxSect, uint8_t* d_heat_pixels) {
5    int x = threadIdx.x + blockDim.x * blockIdx.x;
6    int start = x * maxSect;
7    int max = start + maxSect;
8    chunk_t cc, pc;
9    uint8_t redColor = 0;
10   int df;
11   int size = sizeof(chunk_t);
12
13   for (int i = start; i < max; i++) {
14     cc = ((chunk_t *)current)[i];
15     pc = ((chunk_t *)previous)[i];
16
17     for (int j = 0; j < size; j++) {
18       df = ((uint8_t *)&cc)[j] - ((uint8_t *)&pc)[j];
19
20       if ((df < -LR_THRESHOLDS || df > LR_THRESHOLDS)) {
21         redColor = 255;
22       }
23   }
```

```

24     if ((( i * size ) + j ) % 3 == 2){
25         noise_visualization [ ( i * size ) + j ] = redColor;
26         redColor = 0;
27     } else {
28         noise_visualization [ ( i * size ) + j ] = 0;
29     }
30 }
31 }
```

By iteration on all the color components of all the pixels in the frame, the difference with the previous frame is checked. If it is above the threshold, the value of the red color for that pixel is set to 255.

Being the red component the third in a pixel, it is set accordingly to the `redColor` variable only if the `i*size + j % 3 == 2` (the red color).

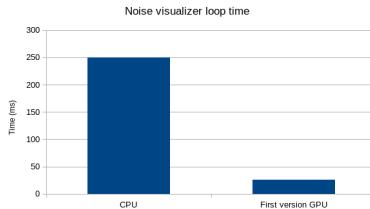
This implementation allows to copy the frame from Host to Device, compute the noise visualization and retrieve back the image from Device to Host memory in around 26ms.

The following table shows some meaningful data extracted from the profiler:

Type	Time (%)	Avg	Min	Max	Name
GPU activities	81.38	21.759ms	19.489ms	32.942ms	kernel_red
	9.59	2.5469ms	2.2134ms	6.8552ms	[CUDA memcpy HtoD]
	9.03	2.4141ms	2.1828ms	3.5408ms	[CUDA memcpy DtoH]
API calls	87.41	13.992ms	2.7579ms	156.01ms	cudaMemcpy
	12.13	194.86ms	3.8501ms	576.73ms	cudaMalloc

Table 15: Profiling result v1.cu for noise visualizer

As you can see, the average time needed to compute the kernel is less with the respect to the heat map. This is due to the simplicity of the function that only need to check the difference against a threshold. As you can see, the GPU acceleration has a huge impact over the computation.



5.0.3 CUDA reduction of global memory accesses

The fact that the only pixels that changes is the red one, the previous CUDA kernel can be optimized even more in order to reduce the number of access to the global memory. In fact, if we set assume that at the beginning all pixels in the `d_heat_map` array are all completely black, the algorithm can reduce the memory access of on third since it will work only on a the red color component. As before, if the difference is above a threshold, the red pixel is turned to red, or to black otherwise.

The resulting kernel is the following:

```

1 __global__ void kernel_red( uint8_t *current , uint8_t *previous ,
2     int maxSect , uint8_t* d_heat_pixels) {
3     int x = threadIdx.x + blockDim.x * blockIdx.x;
4     int start = x * maxSect;
5     int max = start + maxSect;
6     chunk_t cc , pc;
7     uint8_t redColor = 0;
8     int df;
9     int size = sizeof(chunk_t);
10
11    for (int i = start; i < max; i++) {
12        cc = ((chunk_t *)current)[i];
13        pc = ((chunk_t *)previous)[i];
```

```

14     for (int j = 0; j < size; j++) {
15         df = ((uint8_t *)&cc)[j] - ((uint8_t *)&pc)[j];
16
17         if ((df < -LR_THRESHOLDS || df > LR_THRESHOLDS)) {
18             redColor = 255;
19         }
20
21         if (((i * size) + j) % 3 == 2) {
22             noise_visualization[(i * size) + j] = redColor;
23             redColor = 0;
24         }
25     }
26 }
27 }
```

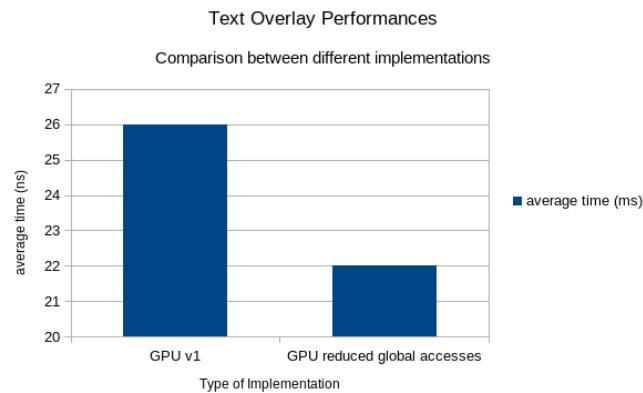
As you can notice, the only difference is that the else condition is the last if statement has been removed. This means that the green and blue component are not modified, thus the global memory is left as is for these two bytes.

This allowed to increase the performance of the kernel even more; as you can see from the table below the kernel now lasts in average 16.684ms, that is a 23% faster. The entire loop is done in about 22ms.

Moreover, the test must be performed in the same conditions because, depending on the video captured, it can show different number of changed pixels with the value of the difference above or below the threshold. For this reason, the *nvprof* has been run while the webcam was covered with a black tape. Due to the fashion of the algorithm, also the minimum and maximum values have been included in the table below.

Type	Time (%)	Avg	Min	Max	Name
GPU activities	75.63	16.684ms	16.372ms	35.942ms	kernel_red
	12.37	2.6294ms	2.2265ms	6.3841ms	[CUDA memcpy HtoD]
	12.00	2.6294ms	2.2265ms	6.3841ms	[CUDA memcpy DtoH]
API calls	86.00	11.610ms	2.6999ms	41.414ms	cudaMemcpy
	13.56	183.71ms	3.7352ms	543.50ms	cudaMalloc

Table 16: Profiling result *v2.cu* for noise visualizer



The code of this and the previous implementation has been disassembled using the NVIDIA tools *cuobjdump* and *nvdasm*; the last one is also needed to get information about the control flow. In order to avoid that the *nvcc* compiler to avoid, with very limited exceptions, dead-code eliminations and register-spilling optimizations. For this reason the code has been compiled with the *-O0* flag.

The generation of the assembly is:

```

1 cuobjdump ./a.out -xelf all
2 nvdasm v1.sm_30.cubin
```

By looking carefully at the assembly code of the first implementation, it's possible to understand the correlation between the assembly code and the instructions. As we would have expected, in the code of the first implementation are present two **ST.E.U8** instructions, used to perform the if and else case.

The following is an extract of the assembly that manage the if statement on the first CUDA implementation.

```

1 // Perform comparison
2 /*0ab0*/ ISETP.EQ.X.AND P0, PT, R5, RZ, PT;
3 /*0ab8*/ PSETP.AND.AND P0, PT, !P0, PT, PT;
4 /*0ac8*/ SSY `(.L_16);

5
6 // Go to if
7 /*0ad0*/ @P0 BRA `(.L_17);

8
9 // Go to else
10 /*0ad8*/ BRA `(.L_18);

11
12
13 // IF
14 .L_17:
15 ...
16 /*0e10*/ ST.E.U8 [R4], R0;
17 /*0e28*/ NOP.S (*"TARGET= .L_16 " *);

18 // ELSE
19 .L_18:
20 ...
21 /*0c50*/ ST.E.U8 [R4], R0;
22 /*0c88*/ NOP.S (*"TARGET= .L_16 " *);

23
24 // Management of the loop control flow
25 .L_16:
26 ...
27
28

```

Figure 21: Assembly extract of the first CUDA implementation of the noise visualizer

So, label **.L_17** and **.L_18** are used to identify the branch that is used to write 0 or 255, accordingly to the fact that it is or not the red component and only if above the threshold.

On the other hand, the assembly of the second implementation shows only one instance of the store byte instruction, as we would have expected:

```

1 // Perform comparison
2 /*0aa8*/ IADD32I RZ.CC, R4, -0x2;
3 /*0ab0*/ ISETP.EQ.X.AND P0, PT, R5, RZ, PT;
4 /*0ab8*/ PSETP.AND.AND P0, PT, !P0, PT, PT;
5 /*0ac8*/ PRMT R0, R29, 0x7610, R0;
6 /*0ad0*/ PSETP.AND.AND P0, PT, P0, PT, PT;
7 /*0ad8*/ PRMT R0, R0, 0x7610, R0;
8 /*0ae0*/ SSY `(.L_16);
9 /*0ae8*/ @P0 NOP.S (* "TARGET= .L_16 *");
10 /*0af0*/ BRA `(.L_17);

11
12
13 // IF
14 .L_17
15 ...
16 /*0e10*/ ST.E.U8 [R4], R0;
17 /*0e28*/ NOP.S (* "TARGET= .L_16 *");

18
19 // Management of the loop control flow
20 .L_16
21 ...

```

Figure 22: Assembly extract of the second CUDA implementation of the noise visualizer

5.0.4 Overlapping implementation

The fact that we can have a frame that changes completely with the respect to the previous one, implies that not only the black pixels must be turned in red if above a certain threshold, but also the red ones that are below must be changed to black. This means that we if we have an image $1920 * 1080 * 3$ we can't go below $1920 * 1080$ access to Global Memory.

A second approach resort to not use the output `d_heat_pixels` as output array but instead write directly over the *previous* frame. In this way, the red pixel will be placed over the frame itself; this is a completely different representation but, if we can modify the original frame (that will be discharged at the end of the loop), this solution is feasible. The result is visible at Figure 23.



Figure 23: Red pixels over the frame

This has been done to overcome the need of setting a pixel to black when the difference between the pixels is below the threshold. With this approach, the number of write to the Global Memory should corresponds to the number of pixels whose difference is above the threshold. Hypothetically, two identical frames should not need any write access to the Global Memory.

The kernel has been modified in the following way, so that the red color is set only if at least one color component of that pixel is above the threshold:

```

1 __global__ void kernel_red(uint8_t *current, uint8_t *previous,
2 int maxSect) {
3     int x = threadIdx.x + blockDim.x * blockIdx.x;
4     int start = x * maxSect;
5     int max = start + maxSect;
6     chunk_t cc, pc;
7
8     bool toUpdate = false;
9     for (int i = start; i < max; i++) {

```

```

10     cc = ((chunk_t *)current)[i];
11     pc = ((chunk_t *)previous)[i];
12     for (int j = 0; j < sizeof cc; j++) {
13         int8_t df = ((uint8_t *)&cc)[j] - ((uint8_t *)&pc)[j];
14
15         if (df < -LR_THRESHOLDS || df > LR_THRESHOLDS) {
16             toUpdate = true;
17         }
18
19         // Update only if needed
20         if (toUpdate && (i*(sizeof cc)+j) % 3 == 2){
21             previous[i*(sizeof cc)+j] = 255;
22             toUpdate = false;
23         }
24     }
25 }
26 }
```

Obviously, also the kernel call must be modified, so that instead of copying back data from the `d_heat_map` the data are retrieved from the previous frame, that has been overwrote during the kernel execution:

```

1 cudaMemcpy(d_current , image2.data , W*H*C * sizeof *image2.data ,
2                                     cudaMemcpyHostToDevice);
3 kernel<<<1, threads>>>(d_current , d_previous , ((W*H*C)/threads)/
4                                     (sizeof(chunk_t)));
5
6 // Rewrite d_previous frame
7 cudaMemcpy(res.data , d_previous , W*H*C * sizeof *res.data ,
8                                     cudaMemcpyDeviceToHost);
```

This implementation, in the same test condition that has been used also for the other two implementations of the noise visualizer, allows to perform the copy of the frame from Host to Device, compute the noise visualization and copy it from the Global memory of the Device to the Host in around 16.5ms.

These are the results from `nvprof`:

Type	Time (%)	Avg	Min	Max	Name
GPU activities	53.05	7.5647ms	4.9421ms	51.426ms	kernel_red
	23.69	3.3450ms	2.4960ms	6.6161ms	[CUDA memcpy HtoD]
	23.26	3.3160ms	2.8359ms	4.2297ms	[CUDA memcpy DtoH]
API calls	81.86	7.7539ms	2.9697ms	57.082	cudaMemcpy
	17.51	166.67ms	5.4298ms	327.91ms	cudaMalloc

Table 17: Profiling result `v3.cu` for noise visualizer, overwriting the red pixels over the previous frame

As you can see from the time for the kernel execution, there is a huge improvement from the black and red implementation of about 25% more.

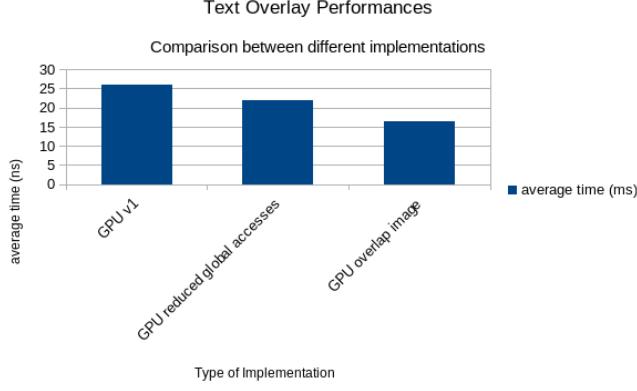


Figure 24: Comparison between GPU versions for noise visualizer

5.1 Fastest implementation

By exploiting the same idea described in the previous paragraph, the algorithm can be even more optimized from the point of view of the execution time. In fact, if instead of taking the two images and compute the difference at pixel level in order to decide if the pixel will be turned red or left as is, the difference computed from the dedicate kernel can be exploited.

By referring to the Section 2, the output of the kernel that computes the difference pixel by pixel in order to send the difference is stored into the `h_diff`. Moreover, an array called `h_xs` contains the locations in which the pixels are changed.

The idea is to exploit the `h_xs` array in order to turn red the corresponding pixels. Its length is defined by an integer value called `h_pos`. Having `nMaxThreads` threads and `h_pos` pixels that must be turned to red, the total threads can work on a portion of `h_xs` that is large:

$$maxSect = \frac{nMaxThreads}{h_xs}$$

This implementation allows to have a huge speedup, since the `h_xs` has its corresponding vector already in the Device, called `d_xs`; for this reason, it's enough to simply pass the point, without previously perform any memory copy from Device to Host.

The final CUDA kernel call is like the following C code:

```

1 // Compute pixel difference
2 kernel2<<<1, nMaxThreads, 0>>>(d_current, d_previous, d_diff, max4, d_pos, d_xs);
3 cudaMemcpyAsync(&pready->h_pos, d_pos, sizeof *d_pos, cudaMemcpyDeviceToHost);
4 cudaDeviceSynchronize();
5
6 // Compute red map over previous image
7 red_black_map_overlap<<<1, nMaxThreads, 0>>>(d_pos, d_xs,
8                                         ceil(pready->h_pos/nMaxThreads), d_previous);

```

The kernel code will turn the pixels defined in the `xs` array into red. As said before, each thread will work on a portion of data long `maxSect` so that the computation load is split.

The initial problem was to find a way to identify the position of the red component given an arbitrary position of any of the three red, green or blue pixels. The image representation of OpenCV `Mat` object is slightly misleading; in fact, the blue component is at position 0 while the red one at position 2. So, an image 3x3 pixels is build in the following way:

B	G	R	B	G	R	B	G	R
B	G	R	B	G	R	B	G	R
B	G	R	B	G	R	B	G	R

Table 18: Image color component representation, with dimension 3x3

At this point, we know that the red color component is always int the *third* position and so, given an *absolute* index i , the corresponding red color component position is identified as:

$$\text{red_position} = i + (2 - i \% 3)$$

By applying this idea, the corresponding CUDA kernel is the following one

```

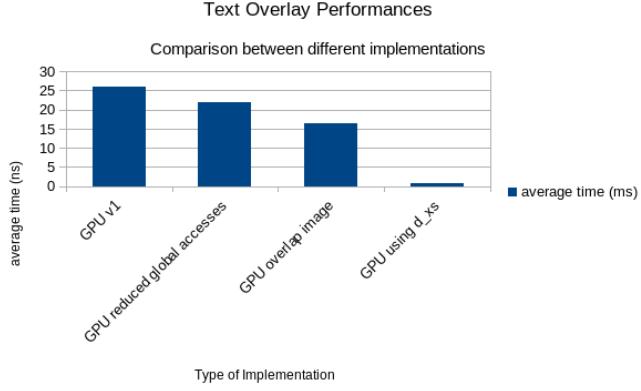
1 __global__ void red_black_map_overlap(unsigned int *pos, int *xs,
2                                     int maxSect, uint8_t *noise_visualization){
3     int x = threadIdx.x + blockDim.x * blockIdx.x;
4     int start = x * maxSect;
5     int max = start + maxSect;
6
7     for (int i = start; i < max && i < *pos; i++) {
8         noise_visualization[xs[i] + (2 - xs[i] % 3)] = 255;
9     }
10 }
```

These are the results from *nvprof*, with $\text{pos} \simeq 70000$ (number of changed pixels is around 70000):

Type	Time (%)	Avg	Min	Max	Name
GPU activities	7.22	915.35us	624.23us	6.7692ms	kernel_red

Table 19: *Profiling result for noise visualizer, overwriting the red pixels over the previous frame using `d_xs` array. The [CUDA memcpy DtoH] and [CUDA memcpy HtoD] are not reported because they can be misleading, since they includes also the cost for the kernel that computes the difference*

By considering only the kernel time, the computation time is reduced by 94.5%, that is pretty impressive.



6 Noise filter

The sensor of an digital camera contains a certain number of pixel, used to acquire the image. The settings of a camera allows to adapt the sensor to the external conditions in order to always get the best shot). In a sunny day, a slow shutter speed, a wide open aperture allows a large number of photons to hit completely the image sensor.

On the other hand, in darker conditions, faster shutter speeds and a closer aperture do not allow enough photons to hit the sensor to generate a suitable voltage to overcome the noise present at the ground. This leads to a low Signal-To-Noise ratio that is directly converted into noise from the internal camera logic. In this case, a correct trim of the of ISO and the shutter aperture time allow to have an image with lower noise; but, in most of the cases, the cheap webcams do not have these advanced configurable feature but perform an auto-adapt procedure.

This is where the *noise filtering* take a major role for this project. Being the data sent based on

the difference with the previous frame, each pixel that is mainly driven by the noise, has higher probability to flip drastically change its value. So, even if the pixel itself does not contain any information, its difference is still sent.

6.1 Convolution filters

Convolution filters are used to perform some perturbation over an image in order to obtain some enhancement. They are based on a *kernel*, a convolution matrix, that is used to perform the convolution between the image and the kernel itself, in the following way:

$$g(x, y) = \sum_{u=-\frac{k}{2}}^{\frac{k}{2}} \sum_{v=-\frac{k}{2}}^{\frac{k}{2}} K[u, v] \cdot I(x - u, y - v)$$

Where K is the kernel and I is the original image; g will be the perturbed resulting image. Depending on the type of selected kernel, the image will be modified in a different way. By defining the size of the kernel as 3×3 , the identity is the following:

$$K = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

In this way, the original image is not modified; let's check it:

$$\begin{aligned} I * K &= \left(\begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} * \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \right) = \\ &= (i \cdot 0) + (h \cdot 0) + (g \cdot 0) + (f \cdot 0) + (e \cdot 1) + (d \cdot 0) + (c \cdot 0) + (b \cdot 0) + (a \cdot 0) \\ &= e \end{aligned}$$

6.2 Mean filter

The idea of the Mean filter, also called *blur box*, is to set the kernel with the elements that have all the same weight but normalized, like in the following 3×3 convolution matrix:

$$K = \frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.11 & 0.11 & 0.11 \\ 0.11 & 0.11 & 0.11 \\ 0.11 & 0.11 & 0.11 \end{bmatrix}$$

It can be simply computed in the following way, assuming the kernel is on $K \times K$:

```

1 float* computeMeanKernel(int K){
2     double sum = 0;
3     float* k = (float*)malloc(K*K*sizeof(float));
4     for (int i = 0; i < K; i++){
5         for (int j = 0; j < K; j++){
6             k[i*K+j] = 1.0/(K*K);
7         }
8     }
9
10    return k;
11 }
```

In this way, the image is basically blurred due to the fact that the pixel at i position is based now also on the values of the surrounding 8 pixels.

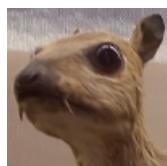


Figure 25: Original

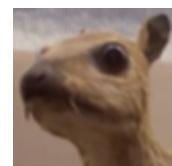


Figure 26: Mean filtering 3×3

This solution can be exploited in order to make the pixels that suffers from the noise problem more smooth with the respect to the neighbors ones. In this way, even if the two consecutive frames are not exactly the same, the difference among the single pixels will be lower, resorting to a lower data sent over the socket. Let's take for simplicity two frames that are in gray scale, between 0 and 255, that are described as two 3×3 matrices.

$$A = \frac{1}{9} \begin{bmatrix} 120 & 131 & 112 \\ 112 & 101 & 82 \\ 44 & 106 & 65 \end{bmatrix} \quad B = \begin{bmatrix} 120 & 139 & 90 \\ 99 & 126 & 106 \\ 46 & 75 & 88 \end{bmatrix}$$

If we simply compute the difference, and we set the threshold to 20, by computing the difference location by location, 5 pixels must be sent as difference. If we apply now the kernel to both the two frames, in order to simulate the behaviour of the algorithm, the two resulting matrices are the following:

$$A = \frac{1}{9} \begin{bmatrix} 51 & 73 & 47 \\ 68 & 96 & 66 \\ 40 & 56 & 39 \end{bmatrix} \quad B = \begin{bmatrix} 53 & 75 & 51 \\ 67 & 98 & 69 \\ 38 & 60 & 43 \end{bmatrix}$$

As you can see now, no pixels need to be sent. This is a huge improvement, since we have saved 100% of the bandwidth if we neglect the overhead to transmit the information that no pixels have changed. Even if this seems an optimal solution, it has one major drawback. The image will become distorted from the details point of view. In fact, a trade-off between image quality and saved bandwidth must be selected.

6.3 Gaussian filter

The second type of filter that has been selected to reduce the noise, is the *Gaussian filter*, that allows to reach a similar result with the respect to the Mean one. The idea is always to distort the image to reduce the noise and make the difference among pixels lower. The matrix of a 3×3 Gaussian filter is the following:

$$K = \frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$

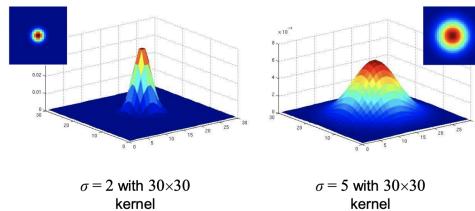
It can be computed using the following formula:

$$K(x, y) = \frac{1}{2 \cdot \pi \cdot \sigma^2} \cdot e^{-\frac{x^2+y^2}{2\sigma^2}}$$

It is called Gaussian filter because of its shape, corresponds to a Gaussian function. A filter of this type requires $6\sigma - 1$ values (to make it sure it's odd), so that for a σ of 3 the kernel must be 17 elements long. The σ value is given in pixels. If we use a square kernel its size can be simply computed as:

$$\text{side}_{\text{length}} = \sqrt{6 \cdot \sigma - 1}$$

Where $\text{side}_{\text{length}}$ is one of the length of the matrix that composes the kernel. A visual example representation is available in the following image.



Gaussian filter weights pixels into a bell-shape curve around the central pixel. This means that the farther the pixel from the center, the lower the weight. On the other hand the Mean filter is just an average between the center and the neighboring pixels. This means that for all the pixels, equal weights are assigned, regardless the distance from the center.



Figure 27: *Original*

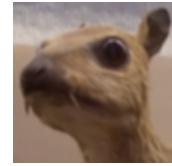


Figure 28: *Gaussian filter 3x3*

The C code, that give σ , computes the entire convolution matrix is the following:

```

1 float* computeGaussianKernel( float s, int K){
2     float sum = 0;
3     float* k = ( float*) malloc(K*K*sizeof( float ));
4     for ( int i = 0; i < K; i++ ){
5         for ( int j = 0; j < K; j++ ){
6             double x = i - (K - 1) / 2.0;
7             double y = j - (K - 1) / 2.0;
8             k[ i*K+j ] = (1.0/(2.0*M_PI*s*s)) * exp(-((x*x + y*y)/(2.0*s*s)));
9             sum += k[ i*K+j ];
10        }
11    }
12    for ( int i = 0; i < K; i++ ) {
13        for ( int j = 0; j < K; j++ ) {
14            k[ i*K+j ] /= sum;
15        }
16    }
17
18    return k;
19 }
```

6.4 Tiled CUDA 2D implementation

In order to test the effectiveness of the algorithm and if it is suitable for this project, a configurable implementation has been developed in CUDA. In order to have all the same test environment, two sequential frames have been taken from the webcam in a semi dark environment with a person in the center.



Figure 29: *First frame*



Figure 30: *Second frame*

As you can see, the two images are pretty similar in term of content; a trained eye can see that the head of the second frame is slightly tilted to the right with the respect to the previous one. The idea is to run now the *noise visualizer*, refer to Section 5, on these two frames and understand where the noise is located.



Figure 31: *Red-black representation of the noise*

The noise, is the one in the background that do not contribute in any way to the video flow. In fact, even if it is fixed, it is sent anyway via the socket, increasing the used bandwidth. The CUDA implementation is based on the 2D tiled version but includes the management for all three color components.

The entire algorithm is based on the usage of the shared memory and the synchronization among the thread. The idea is to make each thread work on single pixel and apply the kernel, by using the portion of the image that is in the shared memory, called *tile*.

The kernel implementation is composed of multiple parts:

- 1. Shared memory initialization and boundary condition management:** As introduced before, each thread, identified by its *ty* and *tx* coordinates should work on a precise pixels, identified as:

```

1 __shared__ int N_ds[BLOCK_SIZE][BLOCK_SIZE*3];
2
3 x = threadIdx.x + blockIdx.x*TILE_SIZE - K/2;
4 y = threadIdx.y + blockIdx.y*TILE_SIZE - K/2;

```

In fact, the idea is to create a block of threads that corresponds to $BLOCK_SIZE \times BLOCK_SIZE$ as dimension, so that each one of the loads into the shared memory the value of its pixel (the multiplication by 3 in the kernel implies that each threads read all three color components). The following code, instead, includes also the write into the shared memory.

```

1 int tx = threadIdx.x; //W
2 int ty = threadIdx.y; //H
3 int row_o = blockIdx.y*TILE_SIZE + ty;
4 int col_o = blockIdx.x*TILE_SIZE + tx;
5 int row_i = row_o - K/2;
6 int col_i = col_o - K/2;
7
8 if(row_i >= 0 && row_i < H && col_i >= 0 && col_i < W){
9     N_ds[ty][tx*3] = image[row_i*W*3+ col_i * C];
10    N_ds[ty][tx*3+1] = image[row_i*W*3+ col_i*3 + 1 ];
11    N_ds[ty][tx*3+2] = image[row_i*W*3+ col_i*3 + 2 ];
12 } else {
13     N_ds[ty][tx*3] = 0;
14     N_ds[ty][tx*3+1] = 0;
15     N_ds[ty][tx*3+2] = 0;
16 }

```

The value *TILE_SIZE* defines the dimension of the tile, in which each the kernel is applied. The problem is that, the first pixel (position [0,0]) in the image, if the kernel has a dimension of 3x3, need to access also the location [-1,-1], [-1, 0], [-1, 1], [0, -1], [1, -1]. In this case, they are not available. This is the reason of the presence of the if condition, that sets 0 at these locations in the shared memory by "shifting" the indexes by *K/2*. The problem is shown in the below picture:

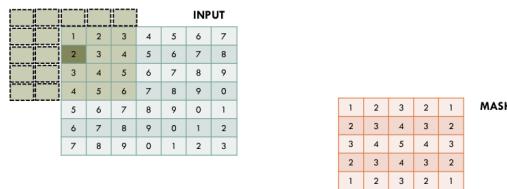


Figure 32: Boundary conditions among the image border

In this way, all the pixels needed to process a tile of *TILE_SIZE*x*TILE_SIZE* dimension, are loaded into the *N_ds* shared memory that is large *BLOCK_SIZE***BLOCK_SIZE**3 in order to contains also the values at the boundary of the tile.

- 2. Filter application:** Foreach element in the tile and for all three color components, the kernel is computed using the shared memory.

The entire kernel is the following:

```

1 __constant__ float dev_k[K*K];
2
3 __global__ void convolution_kernel(uint8_t *image, uint8_t *R)
4 {
5     __shared__ int N_ds[BLOCK_SIZE][BLOCK_SIZE*3];
6
7     int tx = threadIdx.x; //W
8     int ty = threadIdx.y; //H
9     int row_o = blockIdx.x*y*TILE_SIZE + ty;
10    int col_o = blockIdx.x*x*TILE_SIZE + tx;
11    int row_i = row_o - K/2;
12    int col_i = col_o - K/2;
13
14    // Management of the boundary conditions
15    if(row_i >= 0 && row_i < H && col_i >= 0 && col_i < W){
16        N_ds[ty][tx*3] = image[row_i*W*3+ col_i * C];
17        N_ds[ty][tx*3+1] = image[row_i*W*3+ col_i*3 + 1 ];
18        N_ds[ty][tx*3+2] = image[row_i*W*3+ col_i*3 + 2 ];
19    } else {
20        N_ds[ty][tx*3] = 0;
21        N_ds[ty][tx*3+1] = 0;
22        N_ds[ty][tx*3+2] = 0;
23    }
24
25    __syncthreads();
26
27    if(row_o < H && col_o < W){
28        float outputR = 0.0;
29        float outputG = 0.0;
30        float outputB = 0.0;
31        if(ty < TILE_SIZE && tx < TILE_SIZE){
32            for(int i = 0; i < K; i++)
33                for(int j = 0; j < K; j++){
34                    outputR += dev_k[i*K+j] * N_ds[i+ty][(j+tx)*3];
35                    outputG += dev_k[i*K+j] * N_ds[i+ty][(j+tx)*3 + 1];
36                    outputB += dev_k[i*K+j] * N_ds[i+ty][(j+tx)*3 + 2];
37                }
38
39        R[row_o*W*3 + col_o*3] = outputR;
40        R[row_o*W*3 + col_o*3 + 1] = outputG;
41        R[row_o*W*3 + col_o*3 + 2] = outputB;
42    }
43 }
44 }
```

The kernel call needs to work on an image of $1920 \times 1080 \times 3$ elements divided into square tiles of $TILE_SIZE \times TILE_SIZE$ dimensions. Since the tile needs also the surrounding element to compute the kernel, the block size is composed of $BLOCK_SIZE \times BLOCK_SIZE$ elements, where:

$$BLOCK_SIZE = TILE_SIZE + K - 1$$

The grid instead, is configured to generate as many tile as needed from the image. In this case, a grid of $W/TILE_SIZE \times H/TILE_SIZE$ is enough (need to be rounded to the upper value). This is the complete kernel call

```

1 dim3 blockSize, gridSize;
2 blockSize.x = BLOCK_SIZE, blockSize.y = BLOCK_SIZE, blockSize.z = 1;
3 gridSize.x = ceil((float)W/TILE_SIZE),
4 gridSize.y = ceil((float)H/TILE_SIZE),
5 gridSize.z = 1;
6
7 // Kernel call
8 convolution_kernel<<<gridSize, blockSize>>>(d_current, d_current_filtered);
```

This solution includes already optimal solution in order to increase the performance and reduce at minimum the overhead introduced by the computation of the filter. In fact, by having defined how convolution if performed, it can be simple called, after have set the kernel in the constant memory.

Having defined the kernel, that once computed is fixed for the entire program execution, allows to not cope with the cache coherence management, allowing data to be aggressively cached in the L1 cache (the nearest to the SM).

Moreover, access to the memory has been optimized by resorting to the usage of the shared memory. Since the larger the dimension of the tile, the higher the ratio between the manipulated pixels and the memory accesses, an optimal solution is to increase as much as possible the tile size. The limit is the shared memory, that is limited.

The execution of the kernel is directly proportional to the size of the tile and the kernel dimension; the code shown before has been run for multiple sessions in order to be able to plot the kernel time for multiple convolution matrix and multiple tile size.

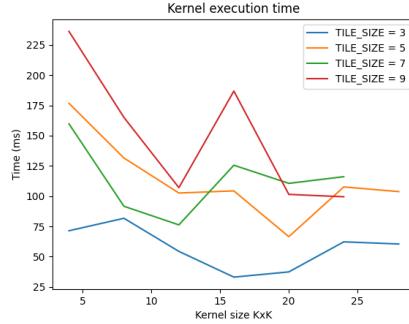


Figure 33: *Kernel execution time from profiler*

Both the Gaussian and the Mean filters have been run over the two frames used as test dataset, the results are available at Table 34. The number of pixels changed, if no filter is applied, is 369350, that is 5.93% of the whole image.

Filter Type	Kernel size	Total saved pixels (%)
Mean	3	3.37
	5	2.31
	7	1.66
	9	1.24
Gaussian	3 ($\sigma = 1$)	3.58
	4 ($\sigma = 2$)	2.87
	5 ($\sigma = 3$)	2.37
	6 ($\sigma = 5$)	1.98
	7 ($\sigma = 8$)	1.66

Figure 34: *Results of both Gaussian and Mean filter, with different kernel size*

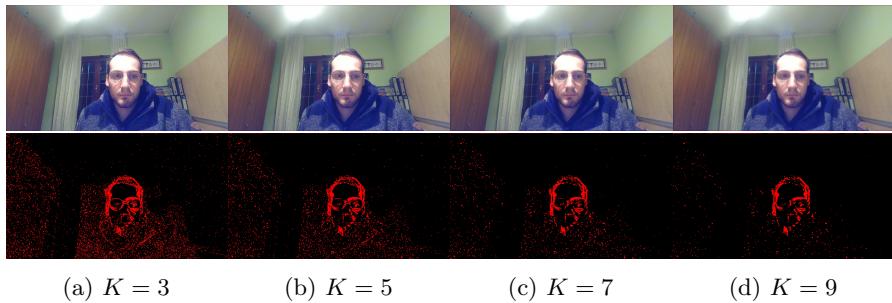


Figure 35: *Mean filtering results*

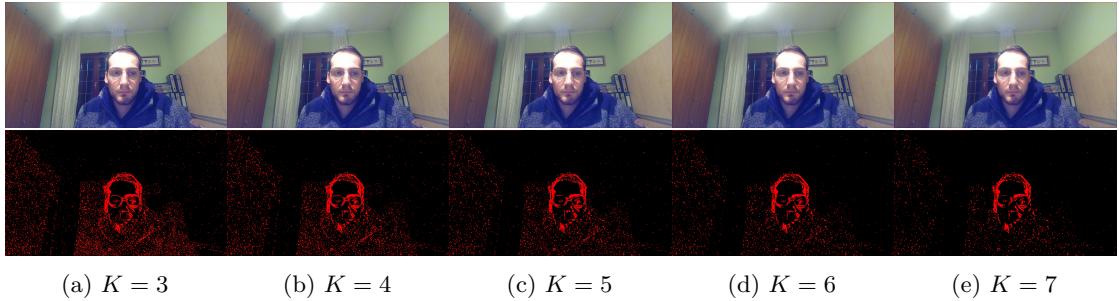


Figure 36: Gaussian filtering results

By visually inspecting the results, by using a kernel that is too large, leads to a non-negligible distortion to the image due to the blur effect. By the way, the team has defined that the best visual result is when $K = 4$, using the Gaussian Filter, that allows to avoid to send about 10600 pixels, with the respect to the selected two frames.

6.5 2D CUDA Tiled Median filter

Another algorithm that has been selected to try to remove the noise from the image is the Medial filter. It works in a completely different way from the Mean and Gaussian filters. The idea behind of it is to preserve the useful details of the image while removing the noise.

Like the Mean filter, the Median one takes N neighbor pixels that are next to the selected one but, instead of replacing it with the average value between all of them it replace that pixel with the *median*.

The median of a plan one dimensional array can be computed by first sorting all elements and then take the one in the middle of the array itself. On the other hand, if we have a matrix, the media con be computed in the same way, but all elements are considered and ordered.

$$\begin{bmatrix} 1 & 7 & 3 \\ 2 & \mathbf{9} & 6 \\ 4 & 1 & 6 \end{bmatrix} \implies [1, 1, 2, 3, \mathbf{4}, 6, 6, 7, 9] \implies 4$$

Table 20: Median of a matrix 3x3

In this case, the example above shows the media computation of a matrix 3x3 that can corresponds to the kernel size of CUDA kernel.

The CUDA algorithm is structured in the same way as the convolution one, in which the image processing is divided into *tiles* of fixed dimensions. In the same way, pixels are loaded into the shared memory in order to reduce the accesses to the Global memory.

So, the first step to do is to generate a CUDA function that is able to extract the media from a given array, of size $K * K$. The following code shows the implementation of the Bubble Sort and take the median element, by assuming that N is odd:

```

1 __device__ uint8_t median(uint8_t *array, int N){
2     bool swapped = true;
3     for (int a = 0; a < N && swapped; a++){
4         swapped = false;
5         for (int i = 0; i < N - 1; i++){
6             if(array[i] > array[i+1]){
7                 uint8_t tmp = array[i];
8                 array[i] = array[i+1];
9                 array[i+1] = tmp;
10                swapped = true;
11            }
12        }
13    }
14    return array[N/2];
15}
16

```

At this point, the kernel must load into the shared memory the image and then, foreach thread three arrays are computed; one for each of the color components. By iterating of the tile, elements are put in the three arrays and the median extracted. At this point, the new image is computed using that median. The following code shows the kernel for the media filtering.

```

1  __global__ void median_kernel(uint8_t *image, uint8_t *R)
2 {
3     __shared__ uint8_t N_ds[BLOCK_SIZE][BLOCK_SIZE*3];
4
5     int tx = threadIdx.x; //W
6     int ty = threadIdx.y; //H
7     int row_o = blockIdx.y*TILE_SIZE + ty;
8     int col_o = blockIdx.x*TILE_SIZE + tx;
9     int row_i = row_o - K/2;
10    int col_i = col_o - K/2;
11
12    if(row_i >= 0 && row_i < H && col_i >= 0 && col_i < W){
13        N_ds[ty][tx*3] = image[row_i*W*3+ col_i * C];
14        N_ds[ty][tx*3+1] = image[row_i*W*3+ col_i*3 + 1];
15        N_ds[ty][tx*3+2] = image[row_i*W*3+ col_i*3 + 2];
16    } else {
17        N_ds[ty][tx*3] = 0;
18        N_ds[ty][tx*3+1] = 0;
19        N_ds[ty][tx*3+2] = 0;
20    }
21
22    __syncthreads();
23
24    if(row_o < H && col_o < W && ty < TILE_SIZE && tx < TILE_SIZE){
25        uint8_t medR[K*K];
26        uint8_t medG[K*K];
27        uint8_t medB[K*K];
28        int r = 0;
29        int g = 0;
30        int b = 0;
31
32        // Copying elements into the three arrays
33        for(int i = 0; i < K; i++)
34            for(int j = 0; j < K; j++){
35                medR[r++] = N_ds[i+ty][(j+tx)*3];
36                medG[g++] = N_ds[i+ty][(j+tx)*3 + 1];
37                medB[b++] = N_ds[i+ty][(j+tx)*3 + 2];
38            }
39
40        // Computation of the median
41        R[row_o*W*3 + col_o*3] = median(medR, K*K);
42        R[row_o*W*3 + col_o*3 + 1] = median(medG, K*K);
43        R[row_o*W*3 + col_o*3 + 2] = median(medB, K*K);
44    }
45}

```

The problem with this implementation is that, it needs to sort an array of $K * K$ elements for each color component in the image, that obviously implies a longer computation. The below images are from the execution of the above algorithm, using $K = 5$ and a frame of $1920*1080*3$ elements. The kernel execution time from the *profiler* is 574.67ms, that is a lot if we want to keep the FPS sufficiently high.



Figure 37: Median filtered frame

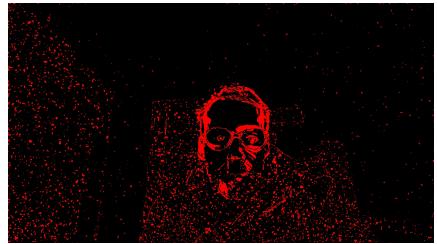


Figure 38: Red-black noise visualization

Figure 39: Example of Median filter with $K = 5$

Let's suppose for a moment that we do not have any constraint from the FPS; *is it a good solution for the noise?* From the point of view of the percentage of the total pixels saved, we have only 2.54%, that is even higher the result from a simpler Mean filter with the same kernel size.

The problem with this filter is coming from the type of noise that is characterizing our frames. In fact, the Median filter works better on *Salt and Pepper* noise, that are some random pixels that have completely different values from the point of view of the magnitude. Let's see an example:

$$\begin{bmatrix} 1 & 1 & 7 & 1 & 3 \\ 2 & 1 & \mathbf{89} & 1 & 6 \\ 4 & 1 & 1 & 1 & 6 \end{bmatrix} \implies [1, 1, 2, 3, \mathbf{4}, 6, 6, 7, 89] \implies 4$$

Table 21: *Median of a matrix 3x3, with salt and pepper noise*



Figure 40: *Salt and pepper noise*

So, due to the fact that the Median filter does not bring any advantage, it has not been optimized and even used. Beside this, the implementation allowed us to better understand the nature of the noise.