

Zadanie 6 - Rozmycie Gaussa GPU

Celem zadania numer 6 było wdrożenie złożonego programu umożliwiającego zastosowanie filtru - rozmycie Gaussa na materiał video. Głównym założeniem tej aplikacji było zaimplementowanie własnego rozwiązania rozmycia z wykorzystaniem mocy obliczeniowej kart graficznych firmy Nvidia, za pomocą architektury CUDA. Do wykonania należało także wykorzystać dostępną bibliotekę OpenCV, która umożliwia manipulację strumieniem wejściowym video, w postaci macierzy zawierających informację o pikselach (obiektu typu *Mat*).

Poniżej została przedstawiona lista kroków do realizacji zadania:

- pobranie argumentów
- obliczenie gridu, bloku i pobranie liczby wątków
- otworzyć materiał filmowy, klatka po klatce
- zastosować filtr Gaussa poprzez dokonanie modyfikacji informacji o pikselach w każdej z klatek - obliczenia zrównoległych z pomocą CUDA
- zapisać kolejno zmodyfikowane klatki
- zapisać całość po zastosowaniu odpowiedniego filtra video

Do implementacji rozmycia Gaussa została użyta tradycyjna metoda bazująca na przekształceniu pikseli przez macierz charakterystyczną dla rozmycia gaussa.

W przypadku zastosowania architektury CUDA do obliczeń równoległych, było odpowiednie przygotowanie gridu i rozdzielenie zadania na wątki.

Część programu odpowiedzialna za podział:

```
1 /**
2  * This method based on run parameters. Calculates the amount of threads/
3  * blocks for the application use.
4  * @param threadCount The amount of threads.
5  */
6 void prepareGrid(unsigned int threadCount)
7 {
8     // Round the number of threads
9     if ((threadCount % 2) == 1)
10     {
11         ++threadCount;
12     }
13
14     // Divide into blocks
15     for (int i = 512; i > 0; i--)
16     {
17         double blocks = (double) threadCount / i;
18
19         if (blocks == (int) (blocks))
20         {
21             double divThreads = sqrt(i);
22             if (divThreads == int(divThreads))
```

```

23     {
24         threadsOnX = int(divThreads);
25         threadsOnY = int(divThreads);
26     }
27     else
28     {
29         threadsOnX = i;
30         threadsOnY = 1;
31     }
32
33     break;
34 }
35 }
36 }

```

Znając specyfikację karty graficznej, na której przeprowadzone były testy programu, należało paramter wejściowy, oznaczający liczbę wątków na której działa aplikacja odpowiednio zinterpretować. Powyższy kod, dzieli na bloki (obliczanie wymiarów tzw. *grid'a*) w zależności od liczby wątków.

Implementacja Kernela CUDA:

```

1  /**
2  * CUDA implementation for Gaussian blur.
3  *
4  * @param imageIn      The pixel to perform.
5  * @param imageOut     The result of pixel conversion.
6  * @param width        The width of frame.
7  * @param height       The height of frame.
8  * @param channels      The channels color for pixel.
9  * @param kernel       The kernel for gaussian blur.
10 * @param kernelSize   The size of kernel.
11 * @param adjustedSize The adjusted size of kernel (usually the half of
12 *                    kernel size rounded down).
13 * @param sum          The sum of all kernel values.
14 */
15 __global__ void gaussBlur(unsigned char *imageIn, unsigned char *imageOut,
16                          int width, int height, int channels, float *kernel, int kernelSize,
17                          int adjustedSize, int sum)
18 {
19     const int index = getThreadId();
20
21     if (index < width * height)
22     {
23         if (!(index % width < adjustedSize
24             || index % width >= width - adjustedSize
25             || index / width < adjustedSize
26             || index / width >= height - adjustedSize))
27         {
28             float x = 0.0f;
29             float y = 0.0f;
30             float z = 0.0f;
31
32             for (int j = 0; j < kernelSize * kernelSize; ++j)
33             {
34                 // Compute index shift to neighboring cords.
35                 int shift = (j / kernelSize - adjustedSize) * width

```

```

36         + j % kernelSize - adjustedSize;
37
38     x += imageIn[(index + shift) * channels] * kernel[j];
39     y += imageIn[(index + shift) * channels + 1] * kernel[j];
40     z += imageIn[(index + shift) * channels + 2] * kernel[j];
41 }
42
43 // Apply to output image and save result.
44 imageOut[index * channels] = (unsigned char) (x / sum);
45 imageOut[index * channels + 1] = (unsigned char) (y / sum);
46 imageOut[index * channels + 2] = (unsigned char) (z / sum);
47 }
48 }
49 }

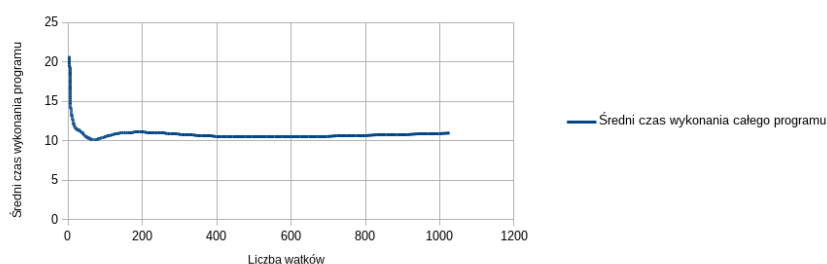
```

Przebieg

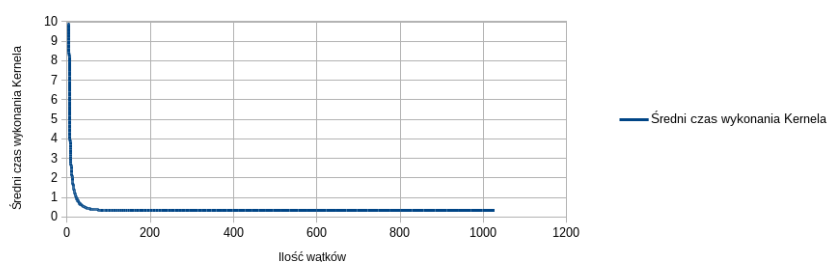
Obliczenia zostały wykonane na serwerze CUDA.

Poniżej wykresy przedstawiający rezultat przeprowadzonych operacji na wątkach.

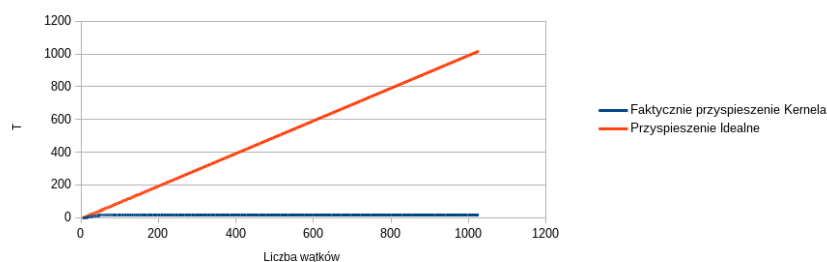
Wykres średniego czasu wykonania programu do ilości wątków



Wykres średniego czasu wykonania Kernelu do ilości wątków



Przyspieszenie wykonania Kernela względem idealnego (liniowego)



Wnioski

Najlepsze przyspieszenie udało się osiągnąć dla 128 wątków (dla filmu Helicopter DivXHT ASP.divx). Jak widać nie da się osiągnąć większego przyspieszenia, bowiem rozmiar macierzy (klatki filmu) jest ograniczony i w pewnym momencie nie ma już możliwości dalszego dzielenia obliczeń pomiędzy wątki. Dodatkowym obciążeniem mogą być konieczności przełączania kontekstów, jednak nie udało nam się znaleźć jednoznacznego potwierdzenia istnienia takiego problemu dla procesorów GPU w literaturze przedmiotu.