



# IRGPU: Harris corner detector

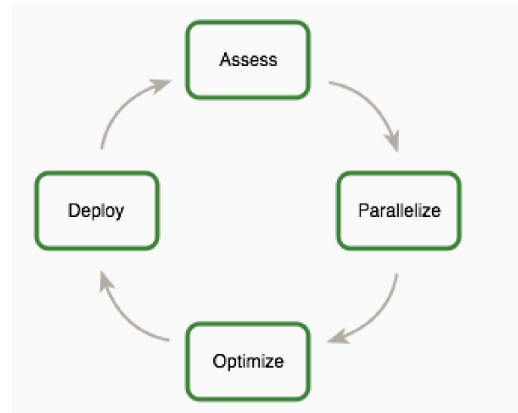
Nathan Cabasso, Ferdinand Mom - IMAGE 2023

# I/ Goal

Implement an Harris Corner detector in CPU & GPU.

## Pipeline:

- Port the Python code from lab session in C++ to get a working base.
- Port the C++ code to CUDA.
- Check all results are correct.
- Identify what can be optimized, implement, measure, repeat.





## A) CPU implementation

The Harris corner detector algorithm can be decomposed in several steps:

- **Step 1:** Convert image to grayscale.
- **Step 2:** Compute harris response on the gray scale image:
  - Compute a gaussian derivative (sobel) kernel.
  - Compute the image derivatives  $I_x$  and  $I_y$ .
  - Compute the structure tensor images  $I_x^2$ ,  $I_y^2$ , and  $I_x * I_y$
  - Convolve a gaussian filter with each structure tensor images.
  - Compute the determinant  $W_{det}$  and the trace  $W_{trace}$  based on previous convolution output.
  - Returns the Harris cornerness response of a given image ( $\frac{W_{det}}{W_{trace}+1}$ ).
- **Step 3:** Perform non-maximal suppression using a morphological opening on Harris corner response.
- **Step 4:** Find top K best keypoints coordinates by sorting the filtered responses.



## A) CPU implementation

```
1  struct Matrix {
2      Matrix();
3      Matrix(int height, int width);
4
5      Matrix operator*(const Matrix& rhs) const;
6      Matrix operator*(const double& rhs) const;
7      Matrix operator+(const Matrix& rhs) const;
8      Matrix operator+(const double& rhs) const;
9      Matrix operator-(const Matrix& rhs) const;
10     Matrix operator/(const Matrix& rhs) const;
11     Matrix operator>(const double& rhs) const;
12     Matrix operator==(const Matrix& rhs) const;
13     double max() const;
14
15     // [...]
16
17     int height, width;
18     std::vector<double> data;
19 };
```

- Abstraction for easier GPU code porting.
- Sequential implementation (No multi-threading)

## B) GPU implementation

### CPU

```
1 Matrix compute_harris_response(const Matrix &image)
2 {
3     int size = 3;
4
5     auto img_x = gauss_derivative(image, size, 1);
6     auto img_y = gauss_derivative(image, size, 0);
7
8     auto gauss = gauss_filter(size);
9
10    auto W_xx = convolution_2D(img_x * img_x, gauss, size);
11    auto W_xy = convolution_2D(img_x * img_y, gauss, size);
12    auto W_yy = convolution_2D(img_y * img_y, gauss, size);
13
14    auto W_det = (W_xx * W_yy) - (W_xy * W_xy);
15    auto W_trace = W_xx + W_yy;
16
17    return W_det / (W_trace + 1.);
18 }
```

### GPU

```
1 MatrixGPU compute_harris_response_gpu(MatrixGPU &image)
2 {
3     int size = 3;
4     auto img_x = sobel_filter_gpu(image, size, 1);
5     auto img_y = sobel_filter_gpu(image, size, 0);
6
7     auto gauss = gauss_filter_gpu(size);
8
9     auto I_xx = img_x * img_x;
10    auto I_xy = img_x * img_y;
11    auto I_yy = img_y * img_y;
12
13    auto W_xx = convolution_2D_gpu(I_xx, gauss);
14    auto W_xy = convolution_2D_gpu(I_xy, gauss);
15    auto W_yy = convolution_2D_gpu(I_yy, gauss);
16
17    auto W_det = (W_xx * W_yy) - (W_xy * W_xy);
18    auto W_trace = W_xx + W_yy;
19
20    return W_det / (W_trace + 1.);
21 }
```



## B) GPU implementation

- **Thrust**: a parallel algorithms library which resembles the C++ Standard Template Library (STL).
- **Kernels implemented**:
  - **grayscale(img)**: takes an 8-bit RGBA image and outputs a floating point grayscale representation.
  - **gauss\_filter(n)**: generates a Gaussian convolution kernel of size  $n$ .
  - **sobel\_filter(n)**: computes the gradient of a gaussian filter of size  $n$ .
  - **convolution\_2D(img, kernel)**: computes a basic 2D convolution of *img* with *kernel*. We implemented 3 variants (**naive**, **tiled**, **tiled with multiple loads per thread**).
  - **morph\_apply\_gpu(img, kernel, mode)**: computes a morphological operation to *img* with *kernel*. The parameter *mode* can be used to switch between erosion and dilatation.
  - **morph\_dilate\_gpu(img, kernel\_size)**: performs an optimized but approximated morphological dilatation using square kernel.
- CPU code for computing the top-K best keypoints coordinates differs from the GPU one as we fully relied on **thrust** library for ease of process.



## II/ Benchmark

- **Google Benchmark:** straightforward to setup and have many options to benchmark such as counting the number of maximum run of our program given a time constraint.
- **nvprof (Nvidia profiler):** enables to see the time spent on each kernel call. Useful to identify the bottlenecks.



## A) Bottleneck optimization

- Morphological operation arithmetic precision
- Numbers of filters for morphology
- Horizontal + vertical filters for morphology
- Tiled convolutions
- Early free (fail)

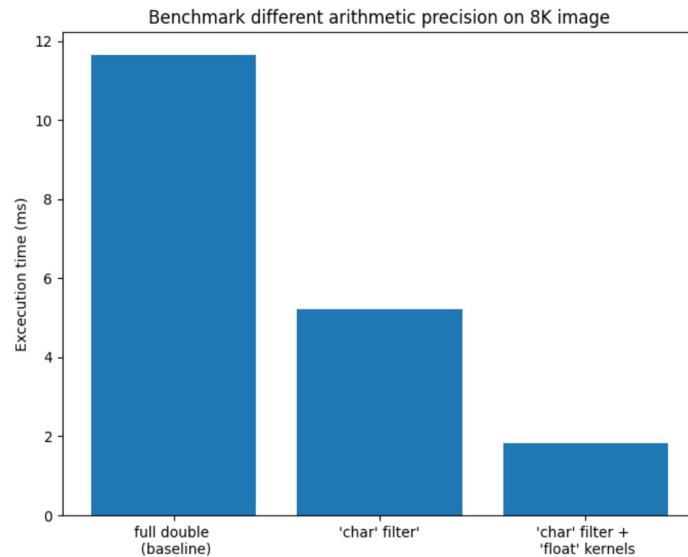




## Morphological operation arithmetic precision

Idea: Since kernel is binary, reduce the size of the kernel by using char instead of double.

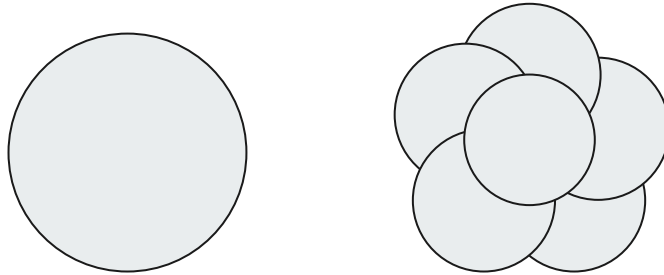
# Morphological operation arithmetic precision



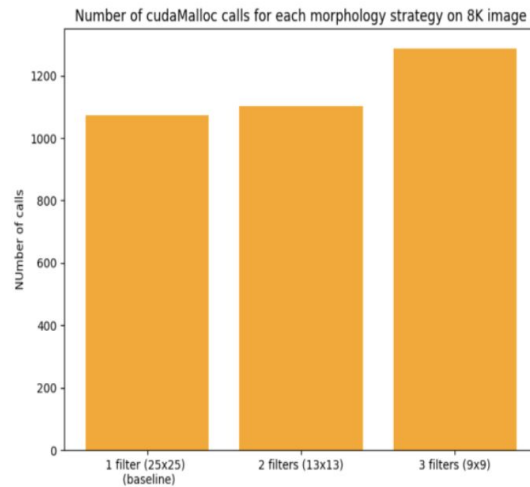
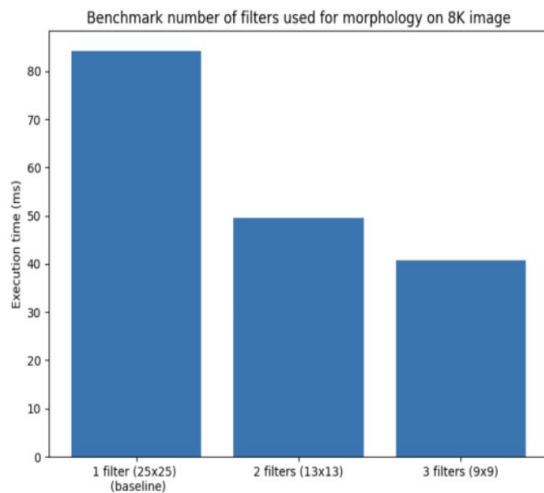


## Numbers of filters for morphology

Idea: instead of using a big 25x25 morphological operation, perform two successive 13x13, or even three 9x9.



# Numbers of filters for morphology

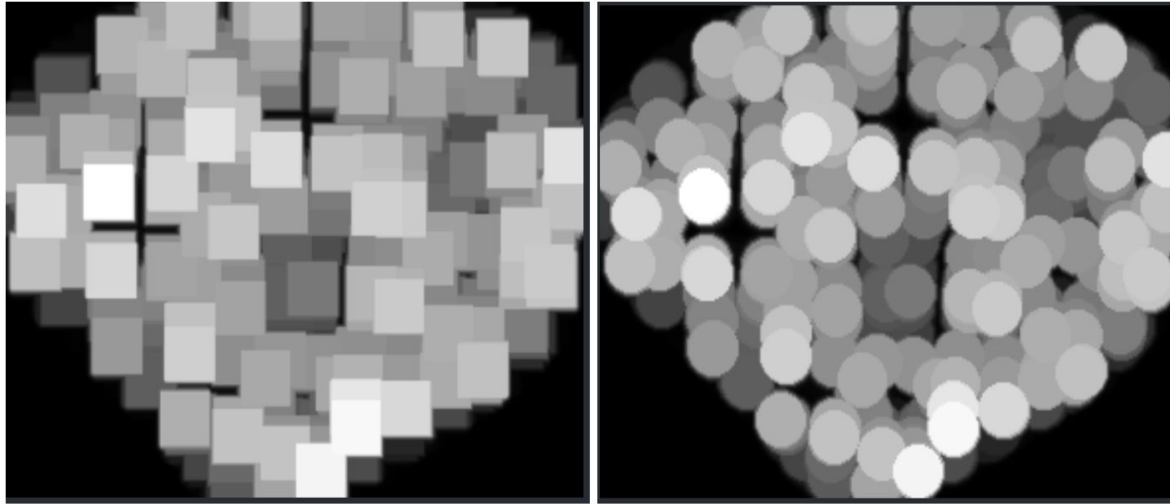




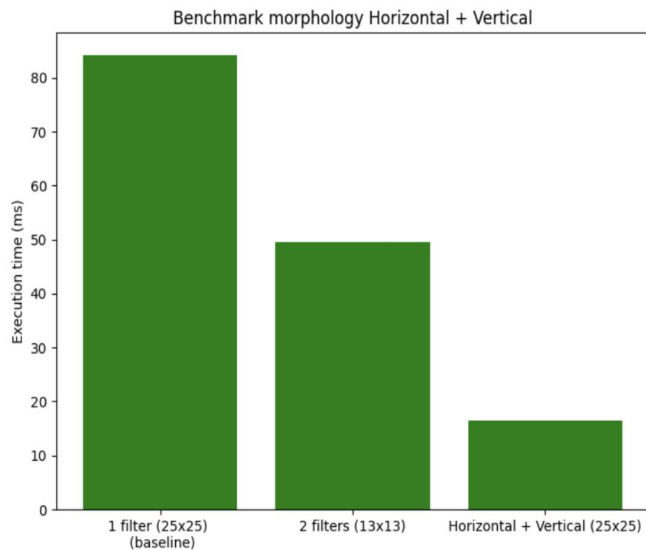
## Horizontal + vertical filters for morphology

Idea: compute the morphology on the two axes separately.

## Horizontal + vertical filters for morphology



# Horizontal + vertical filters for morphology





# Tiled convolutions

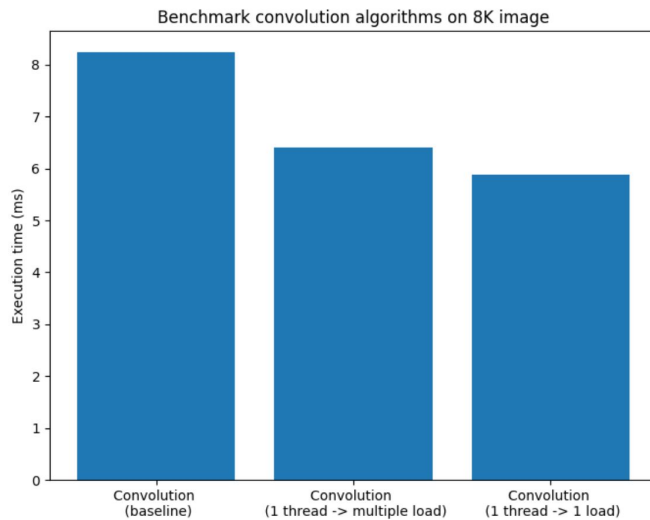
Idea: load tiles in shared memory to allow faster convolution (which was taking 13% of our kernel execution time).

Two solutions:

- One thread  $\leftrightarrow$  one load
- One thread  $\leftrightarrow$  multiple loads



# Tiled convolutions



## Early free (fail)

Idea: cudaMalloc calls are taking a lot of time at the beginning, then CUDA seems to reallocate buffers much faster.

What if we freed buffers as early as possible to allow reuse inside a single run?

Cold run ( `nvprof ./bench` ):

```
Total iteration = 1
--> 99.50% 560.56ms      35 16.016ms 2.3050us 560.39ms cudaMalloc
    0.16% 920.29us      13 70.791us 4.6180us 420.46us cudaDeviceSynchronize
    0.08% 436.55us      56 7.7950us 538ns 130.73us cudaStreamSynchronize
    0.05% 285.28us      75 3.8030us 2.5100us 29.919us cudaLaunchKernel
    0.05% 272.36us      41 6.6420us 1.8300us 71.360us cudaFree
    ...
```

Hot run ( `nvprof ./bench --benchmark_min_time=5` ):

```
Total iterations = 3098
    28.73% 2.34151s    54717 42.793us 1.0890us 456.27us cudaDeviceSynchronize
    14.24% 1.16072s   235704 4.9240us 554ns 456.39us cudaStreamSynchronize
    13.60% 1.10827s   315675 3.5100us 2.3460us 350.29us cudaLaunchKernel
--> 12.54% 1.02178s   147315 6.9360us 1.9210us 115.53ms cudaMalloc
    12.20% 994.06ms   172569 5.7600us 1.7430us 13.622ms cudaFree
    ...
```



## Early free (fail)

```
1 // [...]  
2 auto I_xx = img_x * img_x;  
3 {auto _ = std::move(img_x);} // C++ hack  
4 // [...]
```



## Early free (fail)

Before

Cold run ( nvprof ./bench ):

```
Total iteration = 1
--> 99.50% 560.56ms      35 16.016ms 2.3050us 560.39ms cudaMalloc
    0.16% 920.29us      13 70.791us 4.6180us 420.46us cudaDeviceSynchronize
    0.08% 436.55us      56 7.7950us 538ns 130.73us cudaStreamSynchronize
    0.05% 285.28us      75 3.8030us 2.5100us 29.919us cudaLaunchKernel
    0.05% 272.36us      41 6.6420us 1.8300us 71.360us cudaFree
...
```

After

Cold run ( nvprof ./bench ):

```
99.49% 523.37ms      35 14.953ms 2.3620us 523.21ms cudaMalloc
0.17% 918.63us      13 70.664us 4.3960us 421.54us cudaDeviceSynchronize
0.08% 420.61us      54 7.7890us 647ns 131.63us cudaStreamSynchronize
0.06% 320.38us     101 3.1720us 149ns 149.55us cuDeviceGetAttribute
0.05% 259.89us      73 3.5600us 2.5260us 19.670us cudaLaunchKernel
0.04% 201.67us      41 4.9180us 1.8180us 59.146us cudaFree
```

## B) Summary

