

1. Getting started with GPUs at surfSARA

Use `ssh` to connect to the cluster:

```
ssh username@cartesius.surfsara.nl
```

Introduce the password. If correct, you are now logged in.

If you are a MacOS or Linux user, `ssh` is already available to you in the terminal.

If you are a Windows user, you need to use a `ssh` client for Windows. The easiest option is to use `putty`:

<http://www.chiark.greenend.org.uk/~sgtatham/putty/download.html>

For using the GPU nodes, we need a bit of configuration. Type the following (notet that using “Tab” allows you to partially auto-complete the modules):

```
module load 2020
module load CUDA/11.0.2-GCC-9.3.0
```

You need to do that every time after you log in.

Alternatively, add the same lines to your `.bashrc` file, which is found in your home directory (doing so makes this change permanent and automatically loaded at the start of any `ssh` session; in other words, you don't need to type it every time). If you use the `.bashrc` option, log out (use `exit` in you `ssh` session) and log in again. If this step succeeded, you should be able to run the CUDA compiler now, so please try:

```
nvcc --version
```

This should print:

```
nvcc: NVIDIA (R) Cuda compiler driver
Copyright (c) ...
```

For running jobs, we recommend two main options:

1. Use a job script like "myjob.gpu", like the one in `vector-add` -- see further instructions on how to build such a job script here: <https://userinfo.surfsara.nl/systems/cartesius/usage/batch-usage>

Next, use `sbatch` to launch the job (which contains the name of the executable – again, the example in the folder with the `vector-add` exercise should work out of the box):

```
sbatch myjob.gpu
```

The results will be written upon completion into `slurm-<jobid>.out`.

Use `squeue -u <username>` to check whether your job is still in the queue.

2. Run a job in interactive mode, which effectively means you call `srun` to get a "bash" on one of the nodes, from where you use the code directly (try `nvcc` to test) :

```
srun -n 1 -t 10:00 -p gpu_shared --gres=gpu:1 --pty bash -il
```

.... wait ... then you get access to a `gcnXX` node, like this:

```
[sdemo062@gcn17 ~]$
```

... which you have for 10 minutes (because of “-t 10:00”). Once there, you can run the code from the prompt and get the results – see example below:

```
[sdemo062@gcn17 ~]$ ./vector-add
```

2. Folders and applications

In principle, everything you need for the GPU hands-on session is available on Canvas.

You **MUST** copy this code to your own account on the surfsara cluster.

There are several interesting folders for this practical: `vector-*`, `crypto`, and `difficult`.

We start with `vector-add`, which contains most of the code for a simple vector addition.

Exercise 1:

a. Open the `vector-add.cu` file. Identify the kernel, and add the necessary operation to implement vector addition.

Correct running code should result into something like:

```
vector-add (Sequential): 0.000155 seconds.  
vector-add (kernel): 0.000055 seconds.  
vector-add (memory): 0.000411 seconds.  
results OK!
```

b. Run the code for at 5 different sizes of the array: 256, 1024, 65536, 655360, and 1000000.

Report the execution times for each timer. Comment on the performance difference between these values.

c. Run the code for 5 different grid geometries, and report the results. What is the best performing configuration? Any idea why?

d. Change the operation from addition to subtraction, multiplication, and division, respectively. Compare the performance of these operations for the GPU kernel. Comment on the results.

Notes:

- there is **no kernel code** in the current, given version: the program will run correctly, but the results (for the given code) are reported incorrect.
- you can compile the code with `make`.
- to run the code, revisit page 1 in this document.

Exercise 1 (continued):

e. Take a look at the `vector-add_events.cu`. Notice the difference in measuring performance by comparing the code with the code in `vector-add.cu`. Comment on the (potential) difference between the actual time measurements for the two versions, for vector-add operations, using the 3 largest array sizes. (That is: do "events" lead to significantly different performance measurements?)

f. Take a look at the `vector-add-streams.cu` implementation. Explain what the code does differently, where the performance improvement should come from (if at all), and check whether the performance is actually improved for the streams version. Could you imagine (and implement) a different/better way of using streams for this code? If so, did this new version improve performance?

Exercise 2:

We now work on the `vector-transform` folder.

Please write the kernel that implements the vector transformation from the sequential version (as seen in function `vectorTransformSeq`).

Run the code for at 5 different sizes of the array: 256, 1024, 65536, 655360, and 1000000.

What do you observe, performance-wise?

Compare against the performance of `vector-add`. What do you observe?

3. A Cryptography example

There are many cryptography examples that can be accelerated using parallel processing. The simplest of them is Cesar's code.

In this symmetric encryption/decryption algorithm, one needs to set a numerical key (1 number) that will be added to every character in the text to be encoded.

For example:

Input : ABCDE

Key: 1

Encrypted output: BCDEF

In this assignment you are requested to build a parallel encryption/decryption of a given text file.

The starting code for this example can be found in the folder `crypto`.

Exercise 3:

Please implement a correct encryption and decryption and test it on at least 5 different files. Make a correlation between the size of the files and the performance of the application for both the sequential and the GPU versions. Report speed-up. Comment on the performance you observe.

Notes:

The file names are fixed: `original.data` is the file to be encrypted, `sequential.data` is the reference result for the CPU encryption, and `cuda.data` is the result of the GPU encryption. You are recommended to use `recovered.data` for the decryption, which should be identical with `original.data`.

To test whether the files are identical, use the `diff` command:

```
diff file1 file2
```

If no output is produced the files are identical. If there is a list of differences printed on the screen, these are marked by position in the original file(s).

Exercise 4:

A very interesting extension of this encryption algorithm is to use a larger key - i.e., a set of values, applied to consecutive letters in the original text.

For example:

Input: ABCDE

Key : [1,2]

Output: BDDFF

Explain how you parallelize the new version, and implement this encryption/decryption algorithm as an extension to the original version. Your algorithm must work with keys as large as 256 characters. Ideally, the key length is configurable and can be set by the user; however, for simplicity, you can simply hard-code a length and the key content, as long as the algorithm works correctly for different key sizes and values.

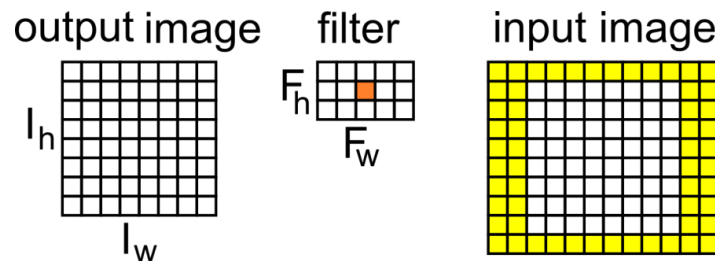
Test this extended version for the same few files as in the previous case, and compare again the results against the sequential version. Report speed-up per file. Comment on the performance you observe.

4. Convolution

2D convolution is an application that finds its use in many domains. The most common is, most likely, image processing, where many filters are based on convolution.

The idea is as follows: given an input image and a filter, the output image is computed such that every pixel is a convolution of the filter and input image. The convolution is a dot product operation.

For example:



A sequential implementation of the convolution is:

```
//for each pixel in the output image
for (y=0; y < image_height; y++) {
  for (x=0; x < image_width; x++) {
    //for each filter weight
    for (i=0; i < filter_height; i++) {
      for (j=0; j < filter_width; j++) {
        output[y][x] += input[y+i][x+j] * filter[i][j];
      }
    }
  }
}
```

The code for convolution is to be found in the folder conv. Note that both the sequential version and a naive kernel implementation are provided.

Exercise 5.

Improve the performance of the convolution kernel using shared memory and any other optimizations you can think of. Please explain how you designed/implemented these optimizations and why you think they should pay off. Conduct a performance analysis study (at least 5 different image sizes), and comment on your findings (e.g.: is there any performance gain? If so, is it as big as you expected it to be?).

5. A more difficult assignment (not graded -- optional)

You are requested to implement an image processing pipeline: the pipeline takes a color image as its input, convert it to grayscale, use the image's histogram to contrast enhance the grayscale image and, eventually, returns a smoothed grayscale version of the input image. For simplicity and accuracy all operations are done in floating point. The program must be benchmarked on the NVIDIA GTX480 GPUs on the DAS-4. A brief description of the four algorithms follows.

Converting a color image to grayscale

Our input images are RGB images; this means that every color is rendered adding together the three components representing Red, Green and Blue. The gray value of a pixel is given by weighting this three values and then summing them together. The formula is:

$$\text{gray} = 0.3 * R + 0.59 * G + 0.11 * B.$$

Histogram Computation The histogram measures how often a value of gray is used in an image. To compute the histogram, simply count the value of every pixel and increment the corresponding counter. There are 256 possible values of gray.

Contrast Enhancement

The computed histogram is used in this phase to determine which are the darkest and lightest gray values actually used in an image - i.e., the lowest (min) and highest (max) gray values that have "scored" in the histogram above a certain threshold. Thus, pixels whose values are lower than min are set to black, pixels whose values are higher than max are set to white, and pixels whose values are inside the interval are scaled.

Smoothing

Smoothing is the process of removing noise from an image. To remove the noise, each point is re-placed by a weighted average of its neighbors. This way, small-scale structures are removed from the image. We are using a triangular smoothing algorithm, i.e. the maximum weight is for the point in the middle and decreases linearly moving from the center. As an example, a 5-point triangular smooth filter in one dimension will use the following weights: 1, 2, 3, 2, 1. In this assignment you will use a two-dimensional 5-point triangular smooth filter.

A sequential version of the application is provided, for your convenience, in the directory "sequential". Please read the code carefully, and try to understand it. A template for the parallel version is also provided, in the "cuda" directory. We recommend that you to start from the template implementation that is provided. You need to parallelize and offload the previously described algorithms to the GPU. The "Kernel" comment in the code indicates the part that must be parallelized.

There are no assumptions about the size of the input images, thus the code must be capable of running with color images of any size. The output must match the sequential version; a compare utility is provided to test for this. The output that you should verify is the final output image, named **smooth.bmp**. You are free to also save the intermediate images, e.g. for debugging, but do not include the time to write this images in the performance measurements.

In the directory "images", 16 different images are provided for testing. You can measure the total execution time of the application, the execution time of the four kernels and the (introduced) memory transfer overheads. This way, you can compute the speedup over the sequential implementation, the achieved GFLOP/s and the utilization.

Performance

Try to optimize (at least) the histogram computation and the smoothing filter (hint: use shared memory).