

1. Getting started with GPUs and the DAS-4

For information about the DAS-4 supercomputer, please go to:
<http://www.cs.vu.nl/das4/>

For information about the special GPU node hardware in the DAS-4 go to the DAS-4 site, then select

“Users → Special Nodes”

For more information on the software for programming GPUs navigate to
“Users → GPUs”

The host name of the VU cluster we are using is: `fs0.das4.cs.vu.nl`

For using the GPU nodes, we need a little bit of configuration. Add the following line to your `.bashrc`:

```
module load cuda55/toolkit prun
```

Now, log out and log in again.

If all is ok, 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) 2005-2013 NVIDIA Corporation
Built on Wed_Jul_17_18:36:13_PDT_2013
Cuda compilation tools, release 5.5, V5.5.0
```

For running jobs, you can use this command:

```
prun -v -np 1 -native '-l gpu=GTX480' <EXECUTABLE >
```

Try, for instance:

```
prun -v -np 1 -native '-l gpu=GTX480'
$CUDA_SDK/bin/x86_64/linux/release/deviceQuery
```

If your job doesn't start immediately, you can check the queue status with
`preserve -long-list`

2. Playing with CUDA

In principle, everything you need for the GPU hands-on session is in:

```
/var/scratch/alvarban/HPC_GPU_2k14
```

The documentation for CUDA is in:

```
/var/scratch/alvarban/HPC_GPU_2k14/cuda-documentation.
```

Especially the CUDA programming guide (`CUDA_C_Programming_Guide.pdf`) is a good starting point and reference for learning and using CUDA. In general, the CUDA documentation is excellent, so use it! You can view it with the “`evince`” program.

You can copy this code to your own account, and play with it:

```
cp -r /var/scratch/alvarban/HPC_GPU_2k14 $HOME
```

In the directory `vector-add` is some example code for a simple vector addition. Let's try to compile and run that code first. You can compile the code with “`make`”.

Now, you should be able to run it with:

```
prun -v -np 1 -native '-l gpu=GTX480' ./vector-add
```

The result should be something like:

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

You can look at the code in `vector-add.cu`, and experiment if you like.

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:

```
/var/scratch/alvarban/HPC_GPU_2k14/crypto
```

Please implement a correct encryption and decryption and test it on several 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.

Note that 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).

Extension

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

For example:

Input: ABCDE

Key : [1,2]

Output: BDDFF

Please implement this encryption/decryption algorithm as an extension to the original version. You can assume the key is already known (fixed), or needs to be transmitted to the kernel as well (user-set).

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.

4. A more difficult assignment

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 these 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 replaced 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.

5. Accelerated Application

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). As an indication, the execution times (in milliseconds) measured for the computation of the largest image in the set (gpu/images/image09.bmp), can be below the following thresholds:

- Grayscale Conversion < 5.5 ms
- Histogram < 68 ms
- Contrast Enhancement < 8.7 ms
- Smoothing < 84 ms
- Total execution < 1023 ms

Compiling and Running Your Application

Please use the provided Makefiles for compiling. Now, you can run your parallel Cuda application with, for example:

```
prun -v -np 1 -native '-l gpu=GTX480' ../bin/cuda
../images/image15.jpg
```

You can look at the results (e.g., smooth.jpg) with the “display” command.
Enjoy!