

CSCE 735 Fall 2020

HW 5: GPU Programming

Due: 11:59pm Wednesday, December 2, 2020

Gaussian Process Regression can be used to predict the values of a function at a point from observations at other points in the domain. As an example, consider an $m \times m$ grid of points on a two-dimensional unit square. Node coordinates are given by

$$(x_i, y_j) = (ih, jh), i = 1, \dots, m, j = 1, \dots, m, \quad (1)$$

where $h = 1/(m+1)$ is the mesh width. Observed data value at the point $r=(x_i, y_j)$ is given by the function

$$f(x_i, y_j) = 1 - \left((x_i - 0.5)^2 + (y_j - 0.5)^2 \right) + d_{ij}, \quad (2)$$

where d_{ij} is a random value between $[-0.05, 0.05]$. The predicted value of the function at $r^*=(x, y)$ is given by

$$f(x, y) = k^T (tI + K)^{-1} f. \quad (3)$$

K is an $n \times n$ matrix with elements $K(r, s) = e^{-\|r-s\|^2}$, where $\|r-s\|^2$ is the distance between grid points r and s (e is the base of the natural logarithm). Note that for an $m \times m$ grid, $n = m^2$. Further, k is an $n \times 1$ vector with elements $k(r) = e^{-\|r-r^*\|^2}$, f is the vector of observed data values given by equation (2), and t is a noise parameter that is set to 0.01.

1. (70 points) In this assignment, you have to develop GPU code to compute $f(x, y)$ for a given point (x, y) . The code should use a **single** processor of the device but should be parallelized to exploit all the cores within the processor. The code should initialize the grid points and observed data values using equations (1) and (2) on the host and move these values to the GPU device. Next, the matrix $A = (tI + K)$ should be computed on the device. This should be followed by LU factorization of A on the device. These factors should be used to compute the solution of the system $Az=f$ using the L and U factors obtained in the previous step. Finally, the predicted value should be computed as $f(x, y) = k^T z$. You must develop your own code to compute the LU factors and to solve the triangular systems. LU factorization can be replaced by Cholesky factorization, which is a more efficient algorithm for symmetric positive definite matrices.
2. (20 points) Describe your strategy to parallelize the algorithm for a **single** multiprocessor of the GPU. Discuss any design choices you made to improve the parallel performance of the code.
3. (10 points) Compute the flop rate you achieve in the factorization routine and in the solver routine. Compare this value with the peak flop rate achievable on the processor, and estimate the speedup obtained over one core and the corresponding efficiency/utilization of the cores on the device. You may choose appropriate values for the grid size to study the features of your implementation.

Submission: You need to upload the following to Canvas:

1. Submit the code you developed.
2. Submit a single PDF or MSWord document that includes the following.
 - Responses to Problem 1, 2, and 3. Response to 1 should consist of a brief description of how to compile and execute the code on the parallel computer

Helpful Information:

1. Source file(s) are available on the shared Google Drive for the class.
2. Load the Intel software stack prior to compiling and executing the code. Use:
`module load intel/2017A CUDA`
3. The run time of a code should be measured when it is executed in dedicated mode. Create a batch file as described on hprc.tamu.edu and use the following command to submit it to the batch system:

```
bsub < batch_file
```