

Programming Assignment 2

CSC410 - Parallel Computing

Aaron G. Alphonsus

November 14, 2018

1 Matrix-Vector Multiplication

A matrix-vector multiplication takes an input Matrix **A**, and a vector **B** and produces an output vector **C**. Each element of the output vector **C** is the dot product of one row of the input Matrix **A** and vector **B**, that is, $C[i] = \sum_j A[i][j] \cdot B[j]$.

This program contains the host program and CUDA kernel to multiply a square matrix with a vector.

2 Performance Analysis

In order to do the performance analysis, we used NVIDIA's profiling tool `nvprof`. This gave us the time taken by the kernel function as well as different CUDA API calls.

The Karp-Flatt metric provides a means to use empirical data to learn something about parallel overhead as well as the amount of inherently sequential computation in a parallel program. The basic idea is to collect data from several runs of the program with increasing numbers of processors.

Definition: Given a parallel computation showing a speedup $\psi(n, p)$ on p processors, where $p > 1$, the experimentally determined serial fraction f_e is defined to be:

$$f_e = \frac{1/\psi(n, p) - 1/p}{1 - 1/p}$$

The speedup $\psi(n, p)$ is calculated as usual - the ratio between sequential execution time and parallel execution time.

$$\psi(n, p) = \frac{T(n, 1)}{T(n, p)}$$

With these formulas, we can benchmark our program. We assumed number of threads to be equal to p . We used a block size of 256 and increased number of blocks by 1 repeatedly. This is why the increase in p is by 256 each time:

| p | 256 | 512 | 768 | 1024 | 1280 | 1536 | 1792 | 2048 |
|--------------|--------|--------|--------|--------|--------|--------|--------|--------|
| $\psi(n, p)$ | 92.95 | 175.93 | 169.81 | 224.40 | 203.94 | 239.84 | 224.87 | 256.82 |
| f_e | 0.0069 | 0.0037 | 0.0046 | 0.0035 | 0.0041 | 0.0035 | 0.0039 | 0.0034 |

What we can see here is that the serial fraction is flattening out indicating that there is limited opportunity for more parallelism. We also notice that some grid size / block size configurations are better than others in achieving speedup even though the total number of threads may be less. We hypothesize that this is due to the design of the hardware.

3 Algorithms and Libraries

Our program follows the typical algorithm for a CUDA C program with a few differences:

1. Declare and allocate host and device memory (we allocated unified memory for the vectors).
2. Initialize vectors on the host. For a quick and easy check of our calculations, we initialized each row of **A** from 1 to **n** and the vector **B** from 1 to **n**. This means that each position in **C** is the sum of squares from 1 to **n** is given by the formula:

$$\frac{n(n+1)(2n+1)}{6}$$

3. Transfer vectors from the host to the device. (This is a usual step i CUDA programs but since we declared unified memory, we did not have do this).
4. Execute the matrix-vector multiplication kernel. We went with coarse-grained parallelism with each thread multiplying a row of matrix **A** with the vector **B**. The dot product summation was done serially.
5. Transfer results from device to host (Not required since we used unified memory).
6. Free the unified memory.

Libraries used:

- The CUDA platform and API
- `<stdio.h>`

4 Functions and Program Structure

The program has 2 functions:

- `main`
- `matvecMul`

4.1 `main`

Arguments: none

Returns: 0 indicating normal termination.

Description:

- Declares vectors and allocates unified memory to them.
- Initializes the vector **B** and each row of the matrix **A** with numbers from 1 to **n**.
- Calls the CUDA kernel after defining grid and block sizes to execute the multiplication in parallel.
- Prints out resultant vector **C**.

4.2 matvecMul

Arguments:

- `double *A`: Square matrix **A** of dimension `n * n`.
- `double *B`: Column vector **B** of length `n`.
- `double *C`: Result of matrix-vector multiplication, row vector **C** of length `n`.
- `int n`: Dimension variable for the matrix and vector (set to be 8192)

Returns: void

Description:

- CUDA kernel code that executes the matrix-vector multiplication on the GPU.
- Coarse-grained parallelism: threads run in parallel for each row of the matrix **A** but execute the dot-product multiplication serially.

5 Compilation and Usage

Compilation: `nvcc -o prog2 prog2.cu` OR `make`

Usage: `./prog2`

The program can be compiled using the command `make`. To get rid of the executable in the folder, run the command `make clean`.

The program multiplies an `n x n` matrix with a vector of length `n`. The vector and each row of the matrix is filled with numbers from 1 to `n`. The resulting vector is the sum of squares from 1 to `n` in each position. `n` is initialized as 8192.

6 Testing and Verification

For verification as we built the program, we used small values of `n` and printed out all three vectors. The answers of the matrix multiplication were verified using Maple. An example can be seen in figure 1.

```
7296837@linux02 program2 >>./prog2
Matrix A
 1.0  2.0  3.0  4.0  5.0
 6.0  7.0  8.0  9.0 10.0
11.0 12.0 13.0 14.0 15.0
16.0 17.0 18.0 19.0 20.0
21.0 22.0 23.0 24.0 25.0
Matrix B
 1.0 2.0 3.0 4.0 5.0
Matrix C
55.0 130.0 205.0 280.0 355.0

7296837@linux02 program2 >>./prog2
Matrix A
 1.0  2.0  3.0  4.0  5.0  6.0  7.0  8.0  9.0 10.0
11.0 12.0 13.0 14.0 15.0 16.0 17.0 18.0 19.0 20.0
21.0 22.0 23.0 24.0 25.0 26.0 27.0 28.0 29.0 30.0
31.0 32.0 33.0 34.0 35.0 36.0 37.0 38.0 39.0 40.0
41.0 42.0 43.0 44.0 45.0 46.0 47.0 48.0 49.0 50.0
51.0 52.0 53.0 54.0 55.0 56.0 57.0 58.0 59.0 60.0
61.0 62.0 63.0 64.0 65.0 66.0 67.0 68.0 69.0 70.0
71.0 72.0 73.0 74.0 75.0 76.0 77.0 78.0 79.0 80.0
81.0 82.0 83.0 84.0 85.0 86.0 87.0 88.0 89.0 90.0
91.0 92.0 93.0 94.0 95.0 96.0 97.0 98.0 99.0 100.0
Matrix B
 1.0 2.0 3.0 4.0 5.0 6.0 7.0 8.0 9.0 10.0
Matrix C
385.0 935.0 1485.0 2035.0 2585.0 3135.0 3685.0 4235.0 4785.0 5335.0
```

Figure 1: Comparison of the static and dynamic schedulers

Once we had a correct solution, we wanted a quick way to test that the program continued to work as we increased the number of elements. To do this we initialized the vector and each row of the matrix with numbers from 1 to n . We did this so that we would have a closed form solution for the resulting vector (sum of squares).

We then calculate the sum of squares, compare it with each element in our resultant matrix, count the mismatches, and print the count and the sum of squares to the screen. The resultant matrix is also printed to the screen along with the grid size and block size for verification.

7 Files Submitted

- prog2.pdf
- Makefile
- prog2.cu
- latex/