

Contents

1	Overview	2
2	CPU Code	2
3	GPU Code	3
3.1	CUDA Code	3
3.2	OpenACC Code	3

1 Overview

The entire code repository can be found at <https://github.com/dgsaf/game-of-life-gpu>. Conway's Game of Life has been accelerated using GPU programming with two models, CUDA and OpenACC, both using C. The code has been derived from the original code which was provided by Cristian Di Pietrantonio, and Maciej Cytowski. It consists of the following items of interest:

- `report/`: The directory containing this `tex` and its resulting `pdf` document.
- `ex1-gol-cuda/`, `ex2-gol-gpu-directives/openacc/`: The CUDA and OpenACC directories have a similar structure, which consists of:
 - `cpu.slurm`: A `slurm` script for submitting CPU jobs on Topaz, for given `n`, `m`, `nsteps`. The CPU code timing output, recorded in [ms], is written to `output/timing-cpu.n-<n>.m-<m>.nsteps-<nsteps>.txt`.
 - `gpu.slurm`: A `slurm` script for submitting GPU CUDA jobs on Topaz, for given `n`, `m`, `nsteps`. The GPU code timing output, recorded in [ms], is written to `output/timing-gpu-cuda.n-<n>.m-<m>.nsteps-<nsteps>.txt` for the CUDA code, and `output/timing-gpu-openacc.n-<n>.m-<m>.nsteps-<nsteps>.txt` for the OpenACC code.
 - `jobs.sh`: A `bash` script which batches a set of jobs, for both the CPU and the GPU codes, on Topaz, for `nsteps` = 100 and `n` = `m` = 1, 2, 4, 8, ..., 16384.
 - `extract.sh`: A `bash` script which, from the jobs batched in `jobs.sh`, `n` = `m` = 1, 2, 4, 8, ..., 16384, extracts the timing output `cpu_elapsed_time`, `cpu_elapsed_time`, `kernel_time`, calculates `speedup`, and writes this performance evaluation to `output/performance.nsteps-<nsteps>.txt`.
 - `output/performance.txt`: A `txt` file which, after the jobs have been submitted and the timing output extracted, contains for each job `n` = `m` = 1, 2, 4, 8, ..., 16384 the performance characteristics `cpu_elapsed_time`, `cpu_elapsed_time`, `speedup`, `kernel_time`.

2 CPU Code

The original code for both the CUDA and OpenACC models has been modified slightly.

- Minor C formatting changes have been made, although only where the original code was modified - unmodified regions of the code remain unadjusted.
- Debugging macros have been utilised to annotate the code for clarity, and can be compiled away to yield performant code.
- The timing methods in `common.c`, `common.h`, have been standardised across the CUDA and OpenACC codes; having originally yielding different return types in each model. The function `float get_elapsed_time(struct timeval start)` now returns the time since `start` was initialised, in [ms], for both codes.
- In `common.c`, `common.h`, the ASCII visualisation has been modified in the following ways: it truncates the grid so that even large grids can be partially visualised in the terminal - allowing for easier verification of grid states across codes, it builds the visualisation output in a

buffer string which is the printed in one call to avoid interference from asynchronous terminal behaviour.

It should be noted that the debugging macros are also ported to the GPU codes, and so, are common to the CPU and GPU codes across both models. The debugging macros are shown in [Listing 1](#), and can be called with a format string, and variable number of arguments similar to how `printf()` is called.

```

3 // debug flags
4 // - 'debug_verbose != 0' will annotate, to stderr, the program as it is
5 //   executed
6 // - 'debug_timing != 0' will annotate, to stderr, the timing variables as
7 //   they are calculated
8 // - 'debug_visual != 0' will annotate, to stderr, the ascii visualisation of
9 //   grid variables as they are initialised and updated
10 const int debug_verbose = 1;
11 const int debug_timing = 1;
12 const int debug_visual = 1;
13
14 // verbose macro
15 #define verbose(format, ...) \
16     if (debug_verbose) { \
17         fprintf(stderr, "[verbose] "format"\n", ##__VA_ARGS__); \
18     }
19
20 // timing macro
21 #define timing(format, ...) \
22     if (debug_timing) { \
23         fprintf(stderr, "[timing] "format"\n", ##__VA_ARGS__); \
24     }
25
26 // visual macro
27 #define visual(current_step, grid, n, m, format, ...) \
28     if (debug_visual) { \
29         fprintf(stderr, "[visual] "format"\n", ##__VA_ARGS__); \
30         visualise_ascii(current_step, grid, n, m); \
31     }

```

Listing 1: Debugging macros from `ex1-gol-cuda/src/game_of_life.c`. Note that these macros are common to both CPU and GPU codes across both models.

3 GPU Code

3.1 CUDA Code

3.2 OpenACC Code

Table 1: Performance characteristics of CUDA code, compiled with -O2, and with debugging statements turned off. All times are presented in units of ms.

n = m	cpu_elapsed_time	gpu_elapsed_time	speedup	kernel_time
1	0.03	240.13	0.0001	1.99
2	0.03	168.12	0.0002	1.14
4	0.05	143.84	0.0004	0.51
8	0.13	149.96	0.0009	0.52
16	0.52	136.58	0.0038	0.56
32	1.17	146.68	0.0080	0.58
64	7.97	132.94	0.0599	0.56
128	16.80	135.33	0.1241	0.60
256	54.00	145.73	0.3706	0.59
512	264.99	131.83	2.0100	0.93
1,024	772.04	135.65	5.6913	1.93
2,048	3,320.08	169.46	19.5919	6.90
4,096	10,886.09	187.74	57.9852	23.89
8,192	79,329.53	377.86	209.9420	94.91
16,384	203,088.16	1,123.27	180.8010	381.04

Table 2: Performance characteristics of OpenACC code, compiled with -O2, and with debugging statements turned off. All times are presented in units of ms.

n = m	cpu_elapsed_time	gpu_elapsed_time	speedup	kernel_time
1	0.02	156.26	0.0001	0.95
2	0.02	162.27	0.0001	0.98
4	0.05	174.83	0.0003	1.17
8	0.06	156.27	0.0004	1.07
16	0.57	177.47	0.0032	1.14
32	0.96	156.14	0.0062	0.99
64	7.29	152.80	0.0477	0.96
128	19.33	175.30	0.1102	1.06
256	59.71	152.44	0.3917	1.16
512	315.81	160.44	1.9684	1.32
1,024	886.82	152.43	5.8178	2.49
2,048	4,088.35	203.10	20.1298	6.75
4,096	12,972.74	198.85	65.2388	22.37
8,192	92,020.30	340.62	270.1540	87.22
16,384	242,173.52	965.25	250.8920	393.68