**Problem 0:**   If necessary, follow the instructions on the
`http://www.ece.lsu.edu/koppel/gpup/proc.html` page for account setup and programming home-
work work flow. For this assignment only edit files `hw07-cuda.cu` and `hw07.cc`. **Find a machine
with GPUs of CC (compute capability) 3.0 or higher,** see the table on
`http://www.ece.lsu.edu/koppel/gpup/sys-status.html`. Compile and run the homework code
unmodified. It should initially show the spring from classroom demo `demo-cuda-04-acc-pat.cc`.

The physics for the spring can be performed using three different sets of code, the one in use
is shown at the beginning of the last line of green text next to Physics. The physics code can be
changed by pressing `a`. Physics code *CPU* uses the CPU and will be very slow. Physics code *CUDA-
M1* and *CUDA-M2* use the GPU running code in file `hw07-cuda.cu`. Code *CUDA-M1* uses routine
`time_step_intersect_1` to detect and resolve intersection of helix segments, and *CUDA-M2* uses
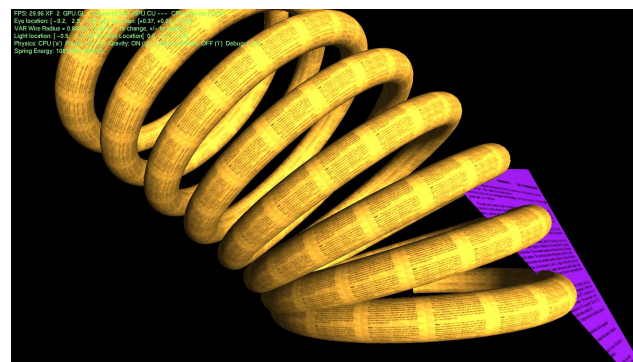routine `time_step_intersect_2`.

The green line starting with "Intersect" shows information about the execution of the chosen
intersection routine. As discussed in class, it is this routine that will dominate performance. On
that line "Bl Sz" is the number of blocks. That can be changed using the `Tab` key to find "Intersect
Block Size" and then pressing `+` and `-` keys. Green text "N Blocks" shows the number of blocks.
The number of blocks is determined by variable "Seg / Block" (`intersect_seg_per_block`). Green
text "Bl/SM" shows the number of blocks per SM in two ways. The first value is the number of
blocks running the intersect routine that *can* fit on an SM. The second number is how many are
assumed running based on the number of blocks and the number of SMs. If the second number
is lower than there are not enough blocks. Green text "Wp/SM" shows the number of warps per
SM. That's based on the block size and the number of blocks per SM. Finally, "Bl Time" shows
the number of cycles it takes to run one block. (The top green line shows timing per frame.)

On the last green line, "Shared Mem" shows how shared memory is being used to buffer the
"b" segments in the intersection routine. "NONE" means that shared memory is not being used,
"SYNC EXPER" is something to be used in Problem 2, "ONE ITER" means that shared memory is
being loaded with "b" segments every `j` iteration (see the intersect routines), and "MULT ITERS"
means that shared memory is loaded with enough data for several `j` iterations—when Problem 3
is solved correctly.

The top green line shows performance. The important items to pay attention to are "GPU.CU"
and "Steps / s". "GPU.CU" shows the time per frame executing CUDA code. Use this to evaluate
the performance of the intersection routines. The last item on the line (you may need to widen
the window), "Steps / s", shows how many physics time steps are performed each second. If
computation is fast enough it will be about 6000, otherwise it will be lower. "FPS" is the frame
rate in frames per second. "XF" shows the frame update interval, an ideal number is 1. A 2
indicates that the frame is updated at half the display's refresh rate. "GPU.GL" is the time the
GPU spends executiong OpenGL code. "CPU GR" is the CPU time for graphics. "CPU PH" is
the CPU time for physics.

Fun Stuff: Pressing `g` will toggle gravity. The free end of the spring can be grabbed by pressing
`b` and it can be moved using the arrow keys.

Initially the arrow keys, PageUp, and
PageDown can be used to move around the
scene. Press (lower-case) `b` and then use the
arrow and page keys to move the free end
of the spring around. (Since the motion is

instantaneous, the spring will act like it was
plucked.) Press `l` to move the light around
and `e` to move the eye (which is what the
arrow keys do when the program starts).

When using the arrow and other keys to move the eye, light, or ball using `Shift` will move by a 5× greater amount and using `Ctrl` will move by one $\frac{1}{5}$ the amount than the motion without either modifier.

Look at the comments in the file `hw07.cc` for documentation on other keys. For documentation see the CUDA C Programming Guide linked to the course `http://www.ece.lsu.edu/koppel/gpup/ref.html` page.

**Problem 1:**   The goal of this assignment is to understand and improve the performance of the intersection routines. The performance of the intersection code is very sensitive to the block size (variable "Intersect Block Size") and to the number of `a` segments per block (variable "Seg / Block", called $m$ in our classroom analysis).

($a$) Indicate the model of GPU you were using (it's shown in the text printed when the program starts). Indicate the number of SMs (shown as MP) for that GPU.

($b$) Adjust these two variables to obtain the best performance on each intersection routine. Do these with shared memory set to NONE. Provide the following information: The settings of those variables. The initial CUDA execution time ("GPU.CU") and the best one obtained.

($c$) Based on our classroom analysis there was less data communication with a larger value of $m$. Are your results consistent with that? If not, explain why.

**Problem 2:**   Set CUDA-M1 to the optimal parameters found in the previous problem. Now, cycle through the different shared memory options. The "ONE ITER" option, which caches enough data for one iteration, should be the slowest. Is it because of the syncthreads?

Modify routine `time_step_intersect_1` so that when `hi.opt_sm_option` is set to value `SMO_sync_experiment` syncthreads is called in a way similar to when `hi.opt_sm_option` is set to `SMO_one_iteration`, but without actually loading or using shared memory. Do so in a way that will determine whether syncthreads is making "ONE ITER" take longer.

Describe what you found by doing the experiment.

**Problem 3:**   Modify `time_step_intersect_1` so that when `hi.opt_sm_option` is set to value `SMO_multiple_iterations` array `pos_cache` is loaded with $B$ items, where $B$ is the block size (value of `blockDim.x`). This should be done inside the `j` loop when data is needed. That is, **don't** load the cache every iteration.

Compare the performance to NONE and ONE ITER.