

Name _____

GPU Microarchitecture
EE 7722
Take-Home Final Examination
Monday, 2 May 2016 to Friday, 6 May 2016

Work on this exam alone. Regular class resources, such as notes, papers, documentation, and code, can be used to find solutions. Do not discuss this exam with classmates or anyone else, except questions or concerns about problems should be directed to Dr. Koppelman.

Problem 1 _____ (35 pts)
Problem 2 _____ (30 pts)
Problem 3 _____ (10 pts)
Problem 4 _____ (10 pts)
Problem 5 _____ (15 pts)

Alias _____

Exam Total _____ (100 pts)

Good Luck!

Problem 1: [35 pts] Appearing below are two versions of a routine, original and CPU-optimized. Both routines compute a running sum but also do something with the first 'low thing' element that they find. Assume that a low thing element will be found by every thread.

Consider their execution on a CC 3.5 GPU. Assume that there are enough threads to hide latency.

Let N denote the value of `data_size` (see the code) and let B and G denote the block and grid sizes (and so the number of threads is BG). Let t_1 be the number of cycles it takes to issue instructions for one thread for one iteration of the loop when low thing is not changed. Let $t_1 + t_s$ be the number of cycles per thread to issue instructions when low thing is set. Assume that $t_s > 100t_1$ due to the cosine and division. Let M denote the number of multiprocessors.

```
__constant__ int data_size, num_threads;
__constant__ float threshold, omega;
__global__ void low_thing_routine(float *dout, float *din) {
    // Original Version
    const int tid = threadIdx.x + blockIdx.x * blockDim.x;
    bool low_thing_set = false;
    float low_thing = 0;
    float sum = 0;
    for ( int i=tid; i<data_size; i += num_threads ) {
        const float d = din[i];
        sum += d;
        if ( !low_thing_set && d < threshold ) {
            low_thing = cos(d*omega) / sum;
            low_thing_set = true;
        }
    }
    dout[tid] = low_thing * sum;
}
```

(a) Using the symbols given above, compute the execution time of the original code for a best-case and worst-case scenario. (The best-case scenario always occurs when $N = BG$, the worst case scenario can't occur unless $N \geq 32BG$.)

- Ignore memory bandwidth limitations.
- Assume that all latency can be hidden.

Execution time for best-case scenario in terms of N , t_1 , etc.

Explain.

Execution time for worst-case scenario in terms of N , t_1 , etc.

Explain.

Problem 1, continued: Appearing below is the so-called CPU-optimized code in which an optimization that makes sense on a CPU has been applied, though the code is still to run on a GPU. This optimization actually fixes one issue with the code on the previous page, but creates two new problems.

```
__constant__ int data_size, num_threads;
__constant__ float threshold, omega;
__global__ void low_thing_routine() {
    // CPU-Optimized Code
    const int tid = threadIdx.x + blockIdx.x * blockDim.x;
    bool low_thing_set = false;
    float low_thing = 0;
    float sum = 0;

    int i = tid;
    for ( ; i < data_size; i += num_threads )
    {
        const float d = din[i];
        sum += d;
        if ( d < threshold ) break;
    }

    low_thing = cos(din[i]*omega) / sum;

    for ( ; i < data_size; i += num_threads ) sum += din[i];

    dout[tid] = low_thing * sum;
}
```

(b) The CPU-optimized code fixes a problem, but creates two new ones. One problem is related to off-chip data bandwidth.

Problem that is fixed (compared to original):

CPU-optimized code accesses more data because:

(c) Again assuming infinite off-chip bandwidth, compute a worst-case execution time. This worst-case time should be based on the fixed problem, and the new one.

Worst-case execution time.

Explain.

Problem 1, continued:

(d) Appearing again below is the Original Version of the code. Re-write it so that it executes efficiently.

Rewrite to avoid inefficiency.

```
__constant__ int data_size, num_threads;
__constant__ float theshold, omega;
__global__ void low_thing_routine(float *dout, float *din) {
    // Original Version

    const int tid = threadIdx.x + blockIdx.x * blockDim.x;

    bool low_thing_set = false;
    float low_thing = 0;
    float sum = 0;

    for ( int i=tid; i<data_size; i += num_threads )
    {

        const float d = din[i];

        sum += d;

        if ( !low_thing_set && d < threshold )
        {
            low_thing = cos(d*omega) / sum;
            low_thing_set = true;
        }

    }

    dout[tid] = low_thing * sum;
}
```

Problem 2: [30 pts] The code fragments below are based on the `vtx-xform-sum` used in class. Remember that in this code each thread computes a value (in routine `cuda_vtx_xform`) and different methods are used to compute a global sum (a sum over all threads).

```
extern "C" __global__ void reduce_atomic_sum_block() {
    const Elt_Type thd_sum = cuda_vtx_xform();
    __shared__ Elt_Type our_sum;
    if ( threadIdx.x == 0 ) our_sum = 0;
    __syncthreads(); // The first syncthreads.
    atomicAdd( &our_sum, thd_sum );
    __syncthreads(); // The second syncthreads.
    if ( threadIdx.x == 0 ) d_app.d_thd_sum[blockIdx.x] = our_sum;
}
```

(a) In the version above an atomic add operating on shared memory is used to compute a block-wide sum. Explain what would happen if each `syncthreads` in the code were removed.

Describe how result would be wrong if first `syncthreads` removed.

Describe how result would be wrong if second `syncthreads` removed.

(b) Appearing below is the assembly code corresponding to the `atomicAdd` from the code above. The add itself is performed by an ordinary add instruction, `FADD`.

```

/*0190*/          BAR.SYNC 0x0;
/*0198*/          PSETP.AND.AND P2, PT, !PT, PT, PT;
/*01a0*/          SSY '(.L_76);
.L_77: /*01c8*/    LDCLK P1, R2, [RZ];
/*01d0*/          @P1 FADD R2, R2, R9;
/*01d8*/          @P1 STSCUL P2, [RZ], R2;
/*01e0*/          @!P2 BRA '(.L_77);
/*01e8*/          NOP.S    (*"TARGET= .L_76 "*);
.L_76: /*01f0*/    BAR.SYNC 0x0;
```

Explain how the other instructions insure that the add is atomic.

Explain why this method of performing an atomic add has poor performance in the CUDA code above.

Problem 2, continued: The code below is based on `reduce_atomic_sum_grid`, but there are three options. Option 0 corresponds to the code used in class. Note that two of the three option lines needs to be commented out for the code to compile.

```
extern "C" __global__ void reduce_atomic_sum_grid() {
    const Evt_Type thd_sum = cuda_vtx_xform();
    const int idx = 0; // Option 0.
    const int idx = blockIdx.x; // Option 1.
    const int idx = threadIdx.x % gridDim.x; // Option 2.

    atomicAdd( &d_app.d_thd_sum[idx], thd_sum );
}
```

Appearing below is the assembly code for the atomic add, with Option 0:

```
/*0170*/      MOV.S R2, c[0x3][0x78]
/*0178*/      MOV R3, c[0x3][0x7c];
/*0188*/      RED.E.ADD.F32.FTZ.RN [R2], R0;
/*0190*/      EXIT;
```

The table below shows the execution time of the three options, and for comparison the sum block code from the previous problem and the tree reduction code (named `reduce_method_2`). The data was collected for a launch of 52 blocks on a GPU with 13 MPs, with an expected occupancy of 4 blocks per MP.

Kernel	Time / μs
<code>reduce_atomic_sum_block</code>	377.152
<code>grid Option 0</code>	337.312
<code>grid Option 1</code>	295.520
<code>grid Option 2</code>	213.888
<code>reduce_method_2</code>	208.576

(c) Provide possible reasons for the performance differences.

Possible reason that Option 1 is faster than Option 0:

Possible reason that Option 2 is faster than Option 1:

Problem 2, continued:

(d) Determine or estimate the amount of data transferred for each of the five kernels above. (The first and last should be simple.) In your answer use B for the block size, G for the number of blocks, and M for the number of multiprocessors.

For the grid kernels, describe how you assume the hardware operates, then provide an estimate based on that assumption.

Amount of data transferred for `reduce_atomic_sum_block`.

Amount of data transferred for `reduce_method_2`.

Describe assumed hardware implementing reduction.

Amount of data transferred for Option 0 (based on assumption).

Amount of data transferred for Option 1 (based on assumption).

Amount of data transferred for Option 2 (based on assumption).

Problem 3: [10 pts] Suppose that most of the area and energy on a Phi chip were used by the vector units and vector registers. Consider an alternative Phi that uses a 1024 b vectors instead of 512 b vectors, but which has half the number of cores, 30 instead of 60. Assume that the alternative and original designs have the same clock frequency and so the same peak FP rate.

(a) Suppose that the cores in the alternative design are the same as in the original design except for the vector unit and registers. Doubling the vector size and halving the core count this way can hurt the performance of some programs.

How might doubling the vector size reduce the performance of some programs?

How might halving the number of cores this way reduce the performance of some other programs? (The performance loss should have a different cause than in the previous answer.)

(b) If the cores in the alternative design are the same except for the vector unit and registers, than the overall cost of the alternative design will be less. Suggest several ways for the cores in the alternative design to be changed to bring cost in line with the original design.

In the alternative design the best thing to do is increase ...

In the alternative design one could also increase ... but it would not be as good as the previous answer.

Problem 4: [10 pts] Answer each question below.

(a) Why does the Phi need branch prediction, but not NVIDIA GPUs? Illustrate with a diagram.

Difference between NVIDIA GPUs and Phi that enables GPUs to omit branch prediction:

(b) In an NVIDIA CC 3.5 device each warp scheduler schedules at most 16 warps. But the block scheduler schedules up to $2^{31} - 1$ blocks.

What would be the approximate cost and performance impact if the warp scheduler had to schedule from a larger number, say 32 warps? (The number of schedulers doesn't change.)

Why can the block scheduler schedule a much larger number of blocks? What makes the block scheduler's job easier?

Problem 5: [15 pts] The code below uses indirection to load from an array named TPairfunc. Further below, is Phi code for the loop body (actually several unrolled iterations).

```
for ( int i=0; i<npairs; i++ )
{
    int othsite = plist[i];
    ans += TPairfunc[othsite] * PsiInv[i];
}

..LN1679: vpaddd    652736(%rdx,%r15){uint16}, %zmm6, %zmm0    #52.11 c1
..L193:   vgatherdpd (%rdi,%zmm0,8), %zmm2{%k1}            #52.11
..LN1689: jkzd     ..L192, %k1    # Prob 50%                #52.11
..LN1690: vgatherdpd (%rdi,%zmm0,8), %zmm2{%k1}            #52.11
..LN1691: jknzd   ..L193, %k1    # Prob 50%                #52.11
..L192:   vfmadd231pd 2368(%rsi,%rcx,8), %zmm2, %zmm8      #52.4 c33
```

(a) The vpaddd instruction is doing two things that an NVIDIA ADD instruction can't do. What are they.

Operations done by vpaddd that would require additional NVIDIA instructions.

(b) Explain how the set of four instructions starting at L193 works.

The four instructions ...

(c) The equivalent NVIDIA assembler code for the instructions starting at L193 would use just use address calculation instructions and one load instruction. Does that mean NVIDIA GPUs perform gathers faster than Phi? Ignore address computation in your answer.

NVIDIA GPUs are faster or about the same speed. (Check one.)

Explain.