

**ETH login ID:**

(Please print in capital letters)

--	--	--	--	--	--	--	--	--	--	--	--

**Full name:**

---

**263-2300: How to Write Fast Numerical Code**

ETH Computer Science, Spring 2016

Midterm Exam

Wednesday, April 20, 2016

**Instructions**

- Make sure that your exam is not missing any sheets, then write your full name and login ID on the front.
- No extra sheets are allowed.
- The exam has a maximum score of 100 points.
- No books, notes, calculators, laptops, cell phones, or other electronic devices are allowed.

Problem 1 (8)	<input type="text"/>
Problem 2 ( $20 = 5 + 3 + 3 + 3 + 6$ )	<input type="text"/>
Problem 3 ( $20 = 4 + 6 + 5 + 5$ )	<input type="text"/>
Problem 4 ( $24 = 12 + 12$ )	<input type="text"/>
Problem 5 (8)	<input type="text"/>
Problem 6 ( $20 = 4 + 8 + 8$ )	<input type="text"/>
<hr/>	
<b>Total (100)</b>	<input type="text"/>

## Problem 1: Peak Performance (8)

Assume a processor that can execute at peak 1 floating point addition and 1 floating point multiplication per cycle. Further assume a function that requires  $a$  floating point additions and  $b$  floating point multiplications. What is the maximal percentage of peak performance that this function can achieve on this processor? Show your work.

**Solution:** The function has a flop count of  $a + b$ . Since the processor can execute 1 addition and 1 multiplication at the same time, the max performance is 2 flops/cycle. The execution time will be bound to the number of additions or multiplications, depending on which one of them is larger. Therefore the performance is:

$$perf = \frac{a + b}{\max(a, b)} \text{ flops/cycle}$$

This function will run at:

$$\frac{perf}{2} \cdot 100\% = \frac{a + b}{2 \cdot \max(a, b)} \cdot 100\%$$

of peak performance.

## Problem 2: Roofline ( $20=5+3+3+3+6$ )

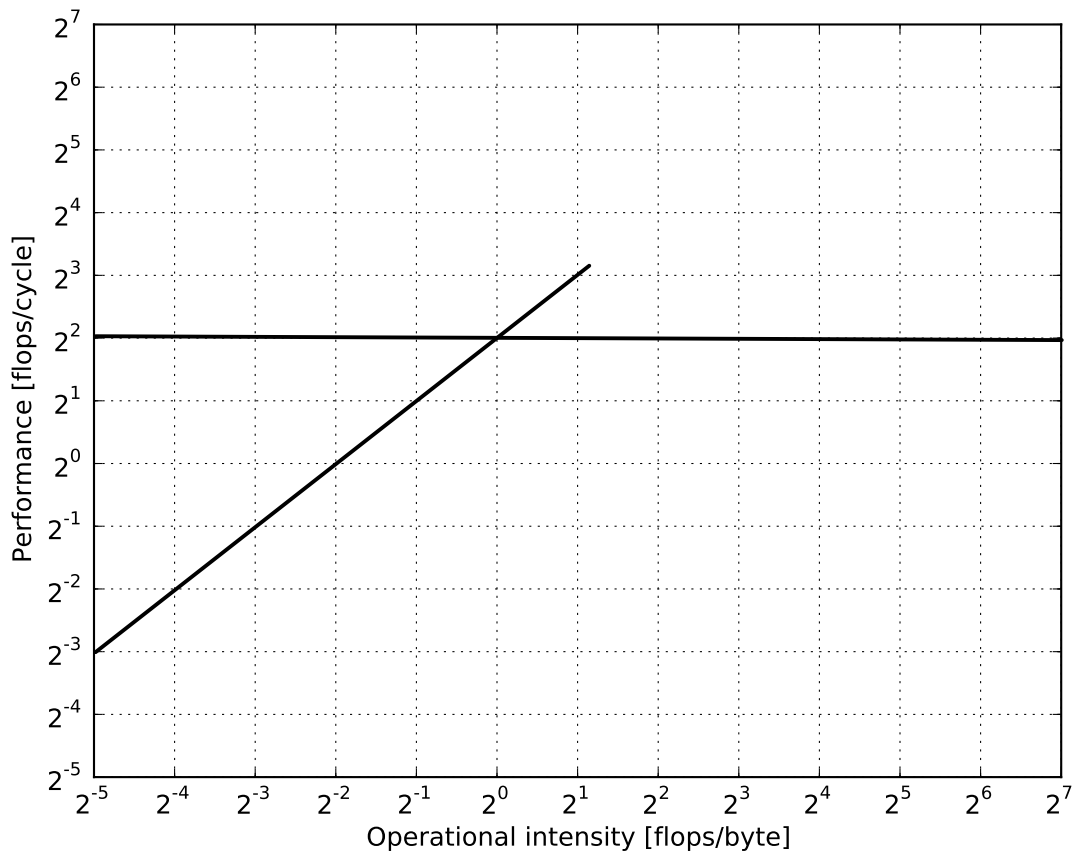
Assume a computer with the following features:

- The CPU can issue, in double precision, 2 (scalar) floating point multiplication and 2 (scalar) floating point additions/subtractions per cycle.
- The total bandwidth between CPU and main memory is 4 bytes/cycle.
- The last level cache is write-back/write-allocate with a size of 2 MB and a block size of 64 bytes.

Further we consider a function  $f$  implementing  $C = AB + C$  in double precision, where each of  $A, B, C$  is an  $n \times n$  triangular matrix, i.e., has  $n(n + 1)/2$  entries, and is stored contiguously in memory. The flop count for this computation is  $n(n + 1)(n + 2)/3$ . The function is executed with cold cache.

**Note:** slight approximations in the calculations are allowed. Explain enough so we see how you got the result.

1. Draw the roofline plot for this system.



2. Determine an upper bound for the operational intensity  $I(n)$  of  $f$ . Consider reads and writes.

**Solution:**

The upper bound for operational intensity arises when there are only compulsory read and write misses, i.e., all the matrices fit in cache. For this case,

$$I(n) \leq \frac{(n+2)}{48}$$

3. For which sizes  $n$  is this bound likely to be the actual intensity?

**Solution:**

If all the matrices fit into the cache the bound calculated above will be the actual intensity. For this,  $12n(n+1) < 2^{21} \implies n < 2^{10}/\sqrt{6}$  (approximately).

4. For which sizes  $n$  is the computation guaranteed memory bound?

**Solution:**

The computation is memory bound for  $I(n) < 1$  flops/byte. For this,  $\frac{(n+2)}{48} < 1 \implies n < 46$ .

5. Now assume the matrices  $A, B, C$  are spread in memory (meaning each is stored non-contiguously with large gaps). How does this change the answers to the three previous questions 2-4?

**Solution:** Since  $A, B, C$  are stored non contiguously, there is guaranteed no spatial locality. More precisely, every load of a double will now fetch a cache block that includes 7 doubles that are not needed, thus transferring 64 bytes instead of 8 bytes.

2. The upper bound now becomes  $I(n) \leq \frac{(n+2)}{384}$

3. All matrices will fit into cache if  $96n(n+1) < 2^{21} \implies n < 2^8/\sqrt{3}$  (approximately).

4.  $\frac{(n+2)}{384} < 1 \implies n < 382$ .

The computation is memory bound for  $\frac{(n+2)}{384} < 1 \implies n < 382$ .

### Problem 3: Performance Analysis (20=4+6+5+5)

The array  $x$  contains  $N$  vectors of length 64 stored consecutively. The function `n2norm` computes the 2-norm for each of these vectors and stores the results in array  $y$ :

```
1 void n2norm(const double * x, double * y, size_t N) {
2     int i;
3     double f1, f2;
4     for (i = 0; i < N; i++) {
5         f1 = 0.;
6         f2 = 0.;
7         for (j = 0; j < 64; j+=2) {
8             f1 += x[i*64+j]*x[i*64+j];
9             f2 += x[i*64+j+1]*x[i*64+j+1];
10        }
11        y[i] = sqrt(f1+f2);
12    }
13 }
```

We make the following assumptions:

- The CPU can issue in each cycle one floating point add and one floating point mult. Both have a latency of one cycle.
- In addition the CPU can issue every 32 cycles a square root. The latency is 32 cycles.
- The system has two levels of cache. Both are write-back/write-allocate.
- L1 cache: size 64 kB, read bandwidth 4 double/cycle.
- L2 cache: size 1 MB, read bandwidth 0.5 double/cycle.
- RAM: read bandwidth 0.25 double/cycle.
- The scalar variables  $i$ ,  $j$ ,  $f1$ ,  $f2$ , and  $N$  are stored in registers.
- A double is 8 bytes.

**Note:** You are allowed to approximate and drop negligible terms.

1. Determine the flop count of this function.

**Solution:**

$$W(N) = N(32 \cdot 4 + 2) = 130N$$

2. Determine the maximal values of  $N$  for which the working set fits into

(a) L1 cache

**Solution:**

$$(64N + N)8 \text{ B} \leq 64k \text{ B}$$
$$N \leq \frac{8k}{65} \approx 128 \text{ doubles}$$

(b) L2 cache

**Solution:**

$$N \leq \frac{1M}{65 \cdot 8} \approx 2k \text{ doubles}$$

3. The performance of `n2norm` is measured as average over many executions. Sketch the expected performance plot for  $N$  up to 100'000.  $N$  is on the x-axis and the y-axis shows the percentage of floating point peak performance (between 0% and 100%) achieved. Provide enough details and also short explanations so we can verify your reasoning.

**Solution:**

- The CPU's peak performance is  $\pi = 1 + 1 + \frac{1}{32} \text{ f/c} \approx 2 \text{ f/c}$ .
- When data fits in L1 the function is compute bound, as we can deduce from the runtimes for computation ( $r_c$ ) and data transfer ( $r_d$ ):

$$r_c = \left(\frac{128}{2} + 1\right)N \text{ cycles} = 65 \text{ cycles}$$
$$r_d = \frac{65N}{4} \text{ cycles} \approx 16 \text{ cycles}$$

Note that the low throughput of `sqrt` is hidden by the innermost loop (64 cycles between two `sqrt` issues). The performance of `n2norm` can be estimated as  $p = \frac{130}{65} \text{ f/c} \approx 2 \text{ f/c}$  (100% of peak).

- When data fits in L2 the function is memory bound:

$$r_d = \frac{65N}{0.5} \text{ cycles} = 130 \text{ cycles.}$$

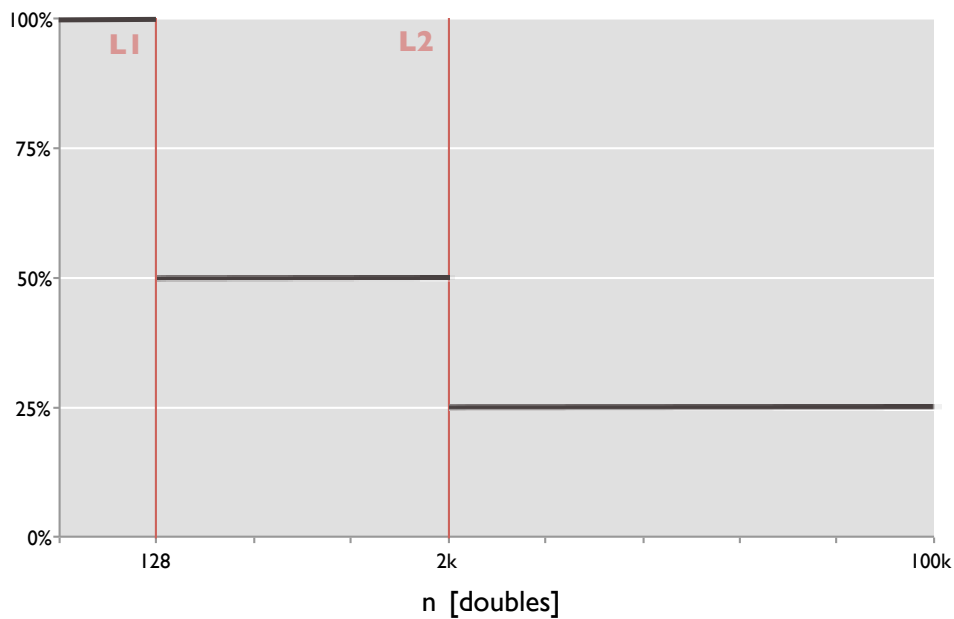
The performance of `n2norm` can be estimated as  $p = \frac{130}{130} \text{ f/c} = 1 \text{ f/c}$  (50% of peak).

- When data does not fit in cache the function is memory bound:

$$r_d = \frac{65N}{0.25} \text{ cycles} = 260 \text{ cycles.}$$

The performance of `n2norm` can be estimated as  $p = \frac{130}{260} \text{ f/c} = 0.5 \text{ f/c}$  (25% of peak).

Percentage of peak performance



4. Now assume that the CPU can issue 2 fused multiply-add (FMA) instructions per cycle, and that each FMA has a latency of 1 cycle. Create a new performance plot analogous to part 3.

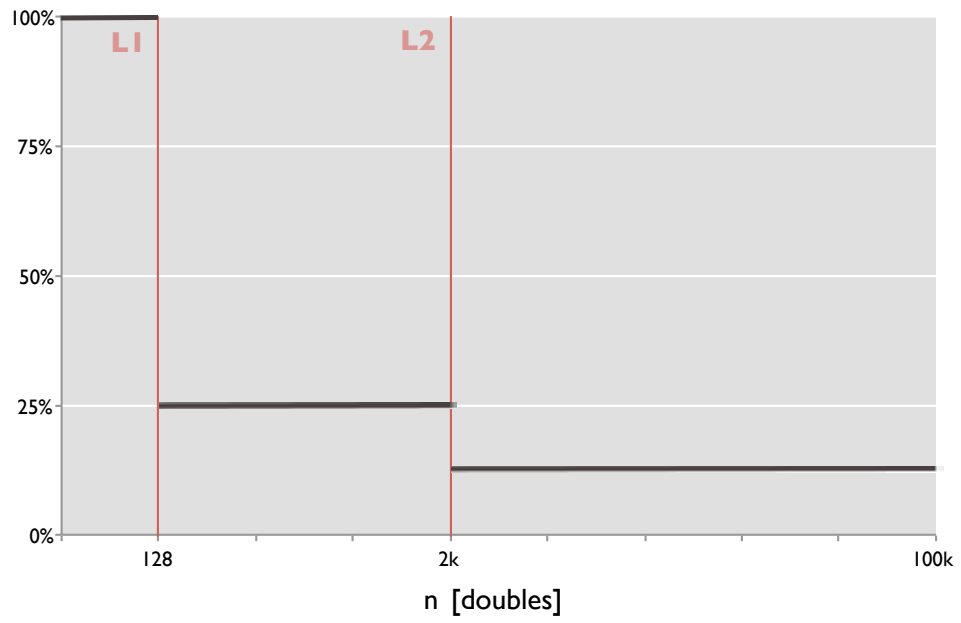
**Solution:** This time the peak performance of the CPU is  $\pi = 4 \text{ f/c}$ . When data fits in L1 the computation runtime can be estimated as

$$r_c = \left(\frac{128}{4} + 1\right)N \text{ cycles} = 33 \text{ cycles.}$$

Again the low throughput of `sqrt` is hidden by the innermost loop (32 cycles between two `sqrt` issues). As in the previous point `n2norm` is compute-bound in this context, and its performance can be estimated as  $p = \frac{130}{33} \text{ f/c} \approx 4 \text{ f/c}$  (100% of peak). In the remaining cases (data fits in L2 and RAM) the function is again memory-bound reaching respectively 25% and 12.5% of peak performance.



Percentage of peak performance



## Problem 4: Cache Mechanics (24=12+12)

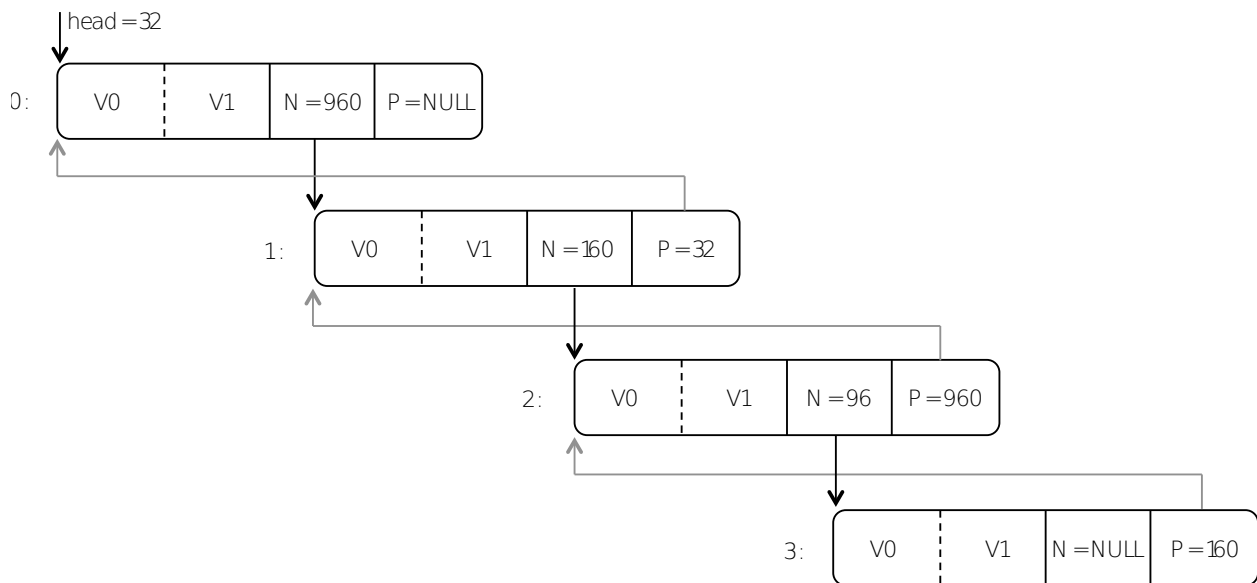
We consider a doubly linked list where nodes have the following structure:

```
1 struct Node {
2   double values[2];    // Labels: V0, V1
3   struct Node* next;   // Label: N
4   struct Node* prev;   // Label: P
5 };
```

Furthermore, we make the following assumptions:

- A double is 8 bytes.
- The system is a 64-bit machine (i.e., pointers hold 8 bytes) with a single cache of size 64 bytes.
- The cache is write-back/write-allocate.
- The cache has (number of sets, associativity, block size) =  $(S, E, 16)$ , i.e., the block size is 16 bytes.
- The cache is initially empty.
- The addresses of all nodes are divisible by 16 (16 byte aligned).
- Memory address 0 maps to the first block of the cache.

Now assume the execution of the code below on the following list:



where every node is labeled with a number from 0 to 3.

```

1 struct Node* p = head; // head == 32
2 while (p->next != NULL) {
3     p->values[0] = p->values[0] + p->values[1];
4     p = p->next; // SHOW CACHE
5 }
6 p->values[0] = p->values[0] + p->values[1]; // SHOW CACHE
7
8 while (p->prev != NULL) {
9     p->values[0] = p->values[0] - p->values[1];
10    p = p->prev; // SHOW CACHE
11 }
12 p->values[0] = p->values[0] - p->values[1]; // SHOW CACHE

```

Assuming that pointers  $p$  and  $head$  are kept in register answer the following questions.

1. Assuming the cache has  $S = 4$  and  $E = 1$ , execute the code above and show the status of the cache whenever requested (right after the execution of lines 4, 6, 10, and 12; so 8 times total). We draw the caches below so you can just fill in. Also provide the hit-miss sequence for the accesses right before the cache status under each cache.

The first snapshot is shown as an example and for the notation you should use. 0 is the node label,  $V_0, V_1$  are the values,  $P, N$  are prev and next. The second snapshot should be filled in the cache to the right of the first one.

**Solution:**

	1.V0 1.V1	1.V0 1.V1	1.V0 1.V1
	1.N 1.P	1.N 1.P	1.N 1.P
0.V0 0.V1	0.V0 0.V1	2.V0 2.V1	3.V0 3.V1
0.N 0.P	0.N 0.P	2.N 2.P	3.N 3.P
<u>MMHHH</u>	<u>MMHHH</u>	<u>MMHHH</u>	<u>MMHH</u>
1.V0 1.V1	1.V0 1.V1	1.V0 1.V1	1.V0 1.V1
1.N 1.P	1.N 1.P	1.N 1.P	1.N 1.P
3.V0 3.V1	2.V0 2.V1	2.V0 2.V1	0.V0 0.V1
3.N 3.P	2.N 2.P	2.N 2.P	0.N 0.P
<u>HHHHH</u>	<u>MMHHH</u>	<u>HHHHH</u>	<u>MMHH</u>

2. Now assume the cache has  $S = 2$  and  $E = 2$  and an LRU replacement policy. Fill the drawn caches and provide the associated hit/miss sequences as before. Again fill the eight caches row by row (meaning, the second snapshot is to the right of the first one). We provide again the first snapshot as example.

**Solution:**

0.V0	0.V1		
0.N	0.P		

MMHHH

0.V0	0.V1	1.V0	1.V1
0.N	0.P	1.N	1.P

MMHHH

2.V0	2.V1	1.V0	1.V1
2.N	2.P	1.N	1.P

MMHHH

2.V0	2.V1	3.V0	3.V1
2.N	2.P	3.N	3.P

MMHH

2.V0	2.V1	3.V0	3.V1
2.N	2.P	3.N	3.P

HHHHH

2.V0	2.V1	3.V0	3.V1
2.N	2.P	3.N	3.P

HHHHH

2.V0	2.V1	1.V0	1.V1
2.N	2.P	1.N	1.P

MMHHH

0.V0	0.V1	1.V0	1.V1
0.N	0.P	1.N	1.P

MMHH

## Problem 5: TLB (8)

We consider the following CPU

- 4 KB page size
- double is 8 bytes
- A fully associative data TLB with 16 entries and LRU replacement.
- Single level, direct mapped cache with 32 sets and 32 bytes line size.

Consider the following code, run with initially cold cache and empty TLB.

```
1 void foo (double *x, int Iterations) { // Iterations is a large number
2   for (int N = 0; N < Iterations; N += 1) {
3     for(int i = 0; i < X; i += S) {
4       x[i] = x[i] * x[i] / 2.0;
5     }
6   }
7 }
```

Replace  $X$  and  $S$  in the code by integers such that this function causes more TLB misses than cache misses. Sketch the cache and explain enough so we see your reasoning. Then derive the TLB misses and cache misses as functions of  $N$  and/or  $Iterations$ . Assume that  $x$  is allocated such that it is big enough for your access pattern.

**Solution:** The page size is 4KB, thus it can fit 512 doubles. To generate more TLB misses than cache misses, we have to make sure that we access it with a stride that is larger than the page size. Thus we infer that  $S \geq 512$ . The TLB is larger than the cache size, thus when traversing virtual pages, we need to avoid conflict misses to reduce the overall cache misses even further. Therefore, we can set:

$$S = 512 + 4$$
$$X = 32 \cdot (512 + 4)$$

This way we make sure that we generate cache misses only on the first iteration of the other loop. On the other hand, we have a TLB miss in every iteration of the innermost loop:

$$misses_{cache} = 32$$
$$misses_{TLB} = 32 * Iterations$$

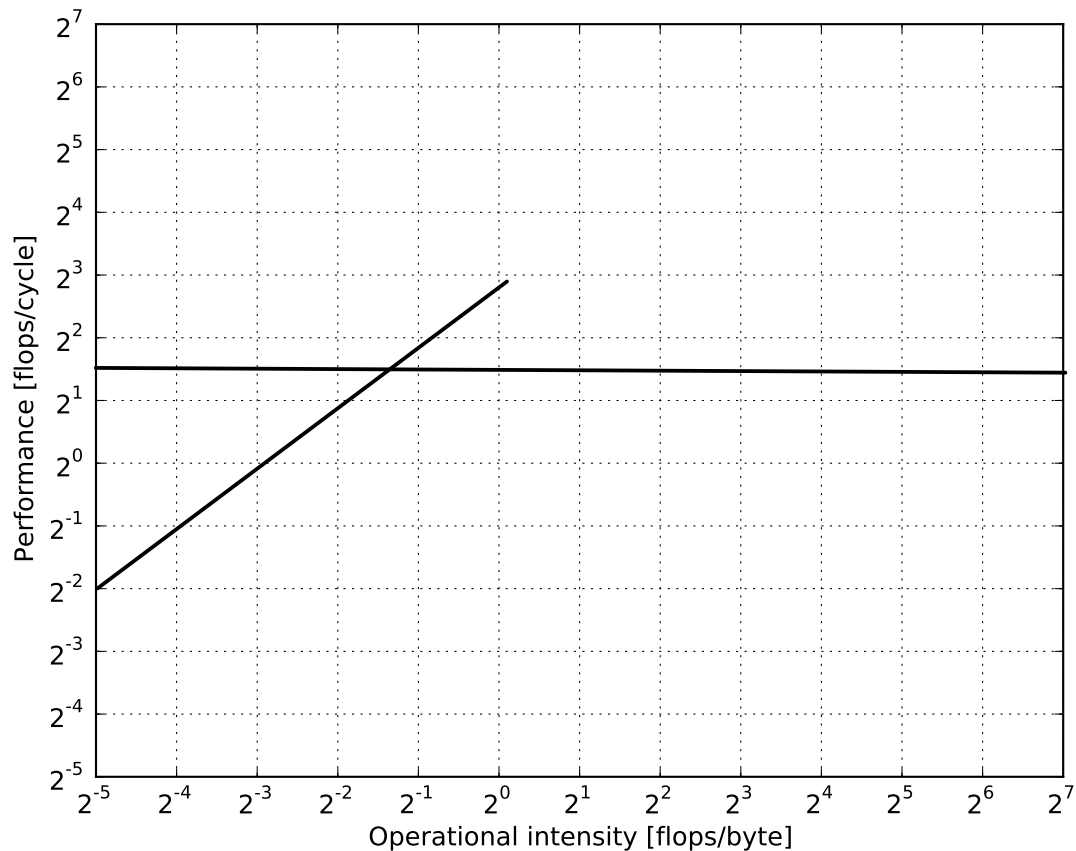
## Problem 6 (20=4+8+8 points)

Assume you are using a system with the following features:

- A CPU that can issue 2 double precision multiplications and 1 double precision addition/subtraction per cycle. Latencies don't matter for this exercise.
- The interconnection between CPU and main memory (size 16 GB) has a maximal bandwidth of 8 bytes/cycle.
- The last level cache is write-allocate/write-back, direct mapped, has size 4 MB and block size of 64 bytes.

Answer the following two questions:

1. Draw the roofline plot for this system:



2. Consider the following code which computes entries in matrix  $g$  using matrix  $f$  ( $f, g$  have size  $n \times n$ ). Assume that  $f$  and  $g$  are allocated sequentially one after the other and first block of  $f$  maps with first cache line:

```
1 void compute(double **f, double **g, int n) {
2     int i, j;
3     double s;
4     for(i = 0; i < n; i++) {
5         for(j = 0; j < n; j++){
6             s = f[0][i]*f[0][i] + f[0][j]*f[0][j];
7             g[i][j] = 0.5*(s - f[i][j]*f[i][j]);
8         }
9     }
10 }
```

Assume a cold cache, that the operators are left associative (expressions are evaluated from left to right), and that a double takes 8 bytes. Now compute for  $n = 256$ ,

- The operational intensity (ignore write-backs).
- An upper bound (as tight as possible) for performance on the specified system.

You are allowed to make minor approximations. Show your work.

**Solution:**

- Both  $f$  and  $g$  fit into cache.  
 $W(n) = 6n^2$ .  
 $Q(n) = 16n^2$ .  
Operational Intensity  $I(n) = \frac{3}{8}$  flops/byte.
- Since  $I(n) = \frac{3}{8}$ , the code is compute bound, the dependencies in computation can be removed by the compiler so that it operates in throughput mode, the instruction mix of the code allows the CPU to perform 3 flops/cycle. Thus, upper bound for performance = 3 flops/cycle.

3. Now consider the following code which computes entries in matrix  $h$  using matrix  $f$  and  $g$  ( $f, g, h$  have size  $n \times n$ ). Assume that  $f, g$  and  $h$  are allocated sequentially one after the other and first block of  $f$  maps with first cache line:

```
1 void compute(double **f, double **g, double **h, int n) {
2     int i, j, k;
3     for(i = 0; i < n; i++){
4         for(j = 0; j < n; j++){
5             for(k = 0; k < n; k++){
6                 h[i][j] = h[i][j] + 0.25*g[i][k] + 0.75*f[k][j];
7             }
8         }
9     }
10 }
```

Assume a cold cache, that the operators are left associative (expressions are evaluated from left to right), and that a double takes 8 bytes. Now compute for  $n = 1024$ ,

- The operational intensity (ignore write-backs).
- An upper bound (as tight as possible) for performance on the specified system.

You are allowed to make minor approximations. Show your work.

**Solution:**

- The size of both matrix  $f$  and  $g$  is twice the size of cache. Each iteration of outer  $i$ -loop loads matrix  $f$  completely which is accessed columnwise. Since matrix  $f$  does not fit in cache this results in cubic cache misses for each access to  $f$  which dominates the memory transfer.  
 $W(n) = 4n^3$ .  
 $Q(n) \approx 64n^3$ .  
 $I(n) \approx \frac{1}{16}$  flops/byte.
- $I(n) < \frac{3}{8}$ , the code is memory bound. Thus, upper bound for performance =  $\frac{1}{2}$  flops/cycle.