

Last name, first name: _____

Student number: _____

263-0007-00L: Advanced Systems Lab

ETH Computer Science, Spring 2022

Midterm Exam

Wednesday, April 27, 2022

Instructions

- Write your full name and student number on the front.
- Make sure that your exam is not missing any sheets.
- 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 ($22 = 2+2+4+4+6+4$)	<input type="text"/>
Problem 2 ($13 = 3+2+4+4$)	<input type="text"/>
Problem 3 ($14 = 6+4+4$)	<input type="text"/>
Problem 4 ($16 = 2+2+6+6$)	<input type="text"/>
Problem 5 ($18 = 2+4+2+4+6$)	<input type="text"/>
Problem 6 ($17 = 6+4+7$)	<input type="text"/>
<hr/>	
Total (100)	<input type="text"/>

Problem 1: Sampler (22 = 2+2+4+4+6+4)

Be brief in your answers, no need to show derivations unless indicated otherwise.

1. Why can a CPU resolve write-after-read (WAR) and write-after-write (WAW) dependencies but not read-after-write (RAW)? How are these dependencies resolved?
2. Can a computation with an $O(1)$ operational intensity benefit from blocking for the caches?
3. Answer the following regarding SIMD intrinsics. Use of pseudo code or descriptive pictures in your answer are both fine.
 - (a) Provide the specification for **one** of the following intrinsics (you can choose which one). We will take the worst answer if two specifications are provided. Use the notation x_i to indicate the i -th element of a vector x .
 - `_mm256_permute_pd(_m256d a, int mask)` or
 - `_mm256_unpacklo_pd(_m256d a, _m256d b)`
 - (b) Explain what `_mm256_set1_ps` does.

4. Consider the computation $\sum_{i=0}^{n-1} \sum_{j=0}^{n-1} A_{ij}$ that sums all elements of a matrix of doubles A of size $n \times n$. Assume that matrix A is sparse with k non-zero elements and full rank. The computation is done with A in CSR (compressed sparse row) format. Indices are represented by (4-byte) integers. Determine a tight upper-bound for the operational intensity of the computation. Show your work.

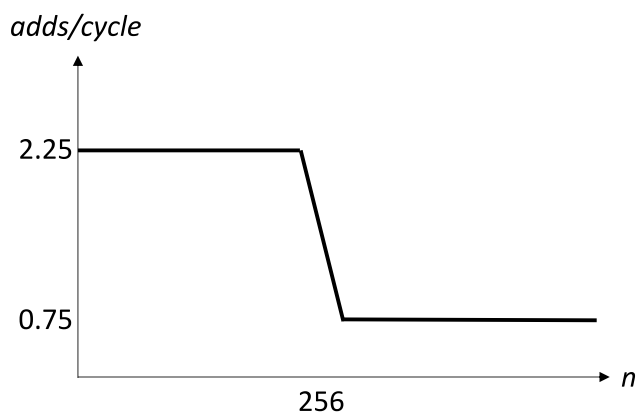
5. Consider a computer with a direct-mapped cache of size 1024 bytes and a block size of 64 bytes. In addition, it has a direct-mapped TLB with 32 entries. The page size is 2MiB. In the following code, assume that the vector x starts at address 0 in memory and that all memory accesses happen in exactly the order that they appear. Determine the smallest value for k that yields a TLB miss for every memory access but only yields compulsory misses in the cache. The cache and TLB are initially empty. Show your work.

```
1 void compute(double *x, unsigned int k) {  
2     double t;  
3     for (int i = 0; i < 256; i += 1) {  
4         t = x[i];  
5         x[i + k] = t + 0.1;  
6     }  
7 }
```

6. Consider the following function. `sizeof(float) = 4`.

```
1 void vecsum(float *a, float *b, float *c, float d*, int n){
2   for (int i = 0; i < n; i += 1) {
3     d[i] = a[i] + b[i] + c[i] + d[i];
4   }
5 }
```

It is run on an Intel-like single-core computer that can perform 3 additions per cycle, without SIMD vector executions. The machine only has one cache. For different input sizes n (starting with very small n), warm-cache measurement yields the following performance plot.



(a) Estimate the size of the cache in bytes. Note that the x-axis shows the input size n which is the length of each vector.

(b) Estimate the read bandwidth β_{cache} in bytes/cycle to the cache.

Problem 2: Bounds (13 = 3+2+4+4)

Consider the following function:

```
1 void compute(float* x, int n, int m){
2     float v1, v2, v3;
3     float c1 = 0.1;
4     float c2 = 0.2;
5     float c3 = 0.3;
6     for(int i=0; i < m-1; i++){
7         x[(i+1)*n] = 1.0;
8         for(int j=0; j < n-2; j++){
9             v1 = x[i*n + j];
10            v2 = x[i*n + j+1];
11            v3 = x[i*n + j+2];
12            x[(i+1)*n+j+1] = (v1+c1)*(v2+c2)*(v3 OP c3); //OP provided in text
13        }
14        x[(i+1)*n + n-1] = x[i*n + n-1];
15    }
16 }
```

Assume that the above code is executed on a computer with the following relevant latency, gap (inverse throughput), and port information:

Instruction	Latency [cycles]	Gap (inverse throughput) [cycles/instruction]	Port
add	3	1	0
mult	3	0.5	0/1
div	5	5	2

The processor does **not** support vector instructions. Further assume that:

1. You can ignore the latency and throughput of loads and stores, i.e., assume they have zero latency and infinite throughput.
2. The compiler does not apply any algebraic transformation: the operations are mapped to assembly instructions as shown.
3. Ignore integer operations.
4. A division counts as one floating-point operation.

Show enough detail with each answer so we understand your reasoning.

1. Determine the maximum theoretical floating-point peak performance in flops/cycle of the computer under consideration.

2. Determine the exact flop count $W(n, m)$ of the `compute` function. Assume that `OP` (in line 12) counts as one floating-point operation.

3. Determine a lower bound (as tight as possible) for the runtime (in cycles) and an associated upper bound for the performance of the `compute` function based on the instruction mix, ignoring dependencies between instructions (i.e., don't consider latencies and assume full throughput). Consider the following two cases:
 - (a) assume that `OP` is a **division** operation.

 - (b) assume that `OP` is a **multiplication** operation.

4. Estimate a lower bound (as tight as possible) for the number of cycles that the computation in line 12 takes to complete. Take latency, throughput and dependency information into account and assume that OP is a **division** operation. Draw the corresponding DAG of the computation performed in line 12.

Problem 3: Operational Intensity (14 = 6+4+4)

Consider the following code implementing a strided matrix vector multiplication ($y = Ax+y$):

```
1 void comp1(double *A, double *x, double *y, int n, int stride){
2   for(int i = 0; i < n; i+= stride)
3     for(int j = 0; j < n; j+= stride)
4       y[i] += A[i*n + j] * x[j]
5 }
```

Assume the following:

- $\text{sizeof}(\text{double}) = 8$.
- A write-back/write-allocate cold cache
- The cache block size is 64 bytes.
- The stride is a power of two.
- n is a multiple of the stride: $n = ms$ for $m \in \mathbb{N}$.
- The flop count is $2 \cdot \left(\frac{n}{s}\right)^2 = 2m^2$

In the derivations you can omit lower order terms (writing \approx instead of $=$). Show your work.

1. Determine a hard upper bound for the operational intensity $I(n, s)$ [flops/byte] in terms of n and the stride denoted as s . Consider only compulsory misses for both reads and writes.

2. You are given a computer which has 2 ports. Each port can execute one addition or one multiplication per cycle (no FMA, and no vector instructions). The memory bandwidth β is 24 bytes per cycle. For which strides is the computation memory bound in the sense of the roofline plot?

3. Consider a version that is strided across only one dimension:

```
1 void comp2(double *A, double *x, double *y, int n, int stride){
2   for(int i = 0; i < n; i+= stride)
3     for(int j = 0; j < n; j++) // no stride here
4       y[i] += A[i*n + j] * x[j]
5 }
```

For the same machine as in Task 2 and assuming that only compulsory misses happen during the computation, can `comp2` ever be memory bound? Justify your answer.

Problem 4: Cache Mechanics (16 = 2+2+6+6)

You are given a write-back/write-allocate cache with 4 sets and LRU replacement policy. Its block size is 12 bytes, and the capacity is 96 bytes. Consider the following code which is the same as in Problem 2. `sizeof(float) = 4`.

```
1 void compute(float* x, int n, int m){
2     float v1, v2, v3, v4;
3     float c1 = 0.1;
4     float c2 = 0.2;
5     float c3 = 0.3;
6     for(int i=0; i < m-1; i++){
7         x[(i+1)*n] = 1.0;
8         for(int j=0; j < n-2; j++){
9             v1 = x[i*n + j];
10            v2 = x[i*n + j+1];
11            v3 = x[i*n + j+2];
12            x[(i+1)*n + j+1] = (v1+c1)*(v2+c2)*(v3+c3);
13        }
14        v4 = x[i*n + n-1];
15        x[(i+1)*n + n-1] = v4;
16    }
17 }
```

Assume that array `x` starts at the memory address 0. Variables `i`, `j`, `c1`, `c2` and `c3` are stored in registers. Memory accesses happen in exactly the order that they appear. Answer the following. Show your work. Hint: It helps to draw the cache.

1. How many floats fit into this cache?
2. What is the associativity of this cache?
3. For each of the following values of m and n do the following two things: i) determine the miss rate; ii) draw the state of the cache at the end of the computation. Show your work.

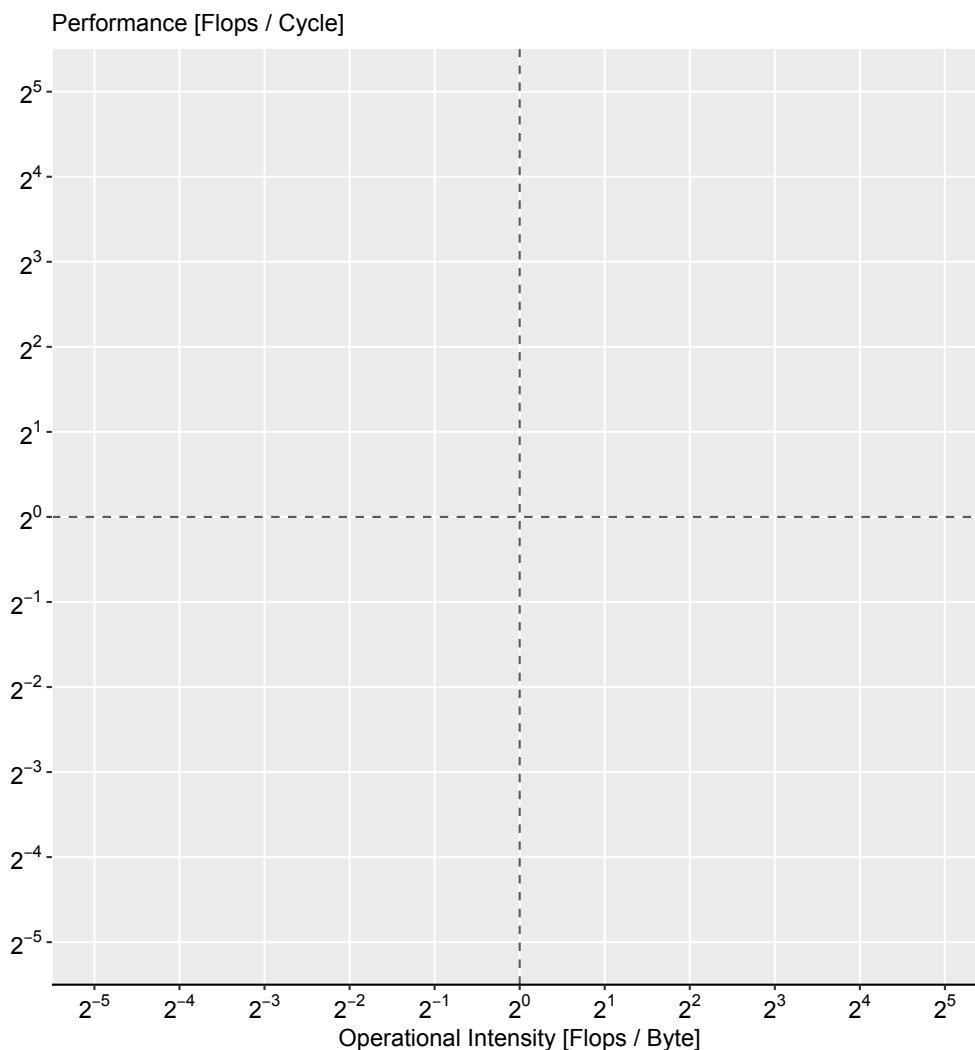
(a) For $m = 2$ and $n = 12$:

(b) For $m = 3$ and $n = 8$:

Problem 5: Roofline (18 = 2+4+2+4+6)

Assume a computer with the following features:

- A CPU with the following ports:
 - Port 1: FMA, MUL, ADD.
 - Port 2: MUL.
 - Port 3: ADD.
- Each of these operations has a throughput of 1 per port and a latency of 4 cycles.
- It does not support any SIMD operations.
- A write-back/write-allocate cache of size 2 MiB with cache block size $B = 64$ bytes. The cache is initially cold.
- The read (memory) bandwidth is 2 doubles per cycle. `sizeof(double) = 8`.



1. Draw the roofline plot for this computer into the above graph. Annotate the lines so we see your reasoning.
2. Consider the following computation where x, y and z are arrays. Assume that x, y , and z are cache-aligned allocated (i.e., the address of an array maps with the first element of a block in the first set of the cache):

```

1 void compute(double* x, double* y, double* z, int n) {
2     for (int i=0; i < n; i++) {
3         y[ i ] = (x[ i ] * y[ i ] + y[i]) + z[i];
4         y[i+1] = (x[i+i] + y[i+1]) * y[i] * z[i+1];
5     }
6 }

```

- (a) Based only on the instruction mix, (i.e., ignoring all type of dependencies), which performance is maximally achievable for this function and why? Draw an associated tighter horizontal roofline into the plot above.
- (b) At what operational intensity $I(n)$ does this new horizontal roofline intersect with the memory roofline?

- (c) What is the performance bound if dependencies are also considered? Assume IEEE 754 arithmetic, i.e., operations cannot be reordered.
3. Assume the cache is fully associative and large enough to fit all the arrays. What is the upper bound for the operational intensity $I(n)$ considering all cache misses? Consider only reads (i.e., ignore write-backs). Based on this $I(n)$, which peak performance is achievable on the specified system taking into account instruction mix and dependencies (i.e., the setting of Task 2c)?

Problem 6: Cache Miss Analysis (17 = 6+4+7)

Consider the following function that uses a i - j - k loop and takes as input matrices A and B of size $n \times n$ and a vector x of size $5n$. A, B, x are not aliased. `sizeof(double) = 8`.

```
1 /* NOTE: Assume that the notation A[i][j] is transformed to A[i*n + j].
2 *       We use the notation A[i][j] for readability only. */
3 void f(double *A, double *B, double *x, int n){
4     double t1, t2;
5     for (int i = 0; i < n-1; i++)
6         for (int j = 0; j < n-1; j++)
7             for (int k = 0; k < n; k++) {
8                 t1 = A[k][ j ] * B[i+1][k];
9                 t2 = A[k][j+1] * B[ i ][k] + x[j + 4*k];
10                A[k][j] = t1 + t2;
11            }
12 }
```

Assume a fully associative write-back/write-allocate cache of size γ bytes with LRU replacement policy, and a cache block size of 64 bytes. Further, assume that n is divisible by 8. Assume an initially cold cache and answer the following. Show your work.

1. Assume that n is much larger than γ (i.e., $n \gg \gamma$) and that γ can fit all data in the innermost loop (i.e., $\gamma > 5 \cdot 64$). Consider cache misses from both reads and writes. Estimate the number of cache misses incurred when accessing each of the arrays as a function of n . In the derivations you can omit lower order terms (writing \approx instead of $=$).

(a) Misses when accessing A :

(b) Misses when accessing B :

(c) Misses when accessing x :

2. Determine the minimum value for γ , as precise as possible, such that the computation only has compulsory misses, i.e., a cache miss only occurs on the first access to a block. For this, assume that LRU replacement is not used and, instead, cache blocks are replaced as effectively as possible to minimize misses.

3. Repeat Tasks 1 and 2 assuming that function f uses a j - k - i loop instead, i.e., the code now looks as follows:

```
1 void f(double *A, double *B, double *x, int n){
2     double t1, t2;
3     for (int j = 0; j < n-1; j++)
4         for (int k = 0; k < n; k++)
5             for (int i = 0; i < n-1; i++) {
6                 t1 = A[k][ j ] * B[i+1][k];
7                 t2 = A[k][j+1] * B[ i ][k] + x[j + 4*k];
8                 A[k][j] = t1 + t2;
9             }
10 }
```

(a) Misses when accessing A :

(b) Misses when accessing B :

(c) Misses when accessing x :

(d) Minimum value of γ such that the computation only has compulsory misses: