

263-0007-00: Advanced Systems Lab

Assignment 4: 120 points

Due Date: April 20th, 17:00

<https://acl.inf.ethz.ch/teaching/fastcode/2023/>

Questions: fastcode@lists.inf.ethz.ch

Academic integrity:

All homeworks in this course are single-student homeworks. The work must be all your own. Do not copy any parts of any of the homeworks from anyone including the web. Do not look at other students' code, papers, or exams. Do not make any parts of your homework available to anyone, and make sure no one can read your files. The university policies on academic integrity will be applied rigorously.

Submission instructions (read carefully):

- (Submission)
Homework is submitted through the [Moodle system](#)
- (Late policy)
You have 3 late days, but can use at most 2 on one homework, meaning submit latest 48 hours after the due time. For example, submitting 1 hour late costs 1 late day. Note that each homework will be available for submission on the system 2 days after the deadline. However, if the accumulated time of the previous homework submissions exceeds 3 days, the homework will not count.
- (Formats)
If you use programs (such as MS-Word or Latex) to create your assignment, convert it to PDF and name it homework.pdf. When submitting more than one file, make sure you create a zip archive that contains all related files, and does not exceed 10 MB. Handwritten parts can be scanned and included.
- (Plots)
For plots/benchmarks, **provide (concise) necessary information for the experimental setup (e.g., compiler and flags) and always briefly discuss the plot and draw conclusions**. Follow (at least to a reasonable extent) the small guide to making plots from the lecture.
- (Neatness)
5 points in a homework are given for neatness.

Exercises

1. Stride Access (15 pts)

Consider the following code executed on a machine with a direct-mapped cache with blocks of size 32 bytes and a total capacity of 1 KiB. Assume that the only memory accesses are to entries of x and occur in the order that they appear (from left to right when in the same line). The cache is initially cold and array x begins at memory address 0.

```
1 double comp(double x[256], int s1, int s2) {
2     double sum = 0.0;
3     for (int i = 0; i < 64; i++) {
4         int j = i + 64;
5         sum += x[s1*i] * x[s2*j];
6     }
7     return sum;
8 }
```

Answer the following. Justify your answers:

- Determine the miss rate when $s_1 = 1$ and $s_2 = 1$.
- Determine the miss rate when $s_1 = 2$ and $s_2 = 2$.
- Determine the miss rate when $s_1 = 2$ and $s_2 = 1$.
- Repeat b) assuming now that the cache is 2-way set associative with a LRU replacement policy. The cache size and block size stay the same.

2. Cache Mechanics (35 pts)

Consider the following code executed on a machine with a direct-mapped write-back/write-allocate cache with blocks of size 8 bytes and a total capacity of 64 bytes. Assume that memory accesses occur in exactly the order that they appear. The variables i, j, m and sum remain in registers and do not cause cache misses. Array x is cache-aligned (first element goes into first cache block) and the first element of y is immediately after the last element of x in memory. Both arrays have size $n = 12$. Assume a cold cache scenario. `sizeof(float) = sizeof(int32_t) = 4 bytes`.

```

1 struct data_t {
2     float a;
3     float b;
4     int32_t u[3];
5 };
6
7 double comp(data_t x[12], data_t y[12]) {
8     float sum = 0;
9     int m = 6;
10    for (int i = 0; i < 3; i++) {
11        for (int j = 0; j <= 9; j+=3) { // j incremented by 3
12            sum += x[(2*i+j)%m].a;
13            sum += y[(4*i+j)%m].b;
14            sum += x[(4*i+j)%m].b;
15        }
16        m = m - 1;
17        // Show state of cache here
18    }
19    return sum;
20 }

```

- (a) Considering the cache misses of the computation, do the following two things for each iteration of the outermost loop: i) determine the miss/hit pattern for x and y (something like x : MMHH... , y : MMMH...); ii) draw the state of the cache at the end of each iteration (i.e. at line 17). See the example below that shows how to draw the cache. Show your work.
- i. Miss/hit pattern and state of the cache in the first iteration ($i = 0$).
 - ii. Miss/hit pattern and state of the cache in the second iteration ($i = 1$).
 - iii. Miss/hit pattern and state of the cache in the third iteration ($i = 2$).
- (b) Repeat the previous task assuming now that the cache is 2-way set associative and uses a LRU replacement policy. The cache size and block size stay the same.

Example. The following example shows how we expect you to draw the cache. The example shows an initially empty cache with $(S, E, B) = (3, 2, 8)$ after $x[0].a$ was accessed. Note that this cache is different from the one specified in the exercise.

State of the cache after accessing $x[0].a$:

Set	0	1
0	$x_0.a, x_0.b$	
1		
2		

3. *Rooflines (40 pt)* Consider a processor with the following hardware parameters (assume 1GB = 10⁹B):

- SIMD vector length of 256 bits.
- The following instruction ports that execute floating point operations:
 - Port 0 (P0): FMA, ADD, MUL
 - Port 1 (P1): ADD, MAX

Each can issue 1 instruction per cycle and each instruction has a latency of 1.

- One write-back/write-allocate cache with blocks of size 64 bytes.
 - Read bandwidth from the main memory is 48 GB/s.
 - Processor frequency is 2 GHz.
- (a) Draw a roofline plot for the machine. Consider only double-precision floating point arithmetic. Consider only reads. Include a roofline for when vector instructions are not used and for when vector instructions are used.
- (b) Consider the following functions. For each, assume that vector instructions are not used, and derive hard upper bounds on its performance and operational intensity (consider only **reads**) based on its **instruction mix** and **compulsory misses**. Ignore the effects of aliasing and assume that no optimizations that change operational intensity are performed (the computation stays as is). FMAs are used to fuse an addition with a multiplication whenever applicable. All arrays are cache-aligned (first element goes into first cache set) and don't overlap in memory. You can further assume that the *max* function is translated into its respective instruction by the compiler and that variables *a*, *b*, *c*, *n* and *i* stay in registers. Assume you write code that attains these bounds, and add the performance to the roofline plot (there should be three dots).

```

1 void comp1(double *x, double *y, double *z, double a, double b, double c, int n) {
2   for (int i = 0; i < n; i++) {
3     z[i] = a * x[i] + y[i] + z[i] * max(x[i] + b, y[i] + c);
4   }
5 }

```

```

1 void comp2(double *x, double *y, double *z, double a, double b, double c, int n) {
2   for (int i = 0; i < n; i++) {
3     z[i] = a + x[i] + y[i] + z[i] + max(x[i] + b, y[i] + c);
4   }
5 }

```

```

1 void comp3(double *x, double *y, double *z, double a, double b, double c, int n) {
2   for (int i = 0; i < n; i++) {
3     z[i] = a * x[i] * y[i] * z[i] * max(x[i] * b, y[i] * c);
4   }
5 }

```

- (c) For each computation, what is the maximum speedup you could achieve by parallelizing it with vector intrinsics?
- (d) Repeat part (b) assuming the following modification in the memory access pattern (strided access). We only show `computation1`, but assume the same modification in `computation2` and `computation3`. Arrays *x* and *y* have an according larger size. Add the new performance of each function to the roofline plot (three additional dots).

```

1 void comp1(double *x, double *y, double *z, double a, double b, double c, int n) {
2   for (int i = 0; i < n; i++) {
3     z[i] = a * x[i] + y[i] + z[i] * max(x[16*i] + b, y[2*i] + c);
4   }
5 }

```

4. Cache Miss Analysis (25 pts)

Consider the following computation that performs a matrix multiplication $C = C + AB^T$ of square matrices A , B and C of size $n \times n$ using a j - i - k loop and .

```
1 void mmm_jik(double *A, double *B, double *C, int n) {
2     for(int j = 0; j < n; j++)
3         for(int i = 0; i < n; i++)
4             for(int k = 0; k < n; k+=2)
5                 C[n*i + j] += A[n*i + k]*B[n*j + k] + A[n*i + k+1]*B[n*j + k+1];
6 }
```

Assume that the code is executed on a machine with a write-back/write-allocate fully-associative cache with blocks of size 64 bytes, a total capacity of γ doubles and with a LRU replacement policy. Assume that n is divisible by 8, cold caches, and that all matrices are cache-aligned. Justify all your answers.

- Assume that $\gamma \ll n$ and determine, as precise as possible, the total number of cache misses that the computation has. For each of the matrices (A , B and C), state also the kind(s) of locality it benefits from to reduce misses.
- Repeat the previous task assuming that we interchange the i -loop and the k -loop, i.e., we have now a j - k - i configuration. Assume that the body of the computation stays the same.
- Using the j - k - i configuration of the previous task, determine the minimum value for γ , as precise as possible, such that the computation only has compulsory misses. For this, assume that LRU replacement is not used and, instead, cache blocks are replaced as effectively as possible to minimize misses.