

1 Data Level Parallelism

<code>__m128i _mm_set1_epi32(int i)</code>	sets the four signed 32-bit integers to i
<code>__m128i _mm_loadu_si128(__m128i *p)</code>	returns 128-bit vector stored at pointer p
<code>__m128i _mm_add_epi32(__m128i a, __m128i b)</code>	returns vector (a0+b0, a1+b1, a2+b2, a3+b3)
<code>__m128i _mm_mullo_epi32(__m128i a, __m128i b)</code>	returns vector (a0*b0, a1*b1, a2*b2, a3*b3)
<code>__m128i _mm_cmpgt_epi32(__m128i a, __m128i b)</code>	returns: vector((a0>b0)?0xffffffff:0, (a1>b1)?0xffffffff:0, (a2>b2)?0xffffffff:0, (a3>b3)?0xffffffff:0)
<code>void _mm_storeu_si128(__m128i *p, __m128i a)</code>	stores 128-bit vector a at pointer p

1.1 The following code uses loop unrolling to improve performance:

```
static void sum_unrolled (int *a) {
    int sum = 0;
    for (int i = 0; i < 16; i += 4) {
        sum += (a + i);
        sum += (a + i + 1);
        sum += (a + i + 2);
        sum += (a + i + 3);
    }
}
```

(a) Implement the following function with a single SIMD instruction:

```
static void sum_unrolled_SIMD (int *a) {
    int sum = 0;
    int result[4];
    __m128i result_v = _mm_set1_epi32(0);
    for (int i = 0; i < 16; i += 4) {
        result_v = _____;
    }
    _mm_storeu_si128((__m128i*)result, result_v);
    sum = result[0] + result[1] + result[2] + result[3];
}

result_v = _mm_add_epi32(result_v, _mm_loadu_si128((__m128i*)a + i));
```

1.2 (a) Implement the following function using SIMD:

```
// Sequential code
static int selective_sum_total (int n, int *a, int c) {
    int sum = 0;
    for (int i = 0; i < n; i += 1) {
        if (a[i] > c) {
            sum += a[i];
        }
    }
    return sum;
}

// SIMD code
static int selective_sum_vectorized (int n, int *a, int c) {
    int result[4];
    __m128i sum_v = _____;
    __m128i cond_v = _____;

    for (int i = 0; i < ____; i += ____) { //Vectorized loop
        __m128i curr_v = _mm_loadu_si128(_____);
        __m128i tmp = _mm_cmpgt_epi32(_____);
        sum_v = _____;
    }
    _mm_storeu_si128(_____);

    for (int i = ____; i < ____; i += 1) { //Tail case
        result[0] += _____;
    }
    return _____;

    __m128i sum_v = _mm_set1_epi32(0);
    __m128i cond_v = _mm_set1_epi32(c);
    for (int i = 0; i < n/4*4; i += 4) { //Vectorized loop
        __m128i curr_v = _mm_loadu_si128((__m128i*)(a+i));
        __m128i tmp = _mm_cmpgt_epi32(curr_v, cond_v);
        tmp = _mm_and_si128(tmp, curr_v);
        sum_v = _mm_add_epi32(sum_v, tmp);
    }
    _mm_storeu_si128((__m128i*)result, sum_v);
    for (int i = n/4*4; i < n; i += 1) { //Tail case
        result[0] += (a[i] > c)? a[i] : 0;
    }
    return result[0] + result[1] + result[2] + result[3];
}
```

- 1.3 (a) You have just begun your virtual internship at Machine Learning Inc. Your boss (since you can't bring her coffee physically) wants you to run data of her preferences through the company's coffee-classifying neural network to decide which flavors she'll like best. You need to get it done before next morning in order to impress her, but you notice that things are running quite slowly - you dig deeper and find that the code's matrix multiplications are not optimized for your machine!

Use SIMD instructions to rewrite the company's matrix-vector multiplication function so that it runs more efficiently on your computer. The matrix has r rows and c columns. You may assume the length of the vector is equal to c .

```
// Company code
static int* matVecMul(int** mat, int* vec, int r, int c) {
    int sum;
    int* result = malloc(r * sizeof(int));
    for (int i = 0; i < r; i++) {
        sum = 0;
        for (int j = 0; j < c; j++) { sum += mat[i][j] * vec[j]; }
        result[i] = sum;
    }
    return result;
}

// SIMD code
static int* efficientMatVecMul(int** mat, int* vec, int r, int c) {
    __m128i sum_v;
    int* result = malloc(_____);
    for (int i = 0; i < r; i++) {
        sum_v = _mm_set1_epi32(0)
        for (int j = 0; j < _____; j++) {
            __m128i mat_v = _mm_loadu_si128(_____);
            __m128i vec_v = _mm_loadu_si128(_____);
            sum_v = _____;
        }
        int* sum_arr = malloc(_____);
        _mm_storeu_si128(sum_arr, sum_v);
        for (int k = _____; k < _____; k++) { sum_arr += _____; }
        result[i] = _____;
    }
    return result;
}
```

```

static int* efficientMatVecMul(int** mat, int* vec, int r, int c) {
    __m128i sum_v;
    int* result = malloc(r * sizeof(int));
    for (int i = 0; i < r; i++) {
        sum_v = _mm_set1_epi32(0)
        for (int j = 0; j < c/4 * 4; j += 4) {
            __m128i mat_v = _mm_loadu_si128((__m128i*) (*mat) + j);
            __m128i vec_v = _mm_loadu_si128((__m128i*) vec + j);
            sum_v = _mm_add_epi32(sum_v, _mm_mullo_epi32(mat_v, vec_v));
        }
        int* sum_arr = malloc(4 * sizeof(int));
        _mm_storeu_si128(sum_arr, sum_v);
        for (int k = c/4 * 4; k < c; k++) {
            sum_arr[0] += mat[i][k] * vec[k];
        }
        result[i] = sum_arr[0] + sum_arr[1] + sum_arr[2] + sum_arr[3];
    }
    return result;
}

```

2 Thread Level Parallelism

Some OpenMP syntax:

<code>#pragma omp parallel{ ... }</code>	Signals the system to spawn threads
<code>#pragma omp for</code>	Split the for loop into equal-sized chunks
<code>int omp_get_num_threads(void);</code>	Returns the number of threads the system has
<code>int omp_get_thread_num(void);</code>	Returns the thread ID of the current thread

2.1 Implement the following function using openMP:

// Sequential code

```
static int selective_product_total (int n, int *a, int c) {
    for (int i = 0; i < n; i += 1) {
        a[i] = (a[i] > c)? 10 : 0;
    }
}
```

- (a) Given $\text{sizeof}(a^*) = \text{sizeof}(\text{int}) * 16$, manually distribute the iterations such that there are four contiguous chunks and no false sharing. Assume the total number of threads is a multiple of 4. `int omp_get_num_threads` \leftarrow returns total number of threads in a team `int omp_get_thread_num` \leftarrow returns current thread ID number

// openMP code

```
static int selective_product_parallelized (int n, int *a, int c) {
    #pragma omp parallel {
        for (int i = 0; i < n; i += 1) {
            if ( _____ ) {
                a[i] = (a[i] > c)? 10 : 0;
            }
        }
    }
}
```

```
if (( i / (n/omp_get_num_threads()) ) == omp_get_thread_num()) {
```

2.2 Implement the following function using openMP:

```
// Sequential code
static int selective_square_total (int n, int *a, int c) {
    for (int i = 0; i < n; i += 1) {
        if (a[i] > c) {
            a[i] *= a[i];
        }
    }
    return product;
}

(a) // openMP for code
static int selective_square_parallelized (int n, int *a, int c) {
    #pragma omp parallel {
        #pragma omp for
        for (int i = ____; i < ____; i += ____) {
            if (_____ > _____) {
                a[i] *= _____;
            }
        }
    }
    return product;
}

    for (int i = 0 ; i < n ; i += 1) {
        if ( a[i] > c ) {
            a[i] *= a[i] ;
        }
    }
```

This is a bit of a trick question: openMP takes care of the parallelism for us by using `#pragma omp parallel / #pragma omp for`. META: Students might get confused on whether they should chunk or alternate threads, because that's usually what the lab covers and/or what lecture covers conceptually.

- 2.3 Assume the code is run with the following function call on a 32-bit machine and access to 8 threads. Also assume that `ptr` is block-aligned:

```
selective_square_parallelized(512, ptr, 61);
```

- (a) Assume the machine uses a 1 KiB direct-mapped cache with 256 B blocks. Is the code correct?

Yes. The intuition here is that over 512 elements, each thread will be given $512 / 8 = 64$ ints to process. $64 \text{ ints} * 4 \text{ B} = 256 \text{ B}$ per thread, which is exactly our block size. There is no problem with false sharing here as each thread will either map to a different index or mismatch in the tag bits. META: This will likely increase conflict misses as the threads switch contexts, because the threads will “fight over” the indexes they map to.

- (b) Assume the machine uses a 1 KiB direct-mapped cache with 1024 B blocks. Is the code correct?

Yes. This question follows the same intuition as the previous part, but now two thread’s worth of data fits into a single block. META: Although this looks like (and will) reduce compulsory misses, the cache management will decrease performance to ensure that false sharing does not happen by preventing threads from writing to the same block at the same time.

3 Data Race!

3.1 Suppose we have `int *A` that points to the head of an array of length `len`. Determine which statement (A)-(E) correctly describes the code execution and provide a one or two sentence justification.

(a) Consider the following code:

```
#pragma omp parallel for
for (int x = 0; x < len; x++){
    *A = x;
    A++;
}
```

Is the code:

- A) Always Incorrect
- B) Sometimes Incorrect
- C) Always Correct, Slower than Serial
- D) Always Correct, Speed relative to Serial depends on Caching Scheme
- E) Always Correct, Faster than Serial

B

Justification: The for loop work is split across threads, but there is a data race to increment the pointer `A`. However, if the threads happen to complete work in disjoint time intervals and in-order, we may get the correct result.

Detailed justification: This implementation is dependent on the threads processing in order, as thread 0, containing the 0 value of `x`, must occur first for the 0th element of `A` to be assigned 0. If any other thread is executed first, `A` cannot correct itself, as `A` is never decremented. Thread ordering with OpenMP is not deterministic, and is certainly not guaranteed to be in order.

Further, there can be data race issues. If multiple threads are executed such that they all execute the first line, `*A = x;` before the second line, `A++;`, they will clobber each other's outputs by overwriting what the other threads wrote in the same position.

On the slim chance (almost impossible, but possible nonetheless) that the threads are scheduled such that each access occurs sequentially, the output will be correct, so it is only "Sometimes Incorrect".


```
(a) #pragma omp parallel
{
    for (int x = 0; x < len; x++){
        *(A+x) = x;
    }
}
```

Is the code:

- A) Always Incorrect
- B) Sometimes Incorrect
- C) Always Correct, Slower than Serial
- D) Always Correct, Speed relative to Serial depends on Caching Scheme
- E) Always Correct, Faster than Serial

C

Justification: Our code computes the correct result but will be slower than the serial equivalent due to duplication of work.

Detailed justification: This is the naive implementation of parallelism. Each thread does the entire job on its own, so there are no concurrency issues.

However, since the serial solution would be the same as our parallel solution, but with only 1 thread, since we can dedicate more resources to a single thread, this solution will be slower.