

1 Pre-Check

This section is designed as a conceptual check for you to determine if you conceptually understand and have any misconceptions about this topic. Please answer true/false to the following questions, and include an explanation:

- 1.1 SIMD is ideal for flow-control heavy tasks (i.e. tasks with many branches/if statements).

False. Data-level parallelism really shines through when we need to repeatedly perform the same operation on a large amount of data. Flow control statements disrupt the continuous flow of computation, which makes programs with them hard to take advantage of SIMD.

- 1.2 Intel's SIMD intrinsic instructions invoke large registers available on the architecture in order to perform one operation on multiple values at once.

True. For example, we can pack four 32-bit integers in a single 128-bit register and perform the same arithmetic operation on all four integers in one go, using an instruction such as `__m128i _mm_add_epi32(__m128i a, __m128i b)`.

- 1.3 The pipelined datapath is an example of parallelism because it performs different stages of instructions in parallel.

True. While a pipelined datapath doesn't execute multiple instructions at the same time, it makes use of each part of the processor at the same time with different instructions, implementing instruction-level parallelism. This can be contrasted with data-level parallelism, which takes advantage of larger registers to do simultaneous memory accesses, and thread-level parallelism, which forks into multiple parallel threads and joins the tasks together.

- 1.4 The most effective way of increasing performance on a modern PC is to increase its clock speed.

False. Modern clock speeds have almost reached their physical limits, and so there's not much room to improve our performance with faster clock speeds. To improve performance, the current best way is to parallelize onto multiple cores (thread-level parallelism).

- 1.5 In thread-level parallelism, the amount of speedup is directly proportional to the increase in number of cores.

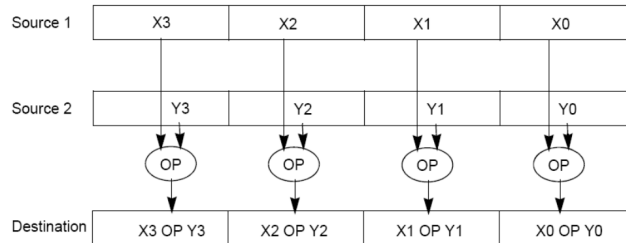
False, usually there is some overhead in parallelizing an operation. Additionally, Amdahl's Law shows that true speedup is affected not only by the number of threads but also by the amount of code that cannot be sped up.

1.6 In thread-level parallelism, threads may run in any order and can start while other threads are partway through their execution.

True. We must ensure that whichever order the threads execute in, the behavior of the program is correct, which includes handling any potential data races.

2 Data-Level Parallelism

The idea central to data level parallelism is vectorized calculation: applying operations to multiple items (which are part of a single vector) at the same time.



Some machines with x86 architectures have special, wider registers, that can hold 128, 256, or even 512 bits. Intel intrinsics (Intel proprietary technology) allow us to use these wider registers to harness the power of DLP in C code.

Below is a small selection of the available Intel intrinsic instructions. All of them perform operations using 128-bit registers. The type `__m128i` is used when these registers hold 4 ints, 8 shorts or 16 chars; `__m128d` is used for 2 double precision floats, and `__m128` is used for 4 single precision floats. Where you see “epiXX”, epi stands for **e**xtended **p**acked **i**nteger, and XX is the number of bits in the integer. “epi32” for example indicates that we are treating the 128-bit register as a pack of 4 32-bit integers.

- `__m128i _mm_set1_epi32(int i)`:
Set the four signed 32-bit integers within the vector to `i`.
- `__m128i _mm_loadu_si128(__m128i *p)`:
Load the 4 successive ints pointed to by `p` into a 128-bit vector.
- `__m128i _mm_mullo_epi32(__m128i a, __m128i b)`:
Return vector $(a_0 \cdot b_0, a_1 \cdot b_1, a_2 \cdot b_2, a_3 \cdot b_3)$.
- `__m128i _mm_add_epi32(__m128i a, __m128i b)`:
Return vector $(a_0 + b_0, a_1 + b_1, a_2 + b_2, a_3 + b_3)$
- `void _mm_storeu_si128(__m128i *p, __m128i a)`:
Store 128-bit vector `a` at pointer `p`.
- `__m128i _mm_and_si128(__m128i a, __m128i b)`:
Perform a bitwise AND of 128 bits in `a` and `b`, and return the result.
- `__m128i _mm_cmpeq_epi32(__m128i a, __m128i b)`:
The `i`th element of the return vector will be set to `0xFFFFFFFF` if the `i`th elements of `a` and `b` are equal, otherwise it'll be set to 0.

- 2.1 SIMD-ize the following function, which returns the product of all of the elements in an array.

```
static int product_naive(int n, int *a) {
    int product = 1;
    for (int i = 0; i < n; i++) {
        product *= a[i];
    }
    return product;
}
```

Things to think about: When iterating through a loop and grabbing elements 4 at a time, how should we update our index for the next iteration? What if our array has a length that isn't a multiple of 4? What can we do to handle this tail case?

```
static int product_vectorized(int n, int *a) {
    int result[4];
    __m128i prod_v = __mm_set1_epi32(1);
    for (int i = 0; i < n/4 * 4; i += 4) { // Vectorized loop
        prod_v = __mm_mullo_epi32(prod_v, __mm_loadu_si128((__m128i *) (a + i)));
    }
    __mm_storeu_si128((__m128i *) result, prod_v);
    for (int i = n/4 * 4; i < n; i++) { // Handle tail case
        result[0] *= a[i];
    }
    return result[0] * result[1] * result[2] * result[3];
}
```

3 Amdahl's Law

In attempting to parallelize a program, the overall performance speedup will always be limited by the fraction of the program that cannot be sped up. This overall speedup can be formulated by Amdahl's Law, which states that

$$\mathbf{Speedup} = \frac{1}{(1 - F) + \frac{F}{S}}$$

Speedup refers to the theoretical speedup of the program compared to its naive implementation. Note that **Speedup** > 1 since we're making our program faster than the original.

F refers to the fraction of the program that can be optimized;

S is the speedup factor for how much that portion of the program can be optimized by, where (**S** > 1)

- 3.1 Derive Amdahl's Law using the ratio: **Speedup** = $t_{\text{naive}}/t_{\text{optimized}}$

First, we can split the overall time a program takes into the time it takes for the part of the program that can be optimized and the rest of it. Letting F represent the fraction that can be sped up, we have:

$$t_{\text{naive}} = F(t_{\text{naive}}) + (1 - F)t_{\text{naive}}$$

Then, we can implement the optimization, known as the speedup factor S into our equation by dividing the optimizable portion to get:

$$t_{\text{optimized}} = \frac{F(t_{\text{naive}})}{S} + (1 - F)t_{\text{naive}}$$

Solving for the ratio **Speedup** = $t_{\text{naive}}/t_{\text{optimized}}$ leads to Amdahl's Law.

- 3.2 Assuming we have infinite threads and resources, what would our overall speedup be for a program with some fraction of our code that can be parallelized F ?

With infinite scaling factor S , our total speedup will approach $\frac{1}{1-F}$. However, in reality there would be some non-zero overhead that is required to properly split up work.

- 3.3 You write code that will search for the phrases "Hello Sean", "Hello Jon", "Hello Dan", "Hello Man", "Bora is the Best!" in text files. With some analysis, you determine you can speed up 40% of the execution by a factor of 2 when parallelizing your code. What is the true speedup?

Using Amdahl's Law with $F=0.4$, $S=2$:

$$\frac{1}{0.6 + \frac{0.4}{2}} = \frac{1}{0.8} = 1.25$$

- 3.4 You run a profiling program on a different program to find out what percent of time within the program each function takes. You get the following results:

Function	% Time
f	30%
g	10%
h	60%

- (a) Assuming that each of these functions can be parallelized by the same speedup factor, which one, if parallelized, would cause the most speedup for the entire program?

h

- (b) What speedup would you get if you parallelized just this function with 8 threads? Assume that work is distributed evenly across threads and there is no overhead for parallelization.

$$1/(0.4 + 0.6/8) \approx 2.1$$

4 Thread-Level Parallelism

OpenMP provides an easy interface for using multithreading within C programs. Some examples of OpenMP directives:

- The `parallel` directive indicates that each thread should run a copy of the code within the block. If a for loop is put within the block, **every** thread will run every iteration of the for loop.

```
#pragma omp parallel
{
    ...
}
```

NOTE: The opening curly brace needs to be on a newline or else there will be a compile-time error!

- The `parallel for` directive will split up iterations of a for loop over various threads. Every thread will run **different** iterations of the for loop. The following two code snippets are equivalent.

```
#pragma omp parallel for          #pragma omp parallel
for (int i = 0; i < n; i++) {    {
    ...                          #pragma omp for
}                                for (int i =0; i < n; i++) { ... }
                                }
```

There are two functions you can call that may be useful to you:

- `int omp_get_thread_num()` will return the number of the thread executing the code
- `int omp_get_num_threads()` will return the number of total hardware threads executing the code

4.1 For each question below, state and justify whether the program is **sometimes incorrect**, **always incorrect**, **slower than serial**, **faster than serial**, or **none of the above**. Assume the default number of threads is greater than 1. Assume no thread will complete before another thread starts executing. Assume `arr` is an `int[]` of length `n`.

(a) // Set element `i` of `arr` to `i`

```
#pragma omp parallel
{
    for (int i = 0; i < n; i++)
        arr[i] = i;
}
```

Slower than serial: There is no `for` directive, so every thread executes this loop in its entirety. `n` threads running `n` loops at the same time will actually execute in the same time as 1 thread running 1 loop. The values should all be correct at the end of the loop since each thread is writing the same values. Furthermore,

the existence of parallel overhead due to the extra number of threads will slow down the execution time.

```
(b) // Set arr to be an array of Fibonacci numbers.
    arr[0] = 0;
    arr[1] = 1;
    #pragma omp parallel for
    for (int i = 2; i < n; i++)
        arr[i] = arr[i-1] + arr[i - 2];
```

Always incorrect (when $n > 4$): Loop has data dependencies: The calculation of all threads but the first one will depend on data from the previous thread. Because we said “assume no thread will complete before another thread starts executing,” this code will always read incorrect values.

```
(c) // Set all elements in arr to 0;
    int i;
    #pragma omp parallel for
    for (i = 0; i < n; i++)
        arr[i] = 0;
```

Faster than serial: The **for** directive automatically makes loop variables (such as the index) private, so this will work properly. The **for** directive splits up the iterations of the loop into continuous chunks for each thread, so there will be no data races.

```
(d) // Set element i of arr to i;
    int i;
    #pragma omp parallel for
    for (i = 0; i < n; i++)
        *arr = i;
        arr++;
```

Sometimes incorrect: Because we are not indexing into the array, there is a data race to increment the array pointer. If multiple threads are executed such that they all execute the first line, `*arr = i;` before the second line, `arr++;`, they will clobber each other’s outputs by overwriting what the other threads wrote in the same position.